```
from google.colab import drive
drive.mount('/content/gdrive')
```

 $\Longrightarrow$  Mounted at /content/gdrive

# Importing the necessary libraries
import time
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

df = pd.read\_csv("/content/gdrive/MyDrive/datathon\_ws/df\_NUEVO (1).csv").copy()

df.head()

<del>_</del>		idListaCobro	idCredito	consecutivoCobro	idBanco_x	montoExigible	montoCobrar	montoCobrado	fechaCobroBanco	idRespue	
	0	155938	738973	41396434	2	622.87	622.87	0.00	NaN		
1	1	155938	739017	41396435	2	1069.11	1069.11	0.00	NaN		
	2	155939	739185	41396436	2	4340.83	4340.83	4340.83	02/01/2025		
	3	155940	732324	41396437	2	2134.21	2134.21	0.00	NaN		
	4	155940	737028	41396438	2	815.76	815.76	0.00	NaN		
	5 ro	ws × 24 columns									
	<b>4</b>									•	

df = pd.read\_csv("/content/gdrive/MyDrive/datathon\_ws/df\_NUEV0 (1).csv").copy()

df.head()

**→** 

	idListaCobro	idCredito	consecutivoCobro	idBanco_x	montoExigible	montoCobrar	montoCobrado	fechaCobroBanco	idRespue
0	155938	738973	41396434	2	622.87	622.87	0.00	NaN	
1	155938	739017	41396435	2	1069.11	1069.11	0.00	NaN	
2	155939	739185	41396436	2	4340.83	4340.83	4340.83	02/01/2025	
3	155940	732324	41396437	2	2134.21	2134.21	0.00	NaN	
4	155940	737028	41396438	2	815.76	815.76	0.00	NaN	
5 rows × 24 columns									
4									

df.info()

<<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2114142 entries, 0 to 2114141
Data columns (total 24 columns):

Data	columns (total 24	columns):
#	Column	Dtype
0	idListaCobro	int64
1	idCredito	int64
2	consecutivoCobro	int64
3	idBanco_x	int64
4	montoExigible	float64
5	montoCobrar	float64
6	montoCobrado	float64
7	fechaCobroBanco	object
8	idRespuestaBanco	float64
9	caso_exitoso	int64
10	fechaCreacionLista	object
11	idBanco_y	int64
12	fechaEnvioCobro	object
13	idEmisora	int64
14	id_estrategia	int64
15	servicio	int64

```
16 tiempo feedback
                          int64
    17 Hora_inicio
                          obiect
    18 Hora fin
                          object
    19 pagare
                          float64
    20 capital
                          float64
    21 fechaAperturaCredito object
    22 CobroExito
                          float64
    23 CobroDevuelta
                          float64
    dtypes: float64(8), int64(10), object(6)
   memory usage: 387.1+ MB
# 1) Imports y carga de datos
import pandas as pd
from pathlib import Path
from xgboost import XGBClassifier
from sklearn.model selection import train test split
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.metrics import roc_auc_score, classification_report
# Carga tu muestra (CSV o XLSX) —> df
# -----
# 2) Feature engineering ligero
 # ⇒ Fechas → datetime
# Fechas → tipo datetime (formato dd/mm/yyyy)
df["fechaEnvioCobro"] = pd.to_datetime(df["fechaEnvioCobro"],
                                    dayfirst=True, errors="coerce")
df["fechaAperturaCredito"] = pd.to_datetime(df["fechaAperturaCredito"],
                                    dayfirst=True, errors="coerce")
# ⇒ Días transcurridos desde origen del crédito
df["dias desde origen"] = (df["fechaEnvioCobro"] - df["fechaAperturaCredito"]).dt.days
# ⇒ Hora del envío de cobro (ej. 0-23)
df["hora_envio"] = pd.to_datetime(df["Hora_ inicio"]).dt.hour
# 3) Selección de columnas (sin fugas de información)
cols drop = [
   "idListaCobro", "montoCobrado", "fechaCobroBanco", "idRespuestaBanco",
   "fechaCreacionLista", "idBanco_x", "idBanco_y", "Hora_ inicio", "Hora_fin",
   "fechaAperturaCredito", "fechaEnvioCobro" # reemplazadas por 'hora_envio' y 'dias_desde_origen'
df = df.drop(columns=cols drop)
# Columnas restantes ▶ variables de entrada
cat_cols = ["idEmisora", "id_estrategia", "servicio"]
num\_cols = [
   "consecutivoCobro", "montoExigible", "montoCobrar",
"tiempo_feedback", "hora_envio", "pagare", "capital",
   "dias_desde_origen", "CobroExito", "CobroDevuelta"
target = "caso exitoso"
# 4) Pre-procesador: One-Hot para categóricas, passthrough numéricas
pre = ColumnTransformer(
   transformers=[
      ("cat", OneHotEncoder(handle unknown="ignore", sparse output=True)
, cat cols)
   remainder="passthrough" # deja las columnas numéricas tal cual
```

```
# 5) Modelo XGBoost + Pipeline
xqb = XGBClassifier(
  n estimators=600,
  max depth=6,
  learning rate=0.08,
  subsample=0.8,
  colsample_bytree=0.8,
  eval metric="logloss",
  random state=42,
  n jobs=-1
)
pipe = Pipeline([
   ("prep", pre),
   ("clf", xgb)
])
# 6) Split train / test y entrenamiento
X_train, X_test, y_train, y_test = train_test_split(
  df[cat_cols + num_cols], df[target],
  test_size=0.20, random_state=42, stratify=df[target]
  <ipython-input-56-9486403f1217>:34: UserWarning: Could not infer format, so each element will be parsed individually, fall
     df["hora_envio"] = pd.to_datetime(df["Hora_ inicio"]).dt.hour
```

## Entrenamiento de modelo

```
# 1) Nuevo modelo con peso de clases
ratio = (y_train == 0).sum() / (y_train == 1).sum()
xgb = XGBClassifier(
    n_estimators=800, max_depth=6, learning_rate=0.06,
    subsample=0.8, colsample bytree=0.8,
    scale_pos_weight=ratio,
    eval_metric='aucpr', random_state=42, n_jobs=-1
pipe.set params(clf=xgb).fit(X train, y train)
# 2) Curva P-R para elegir umbral
from sklearn.metrics import precision_recall_curve, auc
proba = pipe.predict_proba(X_test)[:,1]
prec, rec, thr = precision_recall_curve(y_test, proba)
# Selecciona el threshold donde recall ~0.60 (ejemplo)
target recall = 0.60
idx = next(i for i,r in enumerate(rec) if r < target_recall) - 1</pre>
best_thr = thr[idx]
print("Nuevo threshold:", round(best_thr,3),
      " Precision:", round(prec[idx],2),
      " Recall:", round(rec[idx],2))
# 3) Clasificaciones finales
y pred = (proba >= best thr).astype(int)
print(classification_report(y_test, y_pred))
Nuevo threshold: 0.845
                              Precision: 0.46 Recall: 0.6
                                recall f1-score
                   precision
                                                   support
                0
                        0.98
                                  0.97
                                            0.98
                                                    405177
                1
                        0.46
                                  0.60
                                            0.52
                                                     17652
                                            0.95
        accuracy
                                                    422829
                        0.72
                                  0.78
        macro avg
                                            0.75
                                                    422829
```

# Prueba de modelo para un ID random

```
import numpy as np
import pandas as pd
# 1) Elegir un registro aleatorio
  (si solo quieres segundo cobro, descomenta la línea filtrada)
# ------
# df test case = df[df["consecutivoCobro"] == 2].sample(1, random state=None)
df test case = df.sample(1, random state=None)
# 2) Preparar la fila base SIN 'id_estrategia'
base row = df test case.iloc[0].copy()
if "id_estrategia" not in base_row.index:
   raise KeyError("'id_estrategia' no está en la fila: revisa columnas.")
base row = base row.drop(labels="id estrategia")
# ------
# 3) Lista real de tus 12 estrategias (ajusta a tus IDs)
    Si no estás seguro, puedes extraer las 12 principales del DataFrame:
estrategias_12 = sorted(df["id_estrategia"].unique())[:12]
def predict_for_all_strategies(pipeline, base_row, strategies,
                           cat cols, num cols):
   """Devuelve prob_exito por estrategia para una sola fila."""
   rows = []
   for strat in strategies:
       r = base row.copy()
       r["id estrategia"] = strat
       rows.append(r)
   batch = pd.DataFrame(rows)
   # reindex garantiza orden y crea NaN (numéricas) si faltan columnas
   batch = batch.reindex(columns=cat_cols + num_cols)
   # Probabilidades
   probs = pipeline.predict proba(batch)[:, 1]
   return (
      pd.DataFrame({"id estrategia": strategies, "prob exito": probs})
        .sort_values("prob_exito", ascending=False)
        .reset_index(drop=True)
   )
# 4) Ejecutar y mostrar
resultado = predict_for_all_strategies(
   pipeline=pipe.
   base row=base_row,
   strategies=estrategias 12,
   cat cols=cat cols,
   num_cols=num_cols
print("Crédito ejemplo (idCredito={}):".format(df_test_case["idCredito"].values[0]))
print(resultado)
   Crédito ejemplo (idCredito=664299):
      id estrategia prob exito
    0
                     0.054097
    1
                 2
                     0.052852
    2
                11
                     0.050263
    3
                     0.047007
    4
                 7
                     0.046924
    5
                 6
                     0.045366
                     0.045108
```

#### 5 0.039510

## Prediccion de probabilidades

```
import numpy as np
import pandas as pd
# 1) Elegir un registro aleatorio (o consecutivoCobro == 2 si lo prefieres)
df test case = df.sample(1, random state=None)
# df_test_case = df[df["consecutivoCobro"] == 2].sample(1, random_state=None)
# Extrae el idCredito del registro
credit_id = df_test_case["idCredito"].values[0]
# 2) Preparar la fila base SIN 'id estrategia'
base row = df test case.iloc[0].copy()
if "id estrategia" not in base row.index:
    raise KeyError("'id estrategia' no está en la fila: revisa columnas.")
base row = base row.drop(labels="id estrategia")
# 3) Lista real de tus 12 estrategias
estrategias_12 = sorted(df["id_estrategia"].unique())[:12]
def predict for all strategies(pipeline, base row, strategies,
                               cat_cols, num_cols):
    """Devuelve DataFrame con id_estrategia y prob_exito."""
    for strat in strategies:
        r = base row.copy()
        r["id_estrategia"] = strat
        rows.append(r)
   batch = pd.DataFrame(rows)
   batch = batch.reindex(columns=cat_cols + num_cols)
   probs = pipeline.predict_proba(batch)[:, 1]
    return pd.DataFrame({
        "id estrategia": strategies,
        "prob exito":
                         probs
   })
# 4) Ejecutar y enriquecer con idCredito
resultado = predict_for_all_strategies(
   pipeline=pipe,
   base_row=base_row,
   strategies=estrategias 12,
   cat cols=cat cols,
   num_cols=num_cols
)
# Añade idCredito y reordena columnas
resultado["idCredito"] = credit id
resultado = resultado[["idCredito", "id_estrategia", "prob_exito"]]
# Muestra el DataFrame final
print(resultado)
       idCredito id estrategia prob exito
    0
          745301
                                   0.302873
                              1
          745301
                               2
                                   0.342250
    1
          745301
                                    0.303285
    2
    3
          745301
                                    0.250488
    4
          745301
                                    0.307992
                               6
    5
          745301
                                    0.227003
    6
          745301
                              9
                                    0.277258
          745301
                             11
                                   0.224693
resultado.info()
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 8 entries, 0 to 7
    Data columns (total 2 columns):
                        Non-Null Count Dtype
```

```
0 id_estrategia 8 non-null int64
1 prob_exito 8 non-null float32
dtypes: float32(1), int64(1)
memory usage: 228.0 bytes
```

# Añadir columnas de prediccion de probabilidades a cada registro

```
import numpy as np
import pandas as pd
# 1) Define tus 12 estrategias
estrategias 12 = sorted(df["id estrategia"].unique())[:12]
# 2) Asegúrate de que 'id estrategia' esté en cat cols
if "id estrategia" not in cat cols:
    cat_cols = cat_cols + ["id_estrategia"]
# 3) Para cada estrategia, calcula las probabilidades y añade una columna nueva
for strat in estrategias_12:
    # Crea una copia ligera donde fijas esa estrategia para todas las filas
    tmp = df.copy()
    tmp["id_estrategia"] = strat
    # Reindexa según el orden de columnas que espera el pipeline
    batch = tmp.reindex(columns=cat cols + num cols)
    # Predice la probabilidad de éxito (clase 1)
    probs = pipe.predict_proba(batch)[:, 1]
    # Añade la columna al DataFrame original
    df[f"prob_exito_{strat}"] = probs
# Ahora df tiene, además de sus columnas originales, una columna por cada estrategia:
   prob exito <id estrategia>
print(df.head())
```

df.info()

```
consecutivoCobro
                        int64
    montoExigible
                        float64
    montoCobrar
                        float64
    caso exitoso
                        int64
                        int64
    idEmisora
    id estrategia
                        int64
    servicio
                        int64
    tiempo_feedback
                        int64
    pagare
                        float64
 10 capital
                        float64
 11 CobroExito
                        float64
 12 CobroDevuelta
                        float64
 13 dias_desde_origen int64
 14 hora_envio
                        int32
 15 prob exito 1
                        float32
 16 prob_exito_2
                        float32
 17 prob_exito_4
                        float32
 18 prob exito 5
                        float32
 19 prob_exito_6
                        float32
 20 prob exito 7
                        float32
 21 prob_exito_9
                        float32
22 prob_exito_11
                        float32
dtypes: float32(8), float64(6), int32(1), int64(8)
memory usage: 298.4 MB
```

## Añade dos columnas a df:

```
- best_strategy: id de la estrategia (1,2,4,5,6,7,9,11) con mayor utilidad.
 - best utility: valor de la utilidad correspondiente.
import pandas as pd
def seleccionar mejor estrategia(df):
   Añade dos columnas a df:
     - best strategy: id de la estrategia (1,2,4,5,6,7,9,11) con mayor utilidad.
      - best utility: valor de la utilidad correspondiente.
   M = df['montoCobrar']
   Ce = df['CobroExito']
   Cf = df['CobroDevuelta']
   \lambda = 1.0
   # Lista de estrategias disponibles (coincide con sufijos en prob_exito_*)
   estrategias = [1,2,4,5,6,7,9,11]
   # Calculamos utilidades en forma vectorizada
   utilidades = {}
    for s in estrategias:
        p = df[f'prob_exito_{s}']
        U_s = p * M - \lambda * (p * Ce + (1 - p) * Cf)
       utilidades[s] = U s
   util df = pd.DataFrame(utilidades)
                                                    \# columnas = \{1, 2, 4, ...\}
   df['best_strategy'] = util_df.idxmax(axis=1) # el indice (s) con mayor U_s
   df['best_utility'] = util_df.max(axis=1)
                                                    # valor de U_s máximo
    return df
# df = pd.read csv('tus datos.csv')
df2 = seleccionar_mejor_estrategia(df)
# Ahora df['best_strategy'] y df['best_utility'] están poblados.
```

```
KeyError
                                                Traceback (most recent call last)
    /usr/local/lib/python3.11/dist-packages/pandas/core/indexes/base.py in get_loc(self, key)
       3804
     -> 3805
                         return self. engine.get loc(casted key)
       3806
                     except KeyError as err:
    index.pyx in pandas. libs.index.IndexEngine.get loc()
    index.pyx in pandas._libs.index.IndexEngine.get_loc()
    pandas/_libs/hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item()
    pandas/ libs/hashtable class helper.pxi in pandas. libs.hashtable.PyObjectHashTable.get item()
    KeyError: 'prob exito 1'
    The above exception was the direct cause of the following exception:
    KeyError
                                                Traceback (most recent call last)
                                     🗘 3 frames
    /usr/local/lib/python3.11/dist-packages/pandas/core/indexes/base.py in get_loc(self, key)
        3810
        3811
                             raise InvalidIndexError(key)
     -> 3812
                         raise KeyError(key) from err
       3813
                     except TypeError:
        3814
                         # If we have a listlike key, _check_indexing_error will raise
    KeyError: 'prob_exito_1'
 Next steps:
            Explain error
df2.head()
₹
        idCredito consecutivoCobro montoExigible montoCobrar caso_exitoso idEmisora id_estrategia servicio tiempo_feedbac
     0
                                                                             0
                                                                                         5
           738973
                            41396434
                                              622.87
                                                                                                                  1
                                                                                                                                180
                                                           622.87
     1
           739017
                            41396435
                                             1069.11
                                                          1069.11
                                                                             n
                                                                                         5
                                                                                                        a
                                                                                                                  1
                                                                                                                                180
     2
           739185
                            41396436
                                             4340.83
                                                          4340.83
                                                                             1
                                                                                         5
                                                                                                        9
                                                                                                                                180
     3
                            41396437
                                             2134.21
                                                                                         5
                                                                                                        9
           732324
                                                          2134.21
                                                                             0
                                                                                                                                180
     4
           737028
                            41396438
                                              815.76
                                                           815.76
                                                                             0
                                                                                         5
                                                                                                        9
                                                                                                                                180
    5 rows × 25 columns
import pandas as pd
def seleccionar_mejor_estrategia(df):
    Añade tres columnas a df:
      - best_strategy: id de la estrategia con mayor utilidad.
      - best_utility: utilidad de esa estrategia.
      - real utility: utilidad de la estrategia efectivamente usada (id estrategia).
    M = df['montoCobrar']
    Ce = df['CobroExito']
    Cf = df['CobroDevuelta']
    \lambda = 1.0
    # Estrategias candidatas (coinciden con sufijos en prob_exito_*)
    estrategias = [1,2,4,5,6,7,9,11]
    # 1) Calcular utilidades vectorizadas para cada estrategia
    utilidades = {}
    for s in estrategias:
        p s
             = df[f'prob exito {s}']
             = p_s * M - \lambda * (p_s * Ce + (1 - p_s) * Cf)
        Us
        utilidades[s] = U_s
    util df = pd.DataFrame(utilidades, index=df.index)
    # 2) Mejor estrategia por fila
```

```
df['best_strategy'] = util_df.idxmax(axis=1)
df['best_utility'] = util_df.max(axis=1)

# 3) Utilidad de la estrategia real (id_estrategia)
# Primero, revisamos que todas las reales estén en candidatas
missing = set(df['id_estrategia'].unique()) - set(estrategias)
if missing:
    raise ValueError(f"Estrategias reales no contempladas: {missing}")

# Ahora extraemos de util_df la columna que coincida con id_estrategia
# y la guardamos en real_utility
df['real_utility'] = df.apply(
    lambda row: utilidades[row['id_estrategia']].loc[row.name],
    axis=1
)

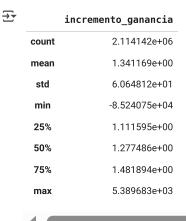
return df

# Ejemplo de uso:
df3 = seleccionar_mejor_estrategia(df)
# Ahora df2 contiene best_strategy, best_utility y real_utility.
```

Se agrega una columna de incremento en ganancia que se calcula al dividir la mejor utilidad entre la utilidad real

```
df3.info()
RangeIndex: 2114142 entries, 0 to 2114141
    Data columns (total 26 columns):
    # Column
                           Dtype
        idCredito
                           int64
         consecutivoCobro
                           int64
         montoExigible
                           float64
                           float64
         montoCobrar
         caso_exitoso
                           int64
        idEmisora
                           int64
         id estrategia
                           int64
         servicio
                           int64
         tiempo_feedback
                           int64
                           float64
         pagare
     10 capital
                           float64
     11 CobroExito
                           float64
     12 CobroDevuelta
                           float64
     13 dias_desde_origen int64
     14 hora envio
                           int32
     15 prob_exito_1
                           float32
     16 prob_exito_2
                           float32
     17 prob exito 4
                           float32
     18 prob_exito_5
                           float32
     19 prob_exito_6
                           float32
     20 prob_exito_7
                           float32
     21 prob_exito_9
                           float32
     22 prob exito 11
                           float32
     23 best_strategy
                           int64
     24 best_utility
                           float64
     25 real utility
                           float64
    dtypes: f\overline{loat32(8)}, float64(8), int32(1), int64(9)
    memory usage: 346.8 MB
import numpy as np
# Evitamos división por cero convirtiendo ceros en NaN
df['utilidad_real_nonzero'] = df['real_utility'].replace(0, np.nan)
# Calculamos el incremento
df['incremento_ganancia'] = df['best_utility'] / df['utilidad_real_nonzero']
# (Opcional) Si quieres que en lugar de NaN aparezca un cero cuando utilidad_real era 0:
df['incremento_ganancia'] = df['incremento_ganancia'].fillna(0)
```

df['incremento\_ganancia'].describe()



df.columns

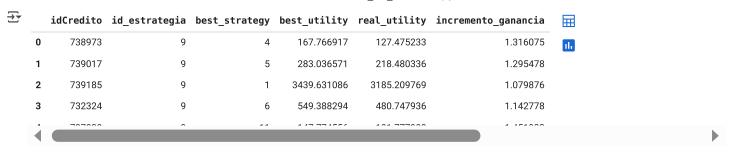
df.head()

<b>→</b>		idCredito	consecutivoCobro	montoExigible	montoCobrar	caso_exitoso	idEmisora	id_estrategia	servicio	tiempo_feedbac
	0	738973	41396434	622.87	622.87	0	5	9	1	180
	1	739017	41396435	1069.11	1069.11	0	5	9	1	180
	2	739185	41396436	4340.83	4340.83	1	5	9	1	180
	3	732324	41396437	2134.21	2134.21	0	5	9	1	180
	4	737028	41396438	815.76	815.76	0	5	9	1	180

5 rows × 28 columns

```
# 1. Selecciona las columnas de interés en un nuevo DataFrame
cols = [
    'idCredito',
    'id_estrategia',
    'best_strategy',
    'best_utility',
    'real_utility',
    'incremento_ganancia'
]
df_export = df[cols].copy()
```

# 2. Comprueba que se ven bien
df export.head()



df\_export.info()

```
RangeIndex: 2114142 entries, 0 to 2114141
   Data columns (total 6 columns):
    # Column
                           Dtype
   _ _ _
       idCredito
                           int64
       id_estrategia
                           int64
        best_strategy
                           int64
       best_utility
real_utility
                           float64
                           float64
    5 incremento ganancia float64
   dtypes: float64(3), int64(3)
   memory usage: 96.8 MB
```

Exportamos un excel con la columna idcredito y las columnas 'best\_strategy',

 'best\_utility', 'real\_utility', 'utilidad\_real\_nonzero', 'incremento\_ganancia' en un nuevo dataframe