

Detección de Fraude por medio de transacciones financieras

Carolina Salas Moreno , Deykel Bernard Salazar, Esteban Ramírez Montero y Kristhel Porras Mata



Introducción

Sistema para detectar fraudes en transacciones financieras usando un dataset sintético que simula comportamientos reales.

El análisis permite aplicar técnicas de aprendizaje automático para identificar fraudes con precisión y rapidez.

Dataset

- Tamaño: 5,000,000 registros
- Atributos: 18 columnas por registro
- Formato: CSV (~800 MB)
- Datos de transacción: ID, fecha, cuentas, monto y tipo.
- Información de comportamiento y metadatos (ubicación, dispositivo, IP).
- Etiquetas de fraude (binarias y tipos), is_fraud como target.
- Disponible en [Kaggle](#)

Desarrollo del proyecto

Tecnologías Utilizadas

- Apache Spark → Procesamiento rápido y distribuido de grandes volúmenes de datos.
- Lenguaje principal: Python
- Bibliotecas clave: Pandas, Scikit-learn, Matplotlib, PySpark
- Delta Lake → Almacenamiento estructurado, confiable y optimizado para consultas eficientes.

Plataforma de Desarrollo

- Databricks (en la nube) – Free edition

Carga de datos

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col
3 spark = SparkSession.builder.getOrCreate()
4
5 # Leer datos desde el volumen en Databricks
6 path = "/Volumes/fraud_project/fraud_detection/fraud_data/"
7 df_spark = spark.read.option("header", True).csv(path)
8 df_spark.printSchema()
9
10 # Convertir a Pandas solo si es necesario para visualización
11 for field in df_spark.schema.fields:
12     if field.dataType.simpleString() == 'string':
13         df_spark = df_spark.withColumn(field.name, col(field.name).cast("string"))
14 df = df_spark.toPandas()
```

```

1 # Exploración básica
2 print("✅ Dimensiones del dataset:", df.shape)
3 print("\n📄 Primeros registros:")
4 display(df.head())
5
6 print("\n📄 Tipos de datos:")
7 print(df.dtypes)
8
9 print("\n🔍 Estadísticas numéricas:")
10 display(df.describe())
11
12 print("\n🔍 Estadísticas categóricas:")
13 #display(df.select_dtypes(include='object').describe())
14 print(df.select_dtypes(include='object').describe())
15
16 # Nulos
17
18 print("\n❓ Valores nulos por columna:")
19 null_counts = df.isnull().sum()
20 null_percentages = (null_counts / len(df)) * 100
21 null_summary = pd.DataFrame({
22     "Nulos": null_counts,
23     "% del total": null_percentages.round(2)
24 })
25 null_summary = null_summary.sort_values(by="Nulos", ascending=False)
26 print(null_summary)
27
28 print("\n❓ Filas duplicadas:", df.duplicated().sum())
29
30 msno.matrix(df, figsize=(15, 4))
31 plt.title("Mapa de valores nulos")
32 plt.show()

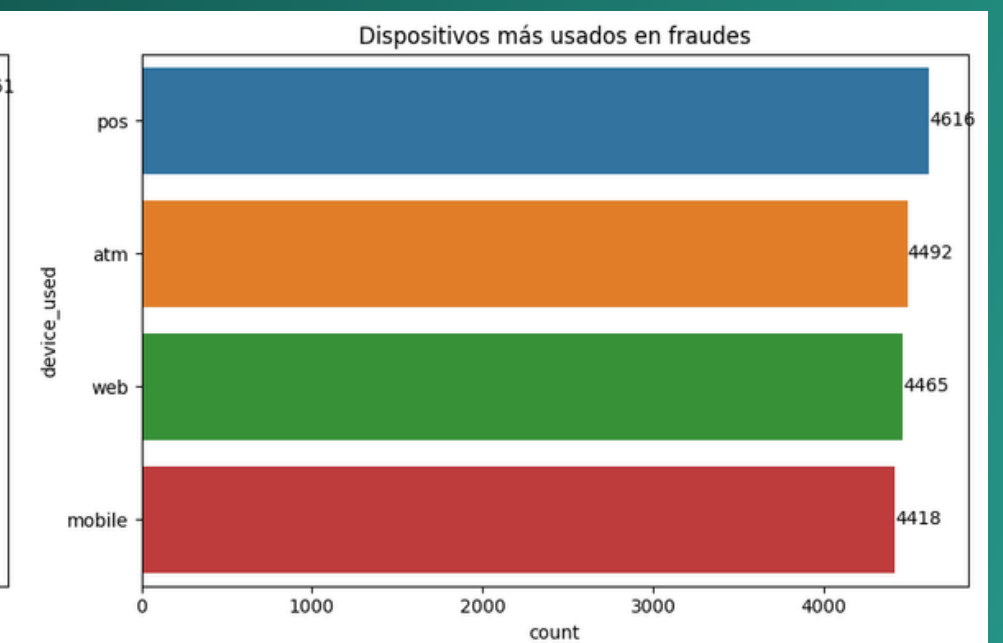
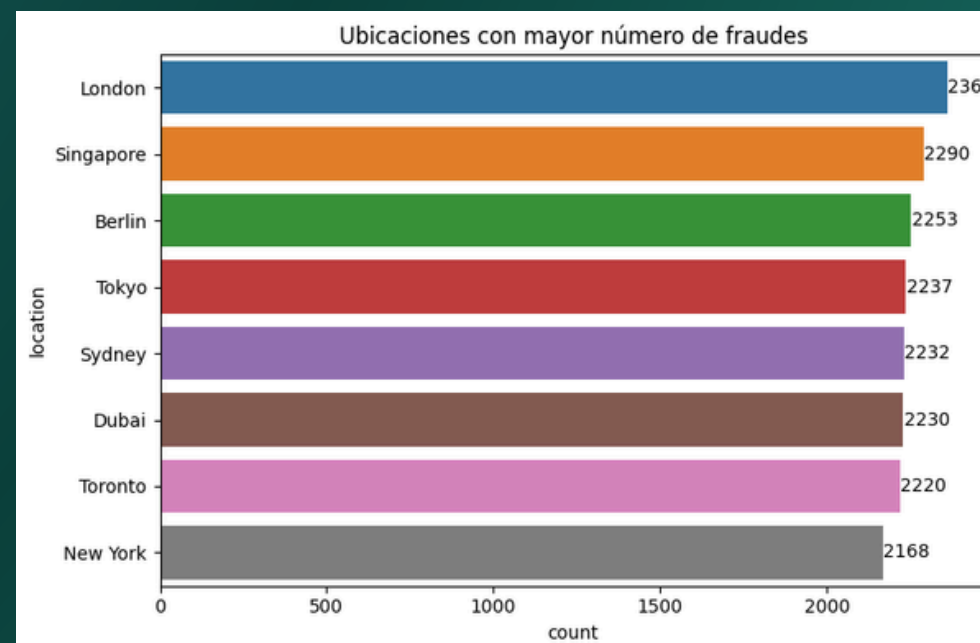
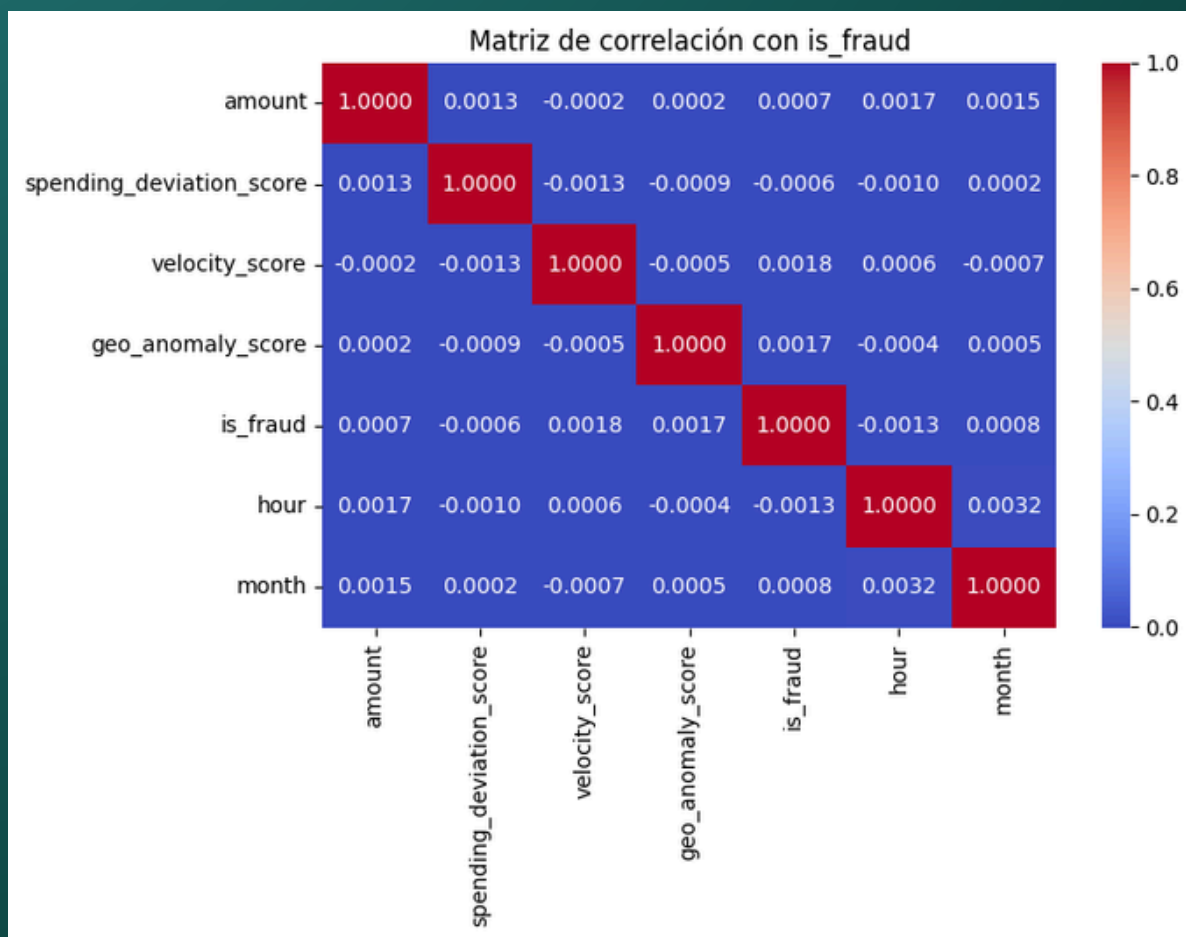
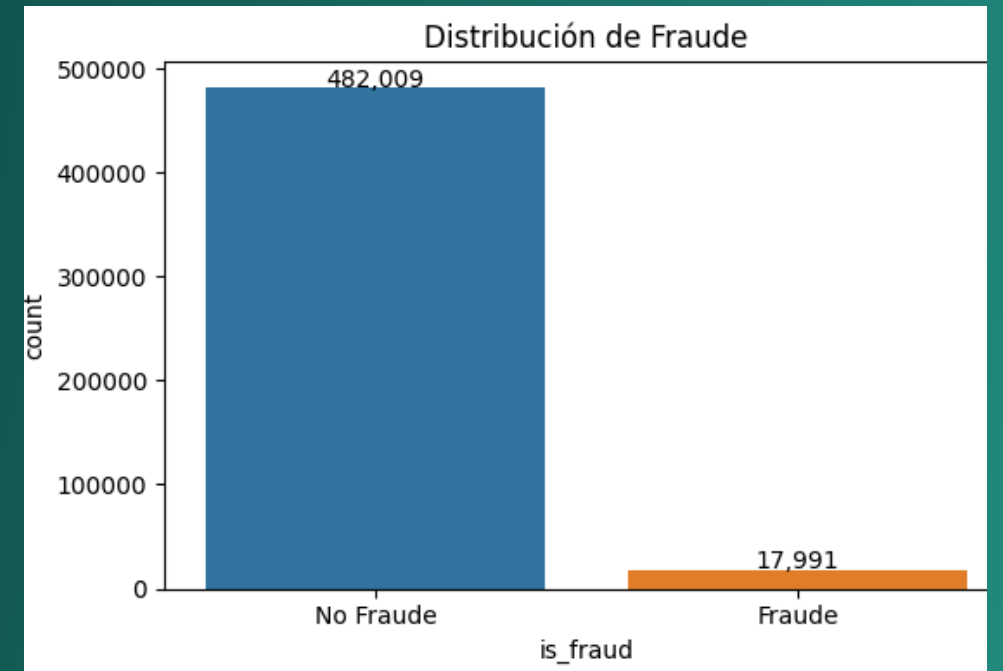
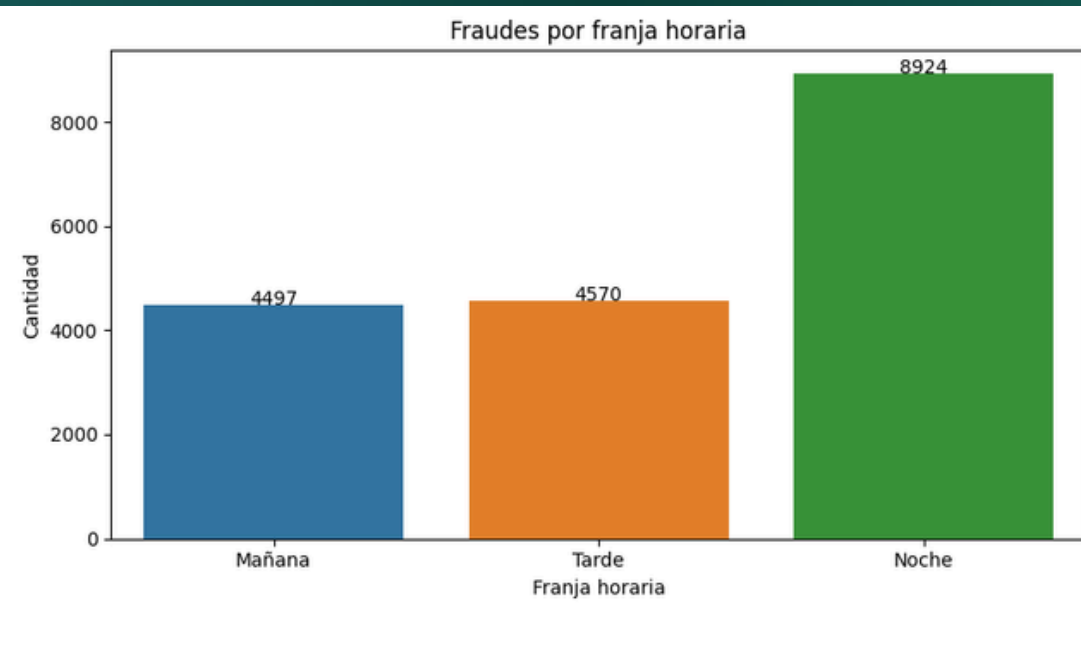
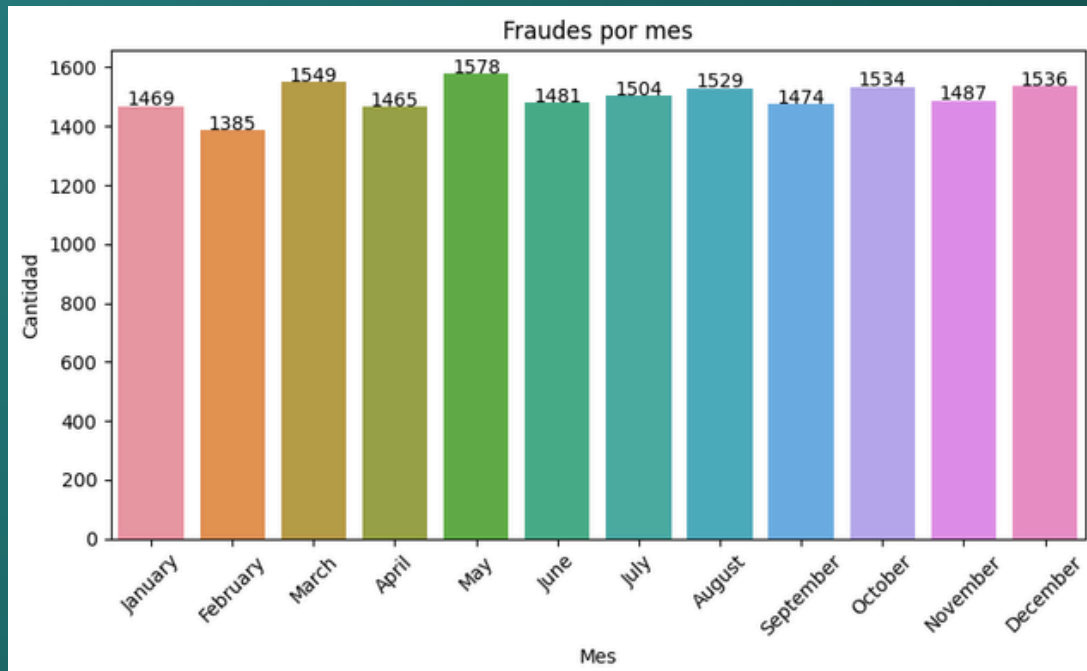
```

```

1 # Conversión de fechas
2 df['timestamp'] = pd.to_datetime(df['timestamp'], errors='coerce')
3
4 # Nueva columna: día, hora, día de la semana, mes y nombre del mes
5 df['date'] = df['timestamp'].dt.date
6 df['hour'] = df['timestamp'].dt.hour
7 df['day_of_week'] = df['timestamp'].dt.day_name()
8 df['month'] = df['timestamp'].dt.month
9 df['month_name'] = df['timestamp'].dt.month_name()
10
11 # Categorizar horas en franjas horarias
12 def categorize_time(hour):
13     if 0 <= hour < 6:
14         return 'Noche'
15     elif 6 <= hour < 12:
16         return 'Mañana'
17     elif 12 <= hour < 18:
18         return 'Tarde'
19     else:
20         return 'Noche'
21
22 df['time_period'] = df['hour'].apply(categorize_time)
23
24 # Orden de los meses
25 month_order = ['January', 'February', 'March', 'April', 'May', 'June',
26               'July', 'August', 'September', 'October', 'November', 'December']
27 time_order = ['Mañana', 'Tarde', 'Noche']
28
29 # Solo fraudes
30 fraudes = df[df['is_fraud'] == True]
31 fig, axes = plt.subplots(1, 2, figsize=(16, 5))
32
33 # Eliminar columnas que no agregan valor
34 df.drop(["fraud_type", "time_since_last_transaction"], axis=1, inplace=True)
35

```


EDA



Pre procesamiento

- Conversión y limpieza básica de datos (fechas y tipos numéricos).
- Creación de flags temporales: fin de semana y feriados.
- Codificación cíclica para hora y mes.
- Indicadores clave.
- Cálculo de tiempos entre transacciones por cuenta remitente.
- Estadísticas agregadas por cuenta (promedio, desviación, conteo).
- One-hot encoding (OHE) en variables categóricas.
- Escalado de variables numéricas.
- Eliminación de columnas irrelevantes.
- División y balanceo del dataset.

PRE PROCESAMIENTO

```
1 def preprocesamiento_fraude_big_data(
2     df: pd.DataFrame,
3     *,
4     target_col: str = "fraud_type",
5     max_low_card: int = 20, # Reducido para eficiencia
6     large_tx_percentile: float = 0.95,
7     test_size: float = 0.20,
8     random_state: int = 42,
9     balance_strategy: str = "undersample", # Mejor para datasets grandes
10    sample_size: int = None, # Para sampling inicial si es necesario
11    optimize_memory: bool = True,
12    verbose: bool = True
13 ):
14     if verbose:
15         print(f" Procesando dataset con {len(df):,} registros")
16         start_time = time.time()
17
18     # === SAMPLING PARA DESARROLLO (OPCIONAL) ===
19     if sample_size and len(df) > sample_size:
20         if verbose:
21             print(f" Tomando muestra estratificada de {sample_size:,} registros para desarrollo...")
22
23         y_full = df[target_col].astype(int)
24         df = df.groupby(target_col, group_keys=False).apply(
25             lambda x: x.sample(min(len(x), sample_size//2), random_state=random_state)
26         ).reset_index(drop=True)
27
28         if verbose:
29             print(f"✓ Muestra tomada: {len(df):,} registros")
30
31     # === OPTIMIZACIONES DE MEMORIA ===
32     if optimize_memory:
33         # Convertir object a category cuando sea eficiente
34         for col in df.select_dtypes(include=['object']).columns:
35             if col != target_col and df[col].nunique() / len(df) < 0.5:
36                 df[col] = df[col].astype('category')
37
38     def canon(s: str) -> str:
39         s = unicodedata.normalize("NFKC", str(s)).strip().lower()
40         return re.sub(r"\s+", " ", s)
41
42     label_rx = re.compile(r"(fraud|label|target|outcome|chargeback)", re.I)
43     df = df.copy()
```

```
1 # === LIMPIEZA NUMÉRICA EFICIENTE ===
2 def to_float_vectorized(series):
3     """Versión vectorizada más eficiente"""
4     if series.dtype in ['int64', 'float64']:
5         return series
6
7     # Convertir a string y limpiar
8     str_series = series.astype(str).str.replace(',', '').str.replace('$', '')
9     # Usar pd.to_numeric que es más eficiente
10    return pd.to_numeric(str_series, errors='coerce')
11
12    numeric_cols = ["amount", "time_since_last_transaction",
13                    "spending_deviation_score", "velocity_score", "geo_anomaly_score"]
14
15    for c in numeric_cols:
16        if c in df.columns:
17            if verbose and len(df) > 100000:
18                print(f" Procesando {c}...")
19            df[c] = to_float_vectorized(df[c])
20
21    # === TARGET PROCESSING ===
22    if target_col not in df.columns:
23        raise ValueError(f"target_col '{target_col}' no está en df.columns")
24
25    if df[target_col].dtype in ("object", "category"):
26        s = df[target_col].astype(str).str.strip().str.lower()
27        negatives = {"", "nan", "none", "null", "0", "false", "no", "normal"}
28        df[target_col] = (~s.isin(negatives)).astype(int)
29    else:
30        df[target_col] = pd.to_numeric(df[target_col], errors="coerce").fillna(0).astype(int)
```

PRE PROCESAMIENTO

```
1 # === CARACTERÍSTICAS TEMPORALES EFICIENTES ===
2 if "timestamp" in df.columns:
3     if verbose:
4         print(" Creando características temporales...")
5
6     df["timestamp"] = pd.to_datetime(df["timestamp"], errors="coerce")
7
8     # Solo las características más importantes para eficiencia
9     df["hour"] = df["timestamp"].dt.hour
10    df["day_of_week"] = df["timestamp"].dt.dayofweek
11
12    # Características binarias eficientes
13    df["is_business_hour"] = ((df["hour"] >= 9) & (df["hour"] <= 17)).astype(np.int8)
14    df["is_weekend"] = (df["day_of_week"] >= 5).astype(np.int8)
15    df["is_late_night"] = ((df["hour"] >= 22) | (df["hour"] <= 6)).astype(np.int8)
16
17    # Solo encoding cíclico esencial
18    df["hour_sin"] = np.sin(2 * np.pi * df["hour"] / 24)
19    df["hour_cos"] = np.cos(2 * np.pi * df["hour"] / 24)
```

```
1 # === SPLIT ESTRATIFICADO TEMPRANO ===
2 if verbose:
3     print(" Realizando split estratificado...")
4
5 y_full = df[target_col].astype(int)
6 idx_train, idx_test = train_test_split(
7     df.index, stratify=y_full, test_size=test_size, random_state=random_state
8 )
9 df_train = df.loc[idx_train].copy()
10 df_test = df.loc[idx_test].copy()
11
12 if verbose:
13     fraud_rate = y_full.mean()
14     print(f"✓ Train: {len(df_train):,} | Test: {len(df_test):,}")
15     print(f"✓ Tasa de fraude: {fraud_rate:.1%}")
```

```
1 # === FEATURE ENGINEERING EFICIENTE (SOLO LO ESENCIAL) ===
2 if {"sender_account", "amount"}.issubset(df.columns):
3     if verbose:
4         print(" Creando estadísticas por cuenta...")
5
6     # Solo estadísticas esenciales para eficiencia
7     stats = df_train.groupby("sender_account")["amount"].agg([
8         "mean", "count", "std"
9     ])
10    stats.columns = ["sender_avg_amount", "sender_tx_count", "sender_std_amount"]
11
12    # Feature derivada simple
13    stats["sender_amount_cv"] = stats["sender_std_amount"] / (stats["sender_avg_amount"] + 1e-6)
14
15    df_train = df_train.join(stats, on="sender_account")
16    df_test = df_test.join(stats, on="sender_account")
17
18    # Z-score de la transacción actual
19    df_train["amount_zscore"] = (df_train["amount"] - df_train["sender_avg_amount"]) / (df_train["sender_std_amount"] + 1e-6)
20    df_test["amount_zscore"] = (df_test["amount"] - df_test["sender_avg_amount"]) / (df_test["sender_std_amount"] + 1e-6)
21
22    # Características de dispositivo/ubicación simplificadas
23    if "device_used" in df.columns and "sender_account" in df.columns:
24        dev_counts = df_train.groupby("sender_account")["device_used"].nunique()
25        df_train["device_diversity"] = df_train["sender_account"].map(dev_counts).fillna(1)
26        df_test["device_diversity"] = df_test["sender_account"].map(dev_counts).fillna(1)
27        df_train["multiple_devices"] = (df_train["device_diversity"] > 1).astype(np.int8)
28        df_test["multiple_devices"] = (df_test["device_diversity"] > 1).astype(np.int8)
```

```
1 # === PREPARACIÓN DE CARACTERÍSTICAS ===
2 base_exclude = {
3     target_col, "is_fraud", "transaction_id", "sender_account", "receiver_account",
4     "ip_address", "device_hash", "timestamp", "device_used", "location"
5 }
6
7 label_like = {c for c in df.columns if label_rx.search(canon(c))}
8 label_like.discard(target_col)
9 no_modelar = (base_exclude | label_like) & set(df.columns)
10
11 drop_cols = [c for c in no_modelar if c in df.columns]
12 X_train = df_train.drop(columns=drop_cols, errors="ignore")
13 X_test = df_test.drop(columns=drop_cols, errors="ignore")
14 y_train = df_train[target_col].astype(int)
15 y_test = df_test[target_col].astype(int)
```

PRE PROCESAMIENTO

```
1 # === PROCESAMIENTO CATEGÓRICO EFICIENTE ===
2 if verbose:
3     print(" Procesando variables categóricas...")
4
5 all_cats = [c for c in X_train.columns if X_train[c].dtype in ("object", "category")]
6 cat_low = [c for c in all_cats if X_train[c].nunique(dropna=True) <= max_low_card]
7 cat_high = [c for c in all_cats if c not in cat_low]
8
9 # Filtrar label-like
10 cat_low = [c for c in cat_low if not label_rx.search(canon(c))]
11 cat_high = [c for c in cat_high if not label_rx.search(canon(c))]
12
13 base_num_cols = [c for c in X_train.columns if c not in all_cats]
14
15 # Imputación eficiente
16 for c in (cat_low + cat_high):
17     X_train[c] = X_train[c].fillna("Unknown")
18     X_test[c] = X_test[c].fillna("Unknown")
19
20 for c in base_num_cols:
21     if X_train[c].dtype in ['object', 'category']:
22         continue
23     median_val = X_train[c].median()
24     X_train[c] = X_train[c].fillna(median_val)
25     X_test[c] = X_test[c].fillna(median_val)
26
27 # Umbral de transacción grande
28 large_tx_threshold = None
29 if "amount" in X_train.columns:
30     large_tx_threshold = X_train["amount"].quantile(large_tx_percentile)
31     X_train["is_large_transaction"] = (X_train["amount"] > large_tx_threshold).astype(np.int8)
32     X_test["is_large_transaction"] = (X_test["amount"] > large_tx_threshold).astype(np.int8)
33     base_num_cols.append("is_large_transaction")
```

```
1 # === ENCODING EFICIENTE ===
2 artefactos = {}
3
4 # One-Hot Encoding solo para baja cardinalidad
5 ohe = None
6 trn_oh = tst_oh = None
7 if cat_low:
8     if verbose:
9         print(f" One-hot encoding para {len(cat_low)} columnas...")
10
11     try:
12         ohe = OneHotEncoder(handle_unknown="ignore", drop="first", sparse_output=True, dtype=np.int8)
13     except TypeError:
14         ohe = OneHotEncoder(handle_unknown="ignore", drop="first", sparse=True, dtype=np.int8)
15
16     ohe.fit(X_train[cat_low])
17     trn_oh = ohe.transform(X_train[cat_low])
18     tst_oh = ohe.transform(X_test[cat_low])
19     artefactos["ohe"] = ohe
20     artefactos["ohe_cols"] = list(cat_low)
21
22 # Frequency encoding para alta cardinalidad (más eficiente que target encoding)
23 fe_maps = {}
24 if cat_high:
25     if verbose:
26         print(f" Frequency encoding para {len(cat_high)} columnas...")
27
28     for c in cat_high:
29         freq = X_train[c].value_counts(normalize=True)
30         fe_maps[c] = freq.to_dict()
31         X_train[c + "_freq"] = X_train[c].map(fe_maps[c]).fillna(0.0)
32         X_test[c + "_freq"] = X_test[c].map(fe_maps[c]).fillna(0.0)
33
34 artefactos["freq_encoding_maps"] = fe_maps
35
36 # Eliminar categóricas originales
37 X_train_drop = X_train.drop(columns=cat_low + cat_high, errors="ignore")
38 X_test_drop = X_test.drop(columns=cat_low + cat_high, errors="ignore")
39
```



```
1 # === ESCALADO EFICIENTE ===
2 if verbose:
3     print(" Escalando características numéricas...")
4
5 scaler = RobustScaler()
6 scale_cols = [c for c in base_num_cols if c in X_train_drop.columns]
7
8 if scale_cols:
9     numeric_mask = X_train_drop[scale_cols].select_dtypes(include=[np.number]).columns
10    scale_cols = [c for c in scale_cols if c in numeric_mask]
11
12    if scale_cols:
13        scaler.fit(X_train_drop[scale_cols])
14        X_train_drop.loc[:, scale_cols] = scaler.transform(X_train_drop[scale_cols])
15        X_test_drop.loc[:, scale_cols] = scaler.transform(X_test_drop[scale_cols])
16
17 artefactos["scaler"] = scaler
18 artefactos["scale_cols"] = scale_cols
19
20 # Limpiar NaN finales
21 X_train_drop = X_train_drop.fillna(0)
22 X_test_drop = X_test_drop.fillna(0)
23
24 # === CREAR MATRICES SPARSE ===
25 X_train_dense = X_train_drop.copy()
26 X_test_dense = X_test_drop.copy()
27
28 X_train_num = sp.csr_matrix(X_train_drop.values, dtype=np.float32) # float32 para memoria
29 X_test_num = sp.csr_matrix(X_test_drop.values, dtype=np.float32)
30
31 if trn_oh is not None:
32     X_train_sparse = sp.hstack([X_train_num, trn_oh], format="csr")
33     X_test_sparse = sp.hstack([X_test_num, tst_oh], format="csr")
34 else:
35     X_train_sparse, X_test_sparse = X_train_num, X_test_num
```

```
1 # === BALANCEO OPTIMIZADO PARA BIG DATA ===
2 if verbose:
3     print(f" Aplicando estrategia de balanceo: {balance_strategy}")
4
5 y_train_res = y_train.copy()
6 X_train_sparse_res = X_train_sparse
7 X_train_dense_res = X_train_dense
8 sample_weight = None
9
10 if balance_strategy == "undersample":
11     # Undersample inteligente preservando diversidad
12     rus = RandomUnderSampler(
13         random_state=random_state,
14         sampling_strategy=0.3 # Ratio más conservador para preservar información
15     )
16     idx = np.arange(len(y_train_res)).reshape(-1, 1)
17     idx_res, y_train_res = rus.fit_resample(idx, y_train_res)
18     sel = idx_res.ravel()
19     X_train_sparse_res = X_train_sparse[sel]
20     X_train_dense_res = X_train_dense.iloc[sel]
21
22     if verbose:
23         print(f" Undersample: {len(y_train_res):,} registros finales")
24
25 elif balance_strategy == "weights":
26     sample_weight = compute_sample_weight(class_weight="balanced", y=y_train_res)
27     if verbose:
28         print(" Class weights aplicados")
29
30 else:
31     if verbose:
32         print(" Sin balanceo aplicado")
33
```

PRE PROCESAMIENTO

```
1  # === ARTEFACTOS FINALES ===
2  artefactos.update({
3      "large_tx_threshold": large_tx_threshold,
4      "base_num_cols": base_num_cols,
5      "cat_low": list(cat_low),
6      "cat_high": list(cat_high),
7      "drop_cols": drop_cols,
8      "sample_weight": sample_weight,
9      "balance_strategy": balance_strategy,
10     "X_train_dense": X_train_dense,
11     "X_test_dense": X_test_dense,
12     "X_train_dense_res": X_train_dense_res,
13 })
14
15 if verbose:
16     elapsed = time.time() - start_time
17     print(f" Preprocesamiento completado en {elapsed:.1f} segundos")
18     print(f" Shape final train: {X_train_sparse_res.shape}")
19     print(f" Shape final test: {X_test_sparse.shape}")
20
21 return X_train_sparse_res, X_test_sparse, y_train_res, y_test, artefactos
```

Optimización

- Se identificó que usar el umbral por defecto (0.5) en detección de fraude genera demasiados falsos positivos, baja precisión (3–5%) y sobrecarga operativa.
- Se diseñaron estrategias de optimización de umbrales con distintos enfoques: precision target, recall target, max F1, conservative y balanced.
- Se estableció un proceso de selección jerárquico, priorizando la precisión (precision target → conservative → max F1 → balanced) debido al alto costo de los falsos positivos frente al ahorro de los verdaderos positivos.
- Se definieron métricas de negocio clave: número total de alertas, porcentaje de precisión y estatus de cumplimiento de objetivos.
- Se modeló el impacto esperado al aplicar umbrales optimizados:
 - Precisión: aumento de 3–5% → 15–30%
 - Reducción de alertas: 70–90% menos
 - Mejora en F1-score
 - Cambio en ROI: de negativo a positivo

Optimización

```
1 def optimizar_precision_fraude(modelos, X_test_sparse, y_test, X_test_dense=None,
2                               target_precision=0.20, target_recall=0.60):
3
4     print(" OPTIMIZACIÓN DE PRECISIÓN Y UMBRAL")
5     print("=" * 60)
6     print(f" Target Precision: {target_precision:.1%}")
7     print(f" Target Recall: {target_recall:.1%}")
8     print()
9
10    resultados_optimizados = []
11    mejores_umbrales = {}
12
13    for nombre, modelo in modelos.items():
14        print(f" Optimizando {nombre}...")
15
16        # Seleccionar datos apropiados
17        if "Gradient Boosting" in nombre and X_test_dense is not None:
18            X_test_model = X_test_dense
19        else:
20            X_test_model = X_test_sparse
21
22        # Obtener probabilidades
23        if hasattr(modelo, 'predict_proba'):
24            y_probs = modelo.predict_proba(X_test_model)[: , 1]
25        else:
26            print(f" {nombre} no tiene predict_proba, usando decisión por defecto")
27            y_pred_default = modelo.predict(X_test_model)
28            precision_default = precision_score(y_test, y_pred_default, zero_division=0)
29            recall_default = recall_score(y_test, y_pred_default, zero_division=0)
30            f1_default = f1_score(y_test, y_pred_default, zero_division=0)
31
32            resultados_optimizados.append({
33                'Modelo': nombre,
34                'Umbral': 'default',
35                'Precision': f"{precision_default:.4f}",
36                'Recall': f"{recall_default:.4f}",
37                'F1-Score': f"{f1_default:.4f}",
38                'Status': 'No optimizable'
39            })
40            continue
41
42        # Calcular curva precision-recall
43        precisions, recalls, thresholds = precision_recall_curve(y_test, y_probs)
44        # Encontrar mejores umbrales según diferentes criterios
45        umbrales_candidatos = {}
```

```
1 # 1. Umbral para precision target
2 mask_precision = precisions >= target_precision
3 if mask_precision.any():
4     idx_precision = np.where(mask_precision)[0]
5     # Entre los que cumplen precision, elegir el de mayor recall
6     best_idx = idx_precision[np.argmax(recalls[idx_precision])]
7     umbrales_candidatos['precision_target'] = {
8         'threshold': thresholds[best_idx] if best_idx < len(thresholds) else 0.5,
9         'precision': precisions[best_idx],
10        'recall': recalls[best_idx],
11        'f1': 2 * precisions[best_idx] * recalls[best_idx] / (precisions[best_idx] + recalls[best_idx])
12    }
13
14 # 2. Umbral para recall target
15 mask_recall = recalls >= target_recall
16 if mask_recall.any():
17     idx_recall = np.where(mask_recall)[0]
18     # Entre los que cumplen recall, elegir el de mayor precision
19     best_idx = idx_recall[np.argmax(precisions[idx_recall])]
20     umbrales_candidatos['recall_target'] = {
21         'threshold': thresholds[best_idx] if best_idx < len(thresholds) else 0.5,
22         'precision': precisions[best_idx],
23         'recall': recalls[best_idx],
24         'f1': 2 * precisions[best_idx] * recalls[best_idx] / (precisions[best_idx] + recalls[best_idx])
25    }
26
27 # 3. Umbral que maximiza F1
28 f1_scores = 2 * precisions * recalls / (precisions + recalls + 1e-8)
29 best_f1_idx = np.argmax(f1_scores)
30 umbrales_candidatos['max_f1'] = {
31     'threshold': thresholds[best_f1_idx] if best_f1_idx < len(thresholds) else 0.5,
32     'precision': precisions[best_f1_idx],
33     'recall': recalls[best_f1_idx],
34     'f1': f1_scores[best_f1_idx]
35 }
36
37 # 4. Umbral balanceado (precision ~ recall)
38 balance_diff = np.abs(precisions - recalls)
39 balanced_idx = np.argmin(balance_diff)
40 umbrales_candidatos['balanced'] = {
41     'threshold': thresholds[balanced_idx] if balanced_idx < len(thresholds) else 0.5,
42     'precision': precisions[balanced_idx],
43     'recall': recalls[balanced_idx],
44     'f1': f1_scores[balanced_idx]
45 }
```

Optimización

```
1 # 5. Umbral conservador (alta precision, recall moderado)
2 conservative_mask = precisions >= 0.15 # Al menos 15% precision
3 if conservative_mask.any():
4     conservative_idx = np.where(conservative_mask)[0]
5     best_idx = conservative_idx[np.argmax(recalls[conservative_idx])]
6     umbrales_candidatos['conservative'] = {
7         'threshold': thresholds[best_idx] if best_idx < len(thresholds) else 0.7,
8         'precision': precisions[best_idx],
9         'recall': recalls[best_idx],
10        'f1': f1_scores[best_idx]
11    }
12
13 # Seleccionar el mejor umbral (priorizar precision target si existe)
14 if 'precision_target' in umbrales_candidatos:
15     mejor_umbral = umbrales_candidatos['precision_target']
16     criterio = 'precision_target'
17 elif 'conservative' in umbrales_candidatos:
18     mejor_umbral = umbrales_candidatos['conservative']
19     criterio = 'conservative'
20 elif 'max_f1' in umbrales_candidatos:
21     mejor_umbral = umbrales_candidatos['max_f1']
22     criterio = 'max_f1'
23 else:
24     mejor_umbral = umbrales_candidatos['balanced']
25     criterio = 'balanced'
26
27 mejores_umbrales[nombre] = mejor_umbral['threshold']
28
29 # Evaluar con el mejor umbral
30 y_pred_opt = (y_probs >= mejor_umbral['threshold']).astype(int)
31
32 # Métricas finales
33 precision_final = precision_score(y_test, y_pred_opt, zero_division=0)
34 recall_final = recall_score(y_test, y_pred_opt, zero_division=0)
35 f1_final = f1_score(y_test, y_pred_opt, zero_division=0)
36 roc_auc = roc_auc_score(y_test, y_probs)
37
38 # Conteos
39 tp = ((y_pred_opt == 1) & (y_test == 1)).sum()
40 fp = ((y_pred_opt == 1) & (y_test == 0)).sum()
41 fn = ((y_pred_opt == 0) & (y_test == 1)).sum()
42 tn = ((y_pred_opt == 0) & (y_test == 0)).sum()
```

```
1 # Métricas de negocio
2 precision_pct = precision_final * 100
3 total_alerts = tp + fp
4 fraud_detection_rate = tp / (tp + fn) if (tp + fn) > 0 else 0
5
6 status = " Cumple targets" if precision_final >= target_precision and recall_final >= target_recall else \
7         " Precision baja" if precision_final < target_precision else \
8         " Recall bajo" if recall_final < target_recall else " No cumple"
9
10 resultados_optimizados.append({
11     'Modelo': nombre,
12     'Umbral': f"{mejor_umbral['threshold']:.3f}",
13     'Criterio': criterio,
14     'Precision': f"{precision_final:.4f}",
15     'Recall': f"{recall_final:.4f}",
16     'F1-Score': f"{f1_final:.4f}",
17     'ROC-AUC': f"{roc_auc:.4f}",
18     'TP': tp,
19     'FP': fp,
20     'FN': fn,
21     'Total_Alerts': total_alerts,
22     'Precision_%': f"{precision_pct:.1f}%",
23     'Status': status
24 })
25
26 print(f"    Mejor umbral: {mejor_umbral['threshold']:.3f} ({criterio})")
27 print(f"    Precision: {precision_final:.1%} | Recall: {recall_final:.1%} | F1: {f1_final:.3f}")
28 print(f"    Alertas totales: {total_alerts:,} (TP: {tp:,}, FP: {fp:,})")
29 print()
30
31 return pd.DataFrame(resultados_optimizados), mejores_umbrales
```

Modelos y métricas

Modelos utilizados

- Logistic Regression
- Random Forest
- Gradient Boosting

Métricas utilizadas

- Precisión
- Sensibilidad
- F1-Score

Entrenamiento

```
1 def entrenar_con_class_weights(X_train_sparse, y_train, X_train_dense, sample_weight=None):
2     print(" RE-ENTRENANDO CON CLASS WEIGHTS (sin undersample)...")
3     print("=" * 60)
4
5     from sklearn.linear_model import LogisticRegression
6     from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
7
8     modelos_weights = {}
9
10    # Logistic Regression con class weight
11    print(" Logistic Regression con class_weight='balanced'...")
12    lr_balanced = LogisticRegression(
13        solver="saga",
14        penalty="l2",
15        C=0.01, # Más regularización para reducir overfitting
16        max_iter=500,
17        tol=1e-4,
18        class_weight='balanced', # + Clave para reducir FP
19        random_state=42,
20        n_jobs=-1
21    )
22    modelos_weights["Logistic Regression Balanced"] = lr_balanced.fit(X_train_sparse, y_train)
23
24    # Random Forest con class weight más conservador
25    print(" Random Forest con class_weight='balanced_subsample'...")
26    rf_balanced = RandomForestClassifier(
27        n_estimators=300,
28        max_depth=12, # Más conservador
29        min_samples_split=50, # Más restrictivo
30        min_samples_leaf=20, # Más restrictivo
31        max_features="sqrt",
32        bootstrap=True,
33        max_samples=0.6, # Menos muestras por árbol
34        class_weight='balanced_subsample',
35        n_jobs=-1,
36        random_state=42
37    )
38    modelos_weights["Random Forest Balanced"] = rf_balanced.fit(X_train_sparse, y_train)
39
40    # Gradient Boosting más conservador
41    if X_train_dense is not None:
42        print(" Gradient Boosting conservador...")
43        gbt_conservative = GradientBoostingClassifier(
44            n_estimators=200,
45            learning_rate=0.05, # Menor Learning rate
46            max_depth=4, # Menor profundidad
47            subsample=0.7,
48            max_features="sqrt",
49            min_samples_split=100, # Más restrictivo
50            min_samples_leaf=50, # Más restrictivo
51            random_state=42
52        )
53        modelos_weights["Gradient Boosting Conservative"] = gbt_conservative.fit(X_train_dense, y_train)
54
55    print(f" {len(modelos_weights)} modelos re-entrenados con class weights")
56    return modelos_weights
```

Optimizando Logistic Regression Balanced...

Mejor umbral: 0.500 (precision_target)

Precision: 3.6% | Recall: 33.9% | F1: 0.066

Alertas totales: 33,401 (TP: 1,218, FP: 32,183)

Optimizando Random Forest Balanced...

Mejor umbral: 0.500 (precision_target)

Precision: 4.4% | Recall: 51.9% | F1: 0.080

Alertas totales: 42,896 (TP: 1,866, FP: 41,030)

Optimizando Gradient Boosting Conservative...

Mejor umbral: 0.204 (precision_target)

Precision: 15.4% | Recall: 0.1% | F1: 0.001

Alertas totales: 13 (TP: 2, FP: 11)

Análisis ROI

Conceptos clave

Verdadero Positivo (TP): Fraude detectado correctamente

- Beneficio: Evitamos pérdida por fraude + costos de investigación

Falso Positivo (FP): Transacción legítima marcada como fraude

- Costo: Fricción al cliente + tiempo de revisión + posible pérdida de cliente

Falso Negativo (FN): Fraude no detectado

- Costo implícito: Pérdida total del fraude (no calculado aquí directamente)

Verdadero Negativo (TN): Transacción legítima correctamente clasificada

- Sin costo ni beneficio directo

Análisis ROI

```
1 def analisis_costo_beneficio(resultados_df, costo_fp=100, beneficio_tp=1000):
2     """
3     Análisis de costo-beneficio para seleccionar el mejor modelo
4     """
5     print(f"\n ANÁLISIS COSTO-BENEFICIO")
6     print("=" * 40)
7     print(f"Costo por Falso Positivo: ${costo_fp}")
8     print(f"Beneficio por Verdadero Positivo: ${beneficio_tp}")
9     print()
10
11     analisis = []
12
13     for _, row in resultados_df.iterrows():
14         if row['Status'] == 'No optimizable':
15             continue
16
17         tp = int(row['TP'])
18         fp = int(row['FP'])
19
20         costo_total = fp * costo_fp
21         beneficio_total = tp * beneficio_tp
22         beneficio_neto = beneficio_total - costo_total
23         roi = (beneficio_neto / costo_total * 100) if costo_total > 0 else 0
24
25         analisis.append({
26             'Modelo': row['Modelo'],
27             'Beneficio_Total': f"${beneficio_total:,}",
28             'Costo_Total': f"${costo_total:,}",
29             'Beneficio_Neto': f"${beneficio_neto:,}",
30             'ROI_%': f"{roi:.1f}%",
31             'Precision_%': row['Precision_%']
32         })
33
34     df_analisis = pd.DataFrame(analisis)
35     print(df_analisis.to_string(index=False))
36
37     if len(df_analisis) > 0:
38         mejor_roi_idx = df_analisis['ROI_%'].str.rstrip('%').astype(float).idxmax()
39         mejor_modelo_negocio = df_analisis.iloc[mejor_roi_idx]
40         print(f"\n MEJOR MODELO POR ROI: {mejor_modelo_negocio['Modelo']}")
41         print(f" ROI: {mejor_modelo_negocio['ROI_%']}")
42         print(f" Beneficio Neto: {mejor_modelo_negocio['Beneficio_Neto']}")
43
44     return df_analisis
```

Costo por Falso Positivo: \$100

Beneficio por Verdadero Positivo: \$1000

	Modelo	Beneficio_Total	Costo_Total	Beneficio_Neto	ROI_%	Precision_%
	Logistic Regression Balanced	\$1,218,000	\$3,218,300	\$-2,000,300	-62.2%	3.6%
	Random Forest Balanced	\$1,866,000	\$4,103,000	\$-2,237,000	-54.5%	4.4%
	Gradient Boosting Conservative	\$2,000	\$1,100	\$900	81.8%	15.4%

MEJOR MODELO POR ROI: Gradient Boosting Conservative

ROI: 81.8%

Beneficio Neto: \$900

Flujo completo y resultados

- Pipeline con optimización
- Resultados por fases

Flujo y resultados

```
1 def pipeline_optimizacion_completa(df, target_col="fraud_type"):
2     """
3     Pipeline completo: re-entrenamiento + optimización de umbral
4     """
5
6     print(" PIPELINE COMPLETO DE OPTIMIZACIÓN")
7     print("=" * 70)
8
9     # 1. Re-preprocesamiento con class weights
10    print("\n=== FASE 1: RE-PREPROCESAMIENTO CON CLASS WEIGHTS ===")
11
12    # Importar función de preprocesamiento (asumiendo que ya existe)
13    X_train_res, X_test_tr, y_train_res, y_test, artefactos = preprocesamiento_fraude_big_data(
14        df,
15        target_col=target_col,
16        balance_strategy="weights", # + Cambio clave
17        max_low_card=15,
18        optimize_memory=True,
19        verbose=True
20    )
21
22    # 2. Re-entrenamiento con configuraciones más conservadoras
23    print("\n=== FASE 2: RE-ENTRENAMIENTO CONSERVADOR ===")
24    modelos_optimizados = entrenar_con_class_weights(
25        X_train_sparse=X_train_res,
26        y_train=y_train_res,
27        X_train_dense=artefactos["X_train_dense_res"]
28    )
29
30    # 3. Optimización de umbrales
31    print("\n=== FASE 3: OPTIMIZACIÓN DE UMBRALES ===")
32    resultados_opt, umbrales = optimizar_precision_fraude(
33        modelos_optimizados,
34        X_test_sparse=X_test_tr,
35        y_test=y_test,
36        X_test_dense=artefactos["X_test_dense"],
37        target_precision=0.15, # 15% precision mínima
38        target_recall=0.50    # 50% recall mínimo
39    )
40
41    # 4. Análisis de negocio
42    print("\n=== FASE 4: ANÁLISIS DE NEGOCIO ===")
43    analisis_negocio = analisis_costo_beneficio(resultados_opt)
44
45    # 5. Comparación con resultados anteriores
46    print("\n=== COMPARACIÓN CON RESULTADOS ANTERIORES ===")
47    print(" ANTES (con undersample):")
48    print("   Precision: 0.035-0.044 (3.5-4.4%)")
49    print("   Recall: 0.63-0.69")
50    print("   F1: 0.066-0.089")
51    print("   Falsos Positivos: 46K+")
52    print()
53    print(" AHORA (con class weights + umbral optimizado):")
54    print(resultados_opt[['Modelo', 'Precision_%', 'Recall', 'F1-Score', 'Total_Alerts', 'Status']].to_string(index=False))
55
56    return resultados_opt, umbrales, analisis_negocio
```

=== COMPARACIÓN CON RESULTADOS ANTERIORES ===

ANTES (con undersample):

Precision: 0.035-0.044 (3.5-4.4%)

Recall: 0.63-0.69

F1: 0.066-0.089

Falsos Positivos: 46K+

AHORA (con class weights + umbral optimizado):

	Modelo	Precision_%	Recall	F1-Score	Total_Alerts	Status
	Logistic Regression Balanced	3.6%	0.3385	0.0658	33401	Precision baja
	Random Forest Balanced	4.4%	0.5186	0.0883	42896	Precision baja
	Gradient Boosting Conservative	15.4%	0.0006	0.0011	13	Recall bajo

RECOMENDACIÓN FINAL:

1. Usar el modelo con mejor ROI del análisis de negocio
2. Aplicar el umbral optimizado correspondiente
3. Monitorear precision en producción y ajustar si es necesario

Desafíos y presentación del código

- Manejo de desbalance de los datos
- Decision entre Series de Panda vs Dataframes de Spark
- Limitaciones de recursos para manejar gran cantidad de datos con la version Free Edition

Código completo...
Ingresa aquí