

PROYECTO: ENTREGA N°1

NATALIA ANDREA GARCÍA RÍOS

MATEO VELÁSQUEZ RODRÍGUEZ

1. PROBLEMA PREDICTIVO POR RESOLVER

Se busca predecir si una transacción bancaria es fraudulenta o no, utilizando información detallada sobre la transacción y la identidad del usuario. La información incluye datos como el tiempo transcurrido desde una fecha de referencia, el importe de la transacción, el código del producto, información de la tarjeta de pago, dirección, distancia, dominio de correo electrónico, características de la tarjeta y dirección, así como otras características específicas diseñadas por Vesta. También se cuenta con información sobre la identidad del usuario, como la conexión de red y la firma digital asociada con la transacción. Estos datos son recolectados por el sistema de protección contra fraudes de Vesta y sus socios de seguridad digital.

2. DATASET A UTILIZAR

Se usará el dataset de la competición [IEEE-CIS Fraud Detection](#) de kaggle. Dicho dataset cuenta con cinco archivos:

1. **train_transaction.csv**: Es uno de los dos conjuntos de datos de entrenamiento que se utilizará para construir el modelo predictivo, contiene la información sobre las transacciones, se compone de 50000 movimientos y posee aproximadamente 51 columnas, las cuales son:

- *TransactionDT*: Timedelta a partir de una fecha y hora de referencia dadas.
- *TransactionAMT*: Importe del pago de la transacción en USD.
- *ProductCD*: Código de producto, el producto de cada transacción.
- card1 - card6: Información de la tarjeta de pago, como tipo de tarjeta, categoría de tarjeta, banco emisor, país, etc.
- *addr*: Dirección.
- *dist*: Distancia.
- *P_y (R_) emaildomain*: Dominio de correo electrónico del comprador y del destinatario.
- *C1-C14*: Recuento, como cuántas direcciones se encuentran asociadas a la tarjeta de pago, etc. El significado real está enmascarado.
- *D1-D15*: Los días transcurridos entre la transacción anterior.

- *M1-M9*: Coincidencia, como los nombres de la tarjeta y la dirección, etc.
 - *Vxxx*: Características ricas diseñadas por Vesta, que incluyen clasificación, recuento y otras relaciones entre entidades.
2. **train_identity.csv**: Es uno de los dos conjuntos de datos de entrenamiento que se utilizará para construir el modelo predictivo, contiene datos sobre la información de identidad - información de conexión de red (IP, ISP, Proxy, etc.) y firma digital (UA/navegador/sistema operativo/versión, etc.) asociada con las transacciones. Estos datos son recolectados por el sistema de protección contra fraudes de Vesta y por socios de seguridad digital. Contiene 41 columnas, que son:
- *TransactionID*
 - *id_01 – id_38*
 - *DeviceType*
 - *DeviceInfo*
3. **Test_transaction.csv**: Es uno de los dos conjuntos de datos de prueba que se utilizará para evaluar el modelo predictivo construido. En este conjunto de datos, se proporcionan las mismas características que en el conjunto de datos de entrenamiento, pero se omite la variable objetivo "isFraud".
4. **test_identity.csv**: Es uno de los dos conjuntos de datos de prueba que se utilizará para evaluar el modelo predictivo construido. En este conjunto de datos, se proporcionan las mismas características que en el conjunto de datos de entrenamiento, pero se omite la variable objetivo "isFraud".
5. **sample_submission.csv**: Es un archivo de muestra que muestra el formato correcto para enviar las predicciones del modelo construido para el conjunto de datos de prueba.

3. MÉTRICAS DE DESEMPEÑO

Las métricas de desempeño requeridas para evaluar el modelo de detección de transacciones fraudulentas pueden ser:

Métricas de machine learning:

- *Accuracy (Precisión)*: Mide la proporción de transacciones identificadas como fraudulentas que son realmente fraudulentas.

$$accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

- *Recall o Sensitivity (Recuperación)*: Mide la proporción de transacciones fraudulentas que el modelo identifica correctamente.

$$recall = \frac{TP}{TP + FN}$$

- *F1 Score*: Combina la precisión y la recuperación para proporcionar una métrica general del rendimiento del modelo.

$$F1\ Score = \frac{2 \cdot accuracy \cdot recall}{(accuracy + recall)}$$

- *Specifity*: Proporciona la probabilidad de que un se dé un resultado no fraudulento, condicionado a que sea realmente fraudulento.

$$Specifity = \frac{TN}{TN + FN}$$

Métricas de negocio:

- *Tasa de falsos positivos (False Positive Rate)*: Mide la proporción de transacciones que son identificadas erróneamente como fraudulentas, lo que puede generar un impacto negativo en la experiencia del usuario y en la reputación de la empresa.

$$False\ Positive\ Rate = \frac{FP}{FP + TN}$$

- *Tasa de verdaderos positivos (True Positive Rate)*: Mide la proporción de transacciones fraudulentas que son identificadas correctamente, lo que ayuda a prevenir pérdidas financieras.

$$True\ Positive\ Rate = \frac{TP}{TP + FN}$$

4. DESEMPEÑO DESEABLE EN PRODUCCIÓN

El desempeño deseable en producción depende mucho de las necesidades de la entidad, los niveles de riesgo que se esté dispuesto a asumir y los costos de los falsos positivos y falsos negativos. Ejemplificando, si tenemos un costo alto de falsos positivos, lo principal sería maximizar la precisión, mientras que, si el costo de los falsos negativos es alto, se buscará maximizar la recuperación, lo que conlleva a aplicar el modelo entrenado a un nuevo conjunto de datos y ajustarlo para hacer predicciones puntuales en estos datos.

En resumen, un buen desempeño se puede considerar si el modelo logra un equilibrio adecuado entre la precisión y la recuperación, y si las tasas de falsos positivos y falsos negativos son aceptables para los objetivos de la empresa.

Lo anterior sería un caso hipotetico en donde éste modelo se ajuste a los requerimientos de una empresa.