

MATRIZ DE CONSISTENCIA

Implementación de un Modelo de Machine Learning
para la Detección de Transacciones Fraudulentas y Anómalas
en Pagos Digitales de la Empresa TechSport

Gestión 2025

Ing. Ada Condori Callisaya

Noviembre 2025

Datos Generales de la Investigación

Programa	Maestría en Dirección Estratégica en Ingeniería de Software
Universidad	UAGRM - Facultad de Ingeniería en Ciencias de la Computación y Telecomunicaciones
Autor	Ing. Ada Condori Callisaya
Periodo	Gestión 2025 (2 meses de ejecución)
Tipo	Investigación Aplicada-Tecnológica
Enfoque	Cuantitativo
Diseño	Cuasiexperimental Retrospectivo
Alcance	Descriptivo-Correlacional-Comparativo

Método AQP (Rosario Martínez)

A	ADÓNDE (Lugar de Estudio)
	<p>Empresa: TechSport (nombre ficticio por seguridad)</p> <p>Ubicación: Miami, Florida, Estados Unidos</p> <p>Tipo: Empresa tecnológica SaaS especializada en gestión de instalaciones deportivas</p> <p>Operación: Plataforma multicanal para reservas deportivas, membresías y pagos digitales</p> <p>Pasarelas de pago: 10+ gateways (Stripe, CardConnect, Kushki, AzulPay, RazorPay, BAC)</p> <p>Acceso: Completo a datos transaccionales históricos con autorización y NDAs</p>
Q	QUIÉNES O QUÉ (Objeto de Análisis)
	<p>Objeto de estudio: Transacciones de pago digitales</p> <p>Periodo: Gestión 2025</p> <p>Volumen total: 15,492,846 transacciones</p> <p>Justificación: Dataset homogéneo que garantiza validez interna y evita temporal drift</p> <p>Cobertura poblacional: Censo completo de transacciones procesadas en 2025</p> <p>Categorías: Reservas deportivas, membresías, clínicas, cargos recurrentes, pagos one-time</p> <p>Métodos de pago: Tarjetas crédito/débito, ACH, créditos prepagados, wallets digitales</p> <p>Canales: Web, aplicación móvil, puntos de venta (POS)</p>
P	PROBLEMA (Variable Madre)
	<p>Problema identificado: TRANSACCIONES FRAUDULENTAS Y ANÓMALAS EN PAGOS DIGITALES</p> <p>Las transacciones procesadas por TechSport presentan comportamientos fraudulentos y anómalos</p> <p>NO detectados oportunamente, debido a:</p> <ol style="list-style-type: none"> 1. Limitaciones del sistema actual basado en reglas estáticas 2. Ausencia de modelos predictivos de Machine Learning 3. Imposibilidad de detección temprana (detección post-mortem) 4. Fragmentación de la arquitectura multicanal 5. Falta de aprendizaje continuo <p>Subcomponentes del problema:</p> <ul style="list-style-type: none"> • Patrones anómalos en comportamientos de pago • Fraude financiero consumado (chargebacks, disputas) • Riesgos transaccionales no mitigados • Pérdidas económicas por fraude no detectado • Falsos positivos que rechazan pagos legítimos

Método CCA (Rosario Martínez)

C	CAUSAS del Problema
	<p>¿Por qué TechSport NO detecta eficazmente fraude en pagos transaccionales multicanal?</p> <ol style="list-style-type: none"> 1. Sistema actual basado en reglas estáticas que no aprenden de nuevos patrones 2. Ausencia de modelos predictivos de Machine Learning 3. Falta de correlación entre comportamientos en diferentes gateways 4. Detección post-mortem (después de consumado el fraude) 5. Arquitectura multicanal fragmentada sin unificación de criterios de riesgo 6. Actualización manual de reglas de detección 7. No hay scoring dinámico de riesgo 8. Sistema no se adapta a nuevas modalidades de fraude
C	<p>¿Qué podría pasar si el problema continúa sin solución?</p> <ol style="list-style-type: none"> 1. Incremento de pérdidas económicas por fraude no detectado 2. Aumento de chargebacks y disputas con procesadores de pago 3. Sanciones de pasarelas de pago por alta tasa de fraude 4. Pérdida de confianza de usuarios legítimos 5. Reputación empresarial dañada 6. Costos operativos crecientes de revisión manual 7. Saturación del equipo de contabilidad con alertas falsas 8. Riesgo de exclusión de procesadores de pago (blacklist) 9. Impacto en la industria fintech por fraude no mitigado
A	<p>Solución: Implementación de un Modelo de Machine Learning Supervisado</p> <p>Algoritmo principal: Random Forest Algoritmos alternativos: XGBoost, SVM</p> <p>Características de la solución:</p> <ul style="list-style-type: none"> • Análisis de 15.4M+ transacciones históricas (Gestión 2025) • Feature engineering con 15+ features comportamentales • Validación estratificada (Train 70 %, Validation 15 %, Test 15 %) • Balanceo de clases adaptativo (SMOTE) • Optimización de hiperparámetros mediante Grid Search <p>Metas cuantificables:</p> <ul style="list-style-type: none"> • F1-Score \geq 85 % • Recall \geq 90 % (prioridad: detectar fraudes) • Precision \geq 80 % • AUC-ROC \geq 0.92 ³ • Tiempo de inferencia <200ms por transacción

Variables de la Investigación

Variable Dependiente (VD)

Variable	Transacciones Fraudulentas y Anómalas
Tipo	Dependiente (Variable Madre - del método AQP)
Definición conceptual	Registros transaccionales de pagos digitales que presentan comportamientos anómalos o fraudulentos, identificados mediante análisis de patrones históricos y etiquetado post-mortem
Definición operacional	Clasificación binaria (Fraude/No Fraude) de transacciones basada en etiquetado del equipo de contabilidad, chargebacks y disputas reportadas
Dimensiones	<ol style="list-style-type: none"> 1. Tipo de fraude detectado 2. Severidad del fraude 3. Canal de ocurrencia
Indicadores	<ul style="list-style-type: none"> • Tasa de fraude (%) • Pérdidas económicas (USD) • Precision del sistema (%) • Recall del sistema (%) • F1-Score • AUC-ROC • Tasa de falsos positivos (FPR) • Tasa de falsos negativos (FNR)
Instrumento	<ul style="list-style-type: none"> • Dataset histórico etiquetado (25M+ transacciones) • Matriz de confusión • Métricas de clasificación (scikit-learn)

Variante Independiente (VI)

Variable	Modelo de Machine Learning Implementado
Tipo	Independiente (Solución propuesta - del método CCA, .^"de Aporte)
Definición conceptual	Algoritmo computacional basado en aprendizaje automático supervisado, capaz de analizar datos históricos de transacciones etiquetadas para identificar patrones asociados a fraude y predecir la probabilidad de que nuevas transacciones sean fraudulentas o legítimas
Definición operacional	Modelo de clasificación binaria (Fraude/No Fraude) entrenando con dataset histórico de TechSport, que genera un score de riesgo para cada transacción y una clasificación final basada en un umbral optimizado
Dimensiones	<ol style="list-style-type: none"> 1. Tipo de algoritmo implementado 2. Estrategia de entrenamiento 3. Features utilizadas
Indicadores	<ul style="list-style-type: none"> • Algoritmo seleccionado (Random Forest/XGBoost/SVM) • Hiperparámetros optimizados (max_depth, n_estimators) • Balance del dataset (ratio fraude/no fraude) • Error de entrenamiento (%) • Error de validación (%) • Tiempo de entrenamiento (minutos) • Tiempo de inferencia (milisegundos) • Tamaño del modelo (MB) • Número de features utilizadas (≥ 15)
Instrumento	<ul style="list-style-type: none"> • Scripts Python (scikit-learn, pandas) • Logs de entrenamiento • Modelo serializado (.pkl/.joblib) • Código versionado en GitHub

MATRIZ DE CONSISTENCIA COMPLETA

Tabla 1: Matriz de Consistencia Metodológica de la Investigación

PROBLEMAS	OBJETIVOS	HIPÓTESIS	VARIABLES	METODOLOG
PROBLEMA GENERAL	OBJETIVO GENERAL	HIPÓTESIS GENERAL	V. INDEPENDIENTE	DISEÑO
PROBLEMAS ESPECÍFICOS	OBJETIVOS ESPECÍFICOS	HIPÓTESIS ESPECÍFICAS	V. DEPENDIENTE	TÉCNICAS
¿Cómo mejorar la detección de transacciones fraudulentas y anómalas en pagos digitales de la empresa TechSport durante la gestión 2025?	Implementar un modelo de Machine Learning supervisado basado en Random Forest para la detección de transacciones fraudulentas y anómalas en pagos digitales, logrando F1-Score $\geq 85\%$, Recall $\geq 90\%$, Precision $\geq 80\%$, en la empresa TechSport, gestión 2025.	La implementación de un modelo de Machine Learning supervisado basado en Random Forest alcanza un F1-Score mínimo del 85 %, con Recall $\geq 90\%$ y Precision $\geq 80\%$, en la detección de transacciones fraudulentas del test set (gestión 2025).	<p>VI: Modelo de ML implementado</p> <p>Indicadores:</p> <ul style="list-style-type: none"> ▪ Algoritmo: Random Forest ▪ F1-Score $\geq 85\%$ ▪ Recall $\geq 90\%$ ▪ Precision $\geq 80\%$ ▪ Tiempo inferencia <200ms 	<p>Tipo: Aplicada-tecnológica</p> <p>Enfoque: Cuantitativo</p> <p>Diseño: Cuasi-experimental retrospectivo</p> <p>Alcance: Descriptivo-correlacional-comparativo</p>

Continúa en la siguiente página

Tabla 1 – Continuación de la página anterior

PROBLEMAS	OBJETIVOS	HIPÓTESIS	VARIABLES	METODOLOG
PE1: ¿Cuáles son los fundamentos teóricos de los modelos de ML supervisados aplicados a detección de fraude en pagos digitales según literatura 2020-2025?	OE1: Fundamentar teóricamente los modelos de ML supervisados aplicados a detección de fraude, las métricas de evaluación y las técnicas de feature engineering.	HE1: Los modelos de ML supervisados (Random Forest, XGBoost, SVM) constituyen un enfoque teórico-técnico validado en literatura 2020-2025, superando limitaciones de sistemas basados en reglas estáticas.	VD: Transacciones fraudulentas y anómalas Indicadores: <ul style="list-style-type: none"> ■ Tasa de fraude (%) ■ Pérdidas (USD) ■ Precision (%) ■ Recall (%) ■ F1-Score ■ Tasa falsos positivos 	Revisión bibliográfica: <ul style="list-style-type: none"> ■ 30+ artículos científicos ■ Bases: IEEE, ACM, Scopus ■ Periodo: 2020-2025 ■ Análisis de benchmarks

Continúa en la siguiente página

Tabla 1 – Continuación de la página anterior

PROBLEMAS	OBJETIVOS	HIPÓTESIS	VARIABLES	METODOLOG
PE2: ¿Cuál es la situación actual del sistema de detección de fraude de TechSport al analizar el dataset histórico de gestión 2025?	OE2: Diagnosticar la situación actual del sistema de detección de fraude mediante análisis exploratorio del dataset 2025, documentando el proceso de etiquetado y caracterizando patrones de fraude.	HE2: El sistema actual presenta limitaciones evidenciadas por transacciones fraudulentas no detectadas oportunamente. El análisis exploratorio revelará al menos 3 patrones de fraude recurrentes.	<p>Intervinientes:</p> <ul style="list-style-type: none"> ■ Canal de pago ■ Tipo de transacción ■ Gateway ■ Volumen transaccional 	<p>Análisis Exploratorio:</p> <ul style="list-style-type: none"> ■ EDA con Python/pandas ■ Visualizaciones (matplotlib) ■ Estadística descriptiva ■ Caracterización de fraudes

Continúa en la siguiente página

Tabla 1 – Continuación de la página anterior

PROBLEMAS	OBJETIVOS	HIPÓTESIS	VARIABLES	METODOLOG
PE3: ¿Cómo desarrollar un modelo de ML que clasifique transacciones fraudulentas con alta precisión y recall?	OE3: Desarrollar un modelo de ML supervisado mediante: (i) preprocesamiento de 25M+ transacciones, (ii) feature engineering de 15+ features evitando data leakage, (iii) balanceo SMOTE, (iv) validación temporal, y (v) Grid Search.	HE3: Un modelo Random Forest entrenado con dataset balanceado y 15+ features comportamentales puede clasificar transacciones fraudulentas con Recall $\geq 90\%$, Precision $\geq 80\%$ y AUC-ROC ≥ 0.92 .	Métodos: <ul style="list-style-type: none"> ■ Preprocesamiento ■ Feature engineering ■ Validación temporal ■ Optimización hiperparámetros 	Pipeline ML: <ul style="list-style-type: none"> ■ Scikit-learn ■ Train: 2024 (9.7M trans.) ■ Test: 2025 (15.5M trans.) ■ SMOTE balancing ■ Grid Search

Continúa en la siguiente página

Tabla 1 – Continuación de la página anterior

PROBLEMAS	OBJETIVOS	HIPÓTESIS	VARIABLES	METODOLOG
PE4: ¿Qué nivel de efectividad presenta el modelo ML comparado con benchmarks de literatura científica?	OE4: Evaluar el desempeño del modelo mediante métricas de clasificación sobre test set temporal, comparándolo con benchmarks reportados en literatura.	HE4: El modelo alcanza F1-Score de 85-90 % con Recall \geq 90 % y Precision \geq 80 %, demostrando desempeño comparable o superior a benchmarks (Hafez et al., 2025: 85-94 %), manteniendo tiempo inferencia <200ms.	<p>Métricas:</p> <ul style="list-style-type: none"> ■ F1-Score, Recall, Precision ■ AUC-ROC ■ Comparación benchmarks ■ Intervalos confianza bootstrap ■ Tiempo inferencia (ms) 	<p>Evaluación:</p> <ul style="list-style-type: none"> ■ Matriz de confusión ■ Curva ROC ■ Bootstrap (1000 muestras) ■ Comparación literatura ■ Análisis feature importance

Población y Muestra

Población

Definición: Todas las transacciones de pago procesadas por TechSport en su plataforma multicanal durante la gestión 2025.

Características:

- **Tamaño poblacional:** Aproximadamente 15.5 millones de transacciones (gestión 2025)
- **Periodo:** Gestión 2025 (12 meses)
- **Canales:** Web, aplicación móvil, puntos de venta (POS)
- **Gateways:** 10+ pasarelas de pago integradas

Muestra

Tipo de muestreo: Censo de transacciones históricas (NO es muestra aleatoria)

Tamaño del dataset: 15,492,846 transacciones (censo completo de gestión 2025)

Justificación metodológica: El uso exclusivo de la gestión 2025 garantiza homogeneidad temporal, evita data drift entre períodos, y permite validación estratificada sin sesgos temporales (Carcillo et al., 2018).

Representatividad:

- Cubre 12 meses de operación (gestión 2025 completa)
- Incluye variaciones estacionales del año
- Incluye transacciones legítimas y fraudulentas etiquetadas

Justificación metodológica (según Sampieri, 2014):

Para estudios cuantitativos con poblaciones grandes, un censo o muestra representativa del 70 %+ es adecuada para inferencias válidas. El dataset utilizado (74.60 %) cumple este criterio y además:

- No es sintético, refleja comportamiento real del sistema
- Incluye diversidad de canales (Web, App, POS)
- Incluye múltiples gateways (10+ pasarelas)
- Contiene etiquetado real de fraude (no simulado)

División del Dataset

Estrategia de validación temporal:

Conjunto	Periodo	Transacciones	% Total
Train set	Gestión 2025	10,845,000	70.0 %
Validation set	Gestión 2025	2,324,000	15.0 %
Test set	Gestión 2025	2,324,000	15.0 %
Total	Gestión 2025	15,492,846	100 %

Justificación de validación estratificada:

El dataset de gestión 2025 se divide mediante estratificación por clase (70/15/15) para garantizar representatividad de transacciones fraudulentas en cada conjunto. Este enfoque es más apropiado que validación temporal cuando se trabaja con un solo período homogéneo, según lo recomendado por Dal Pozzolo et al. (2015) y Carcillo et al. (2018) para datasets de fraude con alta desbalance de clases.

Cronograma de Actividades

Planificación Ejecutiva - 12 Semanas

Tabla 2: Cronograma Detallado de Actividades

Semana	Fechas	Actividad	Entregables
0	Nov 11-17	Setup infraestructura AWS	<ul style="list-style-type: none"> - Servidores AWS configurados - Ambiente Python listo - Acceso dataset verificado
1	Nov 18-24	Corrección perfil + Inicio Cap. 1	<ul style="list-style-type: none"> - Perfil corregido aprobado - Revisión 10 papers - Estructura Cap. 1
2	Nov 25-Dic 1	Continuación Cap. 1	<ul style="list-style-type: none"> - 50% Cap. 1 completo - Revisión 15 papers adicionales - Sección algoritmos ML
3	Dic 2-8	Finalización Cap. 1 + EDA	<ul style="list-style-type: none"> - Cap. 1 completo (100%) - Dataset descargado - Análisis descriptivos
4	Dic 9-15	Cap. 2: Diagnóstico + EDA	<ul style="list-style-type: none"> - Patrones de fraude caracterizados - Proceso etiquetado documentado - Visualizaciones EDA
5	Dic 16-22	Preprocesamiento + Feature Engineering	<ul style="list-style-type: none"> - Dataset limpio (25M trans.) - 15+ features creadas - División temporal train/test
6	Dic 23-29	Entrenamiento modelos	<ul style="list-style-type: none"> - Random Forest entrenado - XGBoost y SVM referencia - Validación temporal ejecutada
7	Dic 30-Ene 5	Optimización + Selección modelo	<ul style="list-style-type: none"> - Grid Search completo - Modelo final seleccionado - Feature importance analizada
8	Ene 6-12	Cap. 3: Desarrollo del Modelo	<ul style="list-style-type: none"> - Capítulo 3 completo - Código en GitHub - Pipeline automatizado
9	Ene 13-19	Cap. 4: Evaluación + Comparación	<ul style="list-style-type: none"> - Métricas finales calculadas - Comparación con literatura - Intervalos confianza bootstrap - Tablas y gráficos
10	Ene 20-26	Redacción: Conclusiones + Intro	<ul style="list-style-type: none"> - Conclusiones completas - Introducción final - Abstract inglés - Borrador completo
11	Ene 27-Feb 2	Correcciones + Preparación defensa	<ul style="list-style-type: none"> - Tesis final lista - Presentación PowerPoint - Script de defensa

Tabla 2 – *Continuación*

Semana	Fechas	Actividad	Entregables
12	Feb 3-9	Ensayo defensa + Ajustes finales	- Defensa ensayada - Respuestas preparadas - Versión final impresa

Total de horas estimadas: 400 horas en 12 semanas = 33 horas/semana promedio

Referencias Metodológicas

Fundamentos del Método AQP/CCA

- **Martínez, R. (2020).** *Los 3 secretos detrás de una tesis: Método AQP/CCA para elaboración de tesis de grado y postgrado.* Editorial Académica.
- **Hernández Sampieri, R., Fernández Collado, C., & Baptista Lucio, P. (2014).** *Metodología de la investigación* (6^a ed.). McGraw-Hill.

Benchmarks de Literatura Científica

- **Hafez, A. I., et al. (2025).** Random Forest for Credit Card Fraud Detection. *Journal of Financial Crime*, F1-Score: 85-94 %.
 - **Feng, L., et al. (2024).** XGBoost for E-commerce Payment Fraud Detection. *IEEE Transactions*, F1-Score: 90-94 %.
 - **Hernández Aros, L. G., et al. (2024).** Machine Learning Models for Payment Fraud Detection. *Scopus Indexed*, Precision: 85-92 %.
-

Firma del Investigador

Ing. Ada Condori Callisaya
Maestrante - UAGRM