



# **Detección de patrones multivariados en ofertas fraudulentas mediante análisis de datos mixtos y sistema de scoring interpretable con visualización interactiva**

**Jharold Alonso Mayorga Villena**

**Orientador:**

*Plan de Tesis presentado la Escuela Profesional Ciencia  
de la Computación como paso previo a la elaboración de  
la Tesis Profesional.*

**UNSA - Universidad Nacional de San Agustín de Arequipa  
Junio de 2025**

---

# 1. Motivación y Contexto

La digitalización global del mercado laboral ha generado nuevas vulnerabilidades sistémicas. Estudios recientes documentan que el 29 % de los buscadores de empleo a nivel global [Commission, 2024a] enfrenta ofertas fraudulentas, con pérdidas económicas estimadas en \$367 millones anuales solo en EE.UU. Estos fraudes impactan especialmente a grupos en situación de vulnerabilidad económica: jóvenes (58 % de casos según [del Trabajo, 2024]), personas con baja cualificación (72 % según [García et al., 2023]), y desempleados de larga duración (41 % según [Alliance, 2024]). Desde la perspectiva computacional, existe una brecha crítica en la literatura actual: modelos como los de [Singh, 2024] muestran limitaciones significativas en el procesamiento de texto no estructurado, con caídas de precisión del 15-22 % en documentos laborales extensos. Esta limitación deriva de dos factores fundamentales: La creciente digitalización del mercado laboral ha generado un terreno fértil para prácticas fraudulentas, con consecuencias devastadoras para la sociedad. Según datos recientes de la Federal Trade Commission [Commission, 2024b], el 29 % de los buscadores de empleo a nivel global ha enfrentado ofertas fraudulentas, con pérdidas económicas que superan los \$367 millones anuales solo en Estados Unidos. Este problema afecta desproporcionadamente a los grupos más vulnerables:

Grupo demográfico	Tasa de afectación
Jóvenes (18-35 años)	58 %
Personas con baja cualificación educativa	72 %
Desempleados de larga duración	41 %

Cuadro 1: Distribución de víctimas de fraudes laborales según [García et al., 2023]

Además del impacto económico, estudios como [Smith, 2023] documentan que el 65 % de las víctimas desarrolla desconfianza permanente hacia las plataformas digitales, lo que limita su acceso a oportunidades laborales legítimas.

### 1.1. Brechas Computacionales

El análisis crítico de 12 estudios recientes (2020-2025) revela tres limitaciones fundamentales en los enfoques existentes:

Estudio	Enfoque	Limitación clave
[Singh, 2024]	Modelos XLM-R para texto	Ignora metadatos estructurales
[Chen et al., 2023]	Fusión TF-IDF + Random Forest	Análisis superficial de patrones lingüísticos
[Kumar and Patel, 2024]	BERT con features híbridas	Alta complejidad computacional
[Wang and Li, 2024]	Redes de grafos (GNN)	Difícil implementación en entornos reales

Cuadro 2: Comparativa de enfoques computacionales existentes

Cuadro 3: Análisis crítico de estudios computacionales recientes (2023-2025)

Autores (Año)	Problema Abordado	Propuesta
Singh (2024)	Baja precisión (82 %) en análisis de documentos largos y complejos	XLM-R para clasificación cross-lingüe
Chen et al. (2023)	Limitaciones en integración de features textuales y estructurales	Fusión TF-IDF + metadatos con Random Forest
Rodríguez (2025)	Detección de patrones evolutivos de fraude	Transformer dinámico con actualización en tiempo real
Wang & Li (2024)	Detección de redes de fraude colaborativo	Redes de Grafos (GNN) para relaciones empresa-puesto
Tanaka et al. (2023)	Identificación de engaños sintácticos complejos	Dependency Parsing + CNN
Kumar & Patel (2024)	Sobreajuste en modelos unimodales	BERT + features estructurales híbridas
Al-Sarem (2023)	Alta varianza en modelos individuales	Ensemble de CNN, LSTM y Transformers
Zhang et al. (2024)	Vulnerabilidad a ataques adversarios	Entrenamiento adversarial robusto
Ortiz et al. (2025)	Alta tasa de falsos negativos en clases minoritarias	Aprendizaje costo-sensitivo con SMOTE-Tomek
Nguyen & Williams (2024)	Subutilización de metadatos no-textuales	Fusión multimodal (texto + imágenes)
Kim et al. (2025)	Limitaciones en detección de patrones novedosos	Few-shot learning con prototipos

---

## 2. Justificación

Los enfoques actuales para detectar ofertas fraudulentas presentan limitaciones significativas que justifican este trabajo. La mayoría de modelos procesan texto y metadatos de forma separada, ignorando relaciones clave entre estos elementos. Soluciones avanzadas basadas en inteligencia artificial funcionan como “cajas negras”, haciendo imposible entender cómo llegan a sus conclusiones. Además, muchos sistemas requieren recursos computacionales que los hacen inviables para su implementación real. Este proyecto propone una alternativa que combina análisis integrado de datos, reglas interpretables y visualización accesible, superando estas tres limitaciones fundamentales. El enfoque permitirá no solo detectar fraudes con alta precisión, sino también explicar claramente los patrones identificados.

## 3. Problemática

El crecimiento exponencial de fraudes laborales representa un problema social y técnico. Desde el punto de vista social, afecta desproporcionadamente a grupos vulnerables y genera desconfianza en las plataformas digitales. Técnicamente, los sistemas actuales tienen precisiones inadecuadas para documentos largos y altas tasas de falsos negativos que permiten el paso de fraudes sofisticados. Los conjuntos de datos públicos están extremadamente desbalanceados, con solo un 5 % de muestras fraudulentas, lo que distorsiona el entrenamiento de modelos. Estos desafíos se ven agravados por la sofisticación creciente de los estafadores, que combinan lenguaje persuasivo con características estructurales específicas para evadir detección.

## 4. Objetivos

El presente trabajo tiene como objetivo principal desarrollar una herramienta de detección temprana de ofertas fraudulentas que supere las limitaciones de los enfoques actuales mediante tres contribuciones clave: Primero, un análisis integrado de datos mixtos que combine técnicas de procesamiento de lenguaje natural (NLP) con minería de reglas de asociación. Este enfoque permitirá identificar patrones composicionales como la co-ocurrencia de ciertos n-gramas en la descripción con valores específicos en los metadatos. Segundo, un sistema de scoring interpretable que asigne pesos a cada señal de riesgo basándose en su poder predictivo medido a través de correlaciones y frecuencia relativa en el dataset. A diferencia de los modelos de caja negra predominantes, este sistema generará explicaciones claras para cada alerta, indicando exactamente qué combinación de factores motivó la clasificación. Tercero, una interfaz visual interactiva diseñada para usuarios no técnicos, que permita explorar los patrones detectados y filtrar ofertas según múltiples criterios.

## Referencias

- [1] Federal Trade Commission, *Consumer Sentinel Network Data Book*, 2024.
- [2] Organización Internacional del Trabajo, *Digital Labor Platforms and Emerging Risks*, 2024.
- [3] García, R. et al., "*Socioeconomic Determinants of Online Fraud Victimization*", Social Science Computer Review, 2023.
- [4] Global Anti-Scam Alliance, *The Evolution of Job Scams in Digital Markets*, 2024.
- [5] Smith, J., "*Psychological Impact of Employment Scams*", Journal of Cybersecurity Psychology, 2023.
- [6] Singh, A., *Cross-lingual Fraud Detection using Transformer Models*", IEEE Transactions on Computational Social Systems, 2024.
- [7] Chen, L. et al., "*Feature Fusion Approaches for Job Scam Detection*", Proc. IEEE Big Data, 2023.
- [8] Shakya, S., *Real/Fake Job Posting Prediction*, Kaggle, 2020.
- [9] Rodríguez, P., "*Dynamic Pattern Recognition in Job Fraud*", Neural Networks, 2025.
- [10] Wang, Y. & Li, Q., "*Graph Neural Networks for Employment Scam Detection*", Proc. ACM KDD, 2024.
- [11] Tanaka, H. et al., "*Deep Syntax Analysis for Fraudulent Job Descriptions*", Proc. ACL, 2023.
- [12] Kumar, R. & Patel, S., "*BERT-based Hybrid Model for Scam Detection*", Expert Systems with Applications, 2024.
- [13] Al-Sarem, M., "*Deep Ensemble Learning for Fake Job Detection*", IEEE Access, 2023.
- [14] Zhang, Y. et al., "*Adversarial Training for Robust Fraud Classification*", Proc. NeurIPS, 2024.
- [15] Ortiz, M. et al., "*Cost-Sensitive Learning for Imbalanced Job Postings*", Machine Learning, 2025.
- [16] Kim et al., "*Few-shot Learning for Job Fraud Detection*", Proc. AAAI, 2025.
- [17] Nguyen & Williams, "*Multimodal Fusion for Job Posting Verification*", Proc. CVPR, 2024.
- [18] Gupta & Sharma, "*Attention Mechanisms for Scam Text Classification*", IEEE Trans. NLP, 2023.
- [19] Rossi et al., "*Unsupervised Anomaly Detection in Job Posts*", Data Mining and Knowledge Discovery, 2024.

- [20] Tanaka & Sato, "Transformer Compression for Real-time Fraud Detection", Proc. ICML, 2025.
- [21] Wei et al., "Cross-platform Fraud Pattern Analysis", ACM Trans. Web, 2024.

## Referencias

- [Alliance, 2024] Alliance, G. A.-S. (2024). The evolution of job scams in digital markets.
- [Chen et al., 2023] Chen, L. et al. (2023). Feature fusion approaches for job scam detection. In *Proc. IEEE Big Data*.
- [Commission, 2024a] Commission, F. T. (2024a). Consumer sentinel network data book.
- [Commission, 2024b] Commission, F. T. (2024b). Consumer sentinel network data book.
- [del Trabajo, 2024] del Trabajo, O. I. (2024). Digital labor platforms and emerging risks.
- [García et al., 2023] García, R. et al. (2023). Socioeconomic determinants of online fraud victimization. *Social Science Computer Review*.
- [Kumar and Patel, 2024] Kumar, R. and Patel, S. (2024). Bert-based hybrid model for scam detection. *Expert Systems with Applications*.
- [Singh, 2024] Singh, A. (2024). Cross-lingual fraud detection using transformer models. *IEEE Transactions on Computational Social Systems*.
- [Smith, 2023] Smith, J. (2023). Psychological impact of employment scams. *Journal of Cybersecurity Psychology*.
- [Wang and Li, 2024] Wang, Y. and Li, Q. (2024). Graph neural networks for employment scam detection. In *Proc. ACM KDD*.