

**Mejora del algoritmo K-means para la detección de fraudes financieros en presencia de  
puntos atípicos**

Jorge Gomez Aguilar

Asesor

Andres Felipe Hernandez Giraldo

Universidad Nacional Abierta y a Distancia UNAD  
Escuela de Ciencias Básicas, Tecnología e Ingeniería ECBTI  
Especialización en Ciencias de Datos y Analítica

2025

**Nota de Aceptación**

---

Nombre Director de Trabajo de Grado

---

Jurado

---

Jurado

### **Dedicatoria**

Dedico este nuevo logro académico a mi hija Emily, te amo mucho y espero siempre estés orgullosa de tu padre, porque logré salir adelante en las condiciones más complicadas para que tu no tengas que pasar por lo mismo. También lo dedico a mi esposa por darme el hogar que siempre anhelé, lleno de amor y apoyo incondicional, y por supuesto, lo dedico a mi madre, que me dio todo el amor que necesité y me hizo saber que las segundas oportunidades existen, ya que nací no de tu vientre sino de tu corazón y por ti soy un ciudadano del mundo.

## Resumen

El presente trabajo de grado se enfoca en la mejora del algoritmo K-means para la detección de fraudes financieros, abordando el problema de los puntos atípicos (outliers) que afectan su precisión. Los outliers pueden distorsionar los resultados del K-means, reduciendo su eficacia para identificar comportamientos fraudulentos en los datos financieros. El objetivo general es evaluar las diferentes soluciones que han sido propuestas para mitigar este problema, mejorando la efectividad del algoritmo en contextos reales. Se revisarán las limitaciones del K-means ante outliers y se evaluarán técnicas alternativas, como el preprocesamiento de datos y algoritmos robustos. Finalmente, se propondrán recomendaciones para optimizar su uso en la detección de fraudes, considerando escenarios con alta presencia de puntos atípicos. Este estudio contribuye a la calidad de los sistemas de detección de fraudes financieros, ofreciendo un enfoque mejorado para el análisis de grandes volúmenes de datos en la industria financiera.

***Palabras claves:*** K-means, cluster, análisis no supervisado, finanzas, outliers

### **Abstract**

The present work focuses on the improvement of the K-means algorithm for financial fraud detection, addressing the problem of outliers that affect its accuracy. Outliers can distort the results of K-means, reducing its effectiveness in identifying fraudulent behavior in financial data. The general objective is to evaluate the different solutions that have been proposed to mitigate this problem, improving the effectiveness of the algorithm in real-world contexts. The limitations of K-means in the face of outliers will be reviewed and alternative techniques, such as data preprocessing and robust algorithms, will be evaluated. Finally, recommendations will be proposed to optimize its use in fraud detection, considering scenarios with high presence of outlier points. This study contributes to the quality of financial fraud detection systems, offering an improved approach for the analysis of large volumes of data in the financial industry.

***Keywords:*** K-means, cluster, unsupervised analysis, financials, outliers

## Tabla de Contenido

Introducción .....	9
Planteamiento del Problema .....	11
Justificación .....	13
Objetivos .....	15
Objetivo General .....	15
Objetivos Específicos .....	15
Marco de Referencia .....	16
Fraude Financiero .....	16
Algoritmo K-means .....	17
Outliers o Puntos Atípicos .....	18
Metodología .....	20
Método .....	20
Tipo de Estudio .....	22
Recolección de Datos .....	23
Limitaciones del Algoritmo K-means .....	25
Técnicas y Enfoques Propuestos para Mejorar la Precisión del K-means .....	33
Recomendaciones para Mejorar el Uso del Algoritmo K-means .....	42
Conclusiones .....	46
Limitaciones .....	48
Investigaciones Futuras .....	49
Recomendaciones .....	50
Referencias .....	52

**Lista de Tablas**

<b>Tabla 1</b> <i>Clasificación del Enfoque para Mejorar el Algoritmo K-means</i> .....	37
<b>Tabla 2</b> <i>Comparación Soluciones Propuestas para la Mejora del Algoritmo k-means</i> .....	38

## Lista de Figuras

<b>Figura 1</b> <i>Proceso Metodológico de Revisión</i> .....	21
<b>Figura 2</b> <i>Comparación de Todos los Algoritmos</i> .....	26
<b>Figura 3</b> <i>Comparativa de Seis Métodos de Detección de Outliers</i> .....	30

## Introducción

La detección de fraudes financieros es uno de los mayores desafíos a los que se enfrentan las instituciones financieras en la actualidad. Los fraudes, que pueden manifestarse de manera sutil o compleja, representan una amenaza significativa para la estabilidad económica de las organizaciones y la confianza de los consumidores. En este contexto, las tecnologías de aprendizaje automático han emergido como herramientas clave para identificar patrones anómalos en grandes volúmenes de datos.

Uno de los algoritmos más utilizados en la minería de datos para este propósito es el K-means, un algoritmo de clustering no supervisado que busca agrupar los datos en categorías basadas en sus similitudes. Sin embargo, una de las principales limitaciones de K-means es su sensibilidad a los puntos atípicos o outliers, que pueden distorsionar el cálculo de los centroides y, en consecuencia, afectar la precisión de la detección de fraudes financieros. Este problema es especialmente crítico en el ámbito financiero, donde las transacciones fraudulentas a menudo representan una pequeña fracción de todos los datos y pueden estar disfrazadas como comportamientos normales.

La presente monografía tiene como objetivo evaluar las soluciones propuestas en la literatura para mitigar los efectos negativos de los puntos atípicos en el rendimiento de K-means. A través de una exhaustiva revisión de los estudios existentes, se identificarán y analizarán las diferentes metodologías alternativas y las mejoras que se han sugerido para robustecer este algoritmo y mejorar su aplicabilidad en la detección de fraudes financieros. El análisis se centrará en soluciones que incluyen técnicas de preprocesamiento de datos, algoritmos híbridos, y modificaciones internas de K-means para hacerlo más robusto frente a los outliers.

A través de esta investigación, se espera contribuir al desarrollo de métodos más eficientes para la detección de fraudes, ofreciendo un análisis detallado de las limitaciones del algoritmo K-means y explorando diversas alternativas que podrían mejorar su desempeño en escenarios donde los datos atípicos son una preocupación constante. Este estudio no solo tiene relevancia en el campo de la minería de datos y el aprendizaje automático, sino que también aporta a la mejora de los sistemas de prevención de fraudes en el sector financiero, fortaleciendo la confianza de los usuarios en los sistemas de seguridad financiera.

## Planteamiento del Problema

El problema central de esta monografía radica en las limitaciones del algoritmo K-means para la detección de fraudes financieros cuando existen puntos atípicos (outliers) en los datos. El K-means es ampliamente utilizado en la detección de patrones anómalos debido a su capacidad para agrupar datos y destacar desviaciones. Sin embargo, este algoritmo asume que los datos son homogéneos y se distribuyen de manera uniforme, lo que no siempre es cierto en conjuntos de datos financieros. La presencia de puntos atípicos puede distorsionar los centroides calculados por K-means, afectando la precisión y resultando en una asignación incorrecta de transacciones a los clústeres. Esto genera falsos positivos, donde transacciones legítimas se clasifican como fraudulentas, y falsos negativos, donde actividades ilícitas pasan desapercibidas.

El problema que se espera analizar es: ¿Cómo afectan los puntos atípicos la eficacia del algoritmo K-means en la detección de fraudes financieros y qué soluciones existen para mitigar sus efectos? La relevancia de abordar esta cuestión es fundamental para mejorar la precisión de los sistemas de detección de fraudes en instituciones financieras, minimizando el impacto negativo tanto de las pérdidas por fraudes no detectados como de los costos asociados a la investigación de falsos positivos.

Actualmente, la sensibilidad de K-means a los outliers se debe a su enfoque basado en la minimización de la varianza dentro de los clústeres, lo que hace que los puntos atípicos influyan desproporcionadamente en la posición de los centroides. Dado que las transacciones financieras pueden incluir datos anómalos por diversas razones, como errores humanos, cambios en el comportamiento del cliente o intentos deliberados de fraude, es esencial entender y superar esta limitación para garantizar un análisis confiable. Las consecuencias de no resolver este problema son significativas: la falta de precisión en la detección de fraudes no solo pone en riesgo los

recursos financieros de las instituciones, sino que también puede llevar a sanciones regulatorias y pérdida de confianza de los clientes. Una gran cantidad de falsos positivos incrementa los costos operativos, al requerir análisis adicionales por parte del equipo de cumplimiento, mientras que los falsos negativos permiten que los fraudes pasen desapercibidos, generando pérdidas sustanciales.

En este contexto, la monografía se enfocará en analizar y evaluar soluciones que han sido propuestas para mitigar el efecto de los outliers en K-means, como el uso de algoritmos robustos y técnicas de preprocesamiento. De esta manera, se contribuirá al desarrollo de métodos más precisos y eficientes para la detección de fraudes financieros, fortaleciendo la capacidad de las instituciones para identificar actividades ilícitas y proteger sus activos.

## Justificación

La detección de fraudes financieros es un desafío crítico para las instituciones bancarias y financieras a nivel mundial. El fraude financiero afecta no solo a las entidades involucradas, sino también a los clientes y la economía en general, generando pérdidas significativas y disminuyendo la confianza del público en el sistema financiero (Carmona Mora & Londoño Morales, 2021). Dada la creciente sofisticación de los métodos utilizados por los defraudadores y el incremento de las transacciones digitales, es esencial contar con sistemas de detección que sean precisos y eficientes.

El algoritmo K-means es una de las técnicas más utilizadas para el análisis de datos en la detección de fraudes debido a su capacidad para identificar patrones y agrupar transacciones (Barragán Garnica, 2022). Sin embargo, como lo menciona Khan et al. (2024) una de sus principales limitaciones es su sensibilidad a los puntos atípicos (outliers), que pueden distorsionar el análisis y reducir la precisión del modelo. Abordar esta problemática es crucial para mejorar la efectividad de los sistemas de detección de fraudes, especialmente en contextos financieros donde los datos suelen incluir transacciones atípicas debido a cambios en el comportamiento del cliente o intentos deliberados de fraude.

La relevancia de esta monografía radica en su contribución al fortalecimiento de los mecanismos de detección de fraudes, lo cual es vital para la estabilidad económica y la protección de los recursos financieros tanto a nivel nacional como internacional. Una detección efectiva del fraude permite a las instituciones financieras minimizar las pérdidas económicas y los costos operativos derivados del análisis de falsos positivos, mejorando la eficiencia operativa y garantizando un entorno más seguro para los usuarios.

Desde una perspectiva social, la investigación propuesta tiene el potencial de contribuir a la protección de los consumidores y al fortalecimiento de la confianza en el sistema financiero. La confianza es un elemento clave para el crecimiento económico, y la capacidad de las instituciones para prevenir actividades fraudulentas es esencial para mantener esta confianza. Además, una mejora en la precisión de los sistemas de detección de fraudes puede ayudar a reducir los costos asociados al cumplimiento regulatorio, lo que resulta en un ahorro significativo de recursos que pueden ser utilizados para mejorar otros servicios financieros.

En el ámbito académico, esta monografía contribuirá al campo de la minería de datos y el aprendizaje automático aplicados a la detección de fraudes. La revisión y evaluación de técnicas para mitigar el impacto de los puntos atípicos en el algoritmo K-means proporcionará un conocimiento valioso sobre las limitaciones de los métodos tradicionales y posibles soluciones para superarlas. Esto no solo beneficiará a la industria financiera, sino que también será útil para otros sectores donde la detección de anomalías es fundamental, como en el análisis de seguridad cibernética y la prevención de lavado de dinero.

En conclusión, la necesidad de mejorar los sistemas de detección de fraudes financieros es indiscutible, y esta monografía se presenta como una respuesta a esta demanda. Al analizar y proponer mejoras al algoritmo K-means para su aplicación en escenarios con puntos atípicos, se espera ofrecer un enfoque innovador y robusto que pueda ser implementado en sistemas de prevención de fraudes. De esta manera, la investigación no solo abordará una problemática actual y relevante, sino que también contribuirá al desarrollo del campo de conocimiento, generando un impacto positivo en la sociedad y en la industria financiera.

## **Objetivos**

### **Objetivo General**

Evaluar las soluciones propuestas para mitigar la sensibilidad del algoritmo K-means frente a la presencia de puntos atípicos en conjuntos de datos financieros, especialmente en la detección de fraudes financieros.

### **Objetivos Específicos**

Identificar las limitaciones clave del algoritmo K-means frente a la presencia de outliers en conjuntos de datos financieros, y cómo estas afectan la precisión en la detección de fraudes.

Comparar las técnicas y enfoques propuestos en la literatura para mejorar la robustez del algoritmo K-means ante datos atípicos.

Proponer recomendaciones para la mejora del uso del algoritmo K-means en sistemas de detección de fraudes financieros en escenarios con puntos atípicos.

## Marco de Referencia

La detección temprana y eficaz del fraude financiero se ha convertido en un eje estratégico dentro del sector financiero y tecnológico, dadas las implicaciones económicas, legales y reputacionales que conlleva. En este contexto, el uso de algoritmos de aprendizaje no supervisado, como el K-means, ha cobrado especial relevancia por su capacidad para agrupar datos sin necesidad de etiquetas previas. Sin embargo, su sensibilidad ante la presencia de puntos atípicos o outliers ha sido ampliamente documentada como una limitación crítica, especialmente en entornos donde los datos fraudulentos representan una fracción mínima del total y poseen patrones atípicos complejos.

Este capítulo presenta los fundamentos teóricos que sustentan el desarrollo del presente estudio.

### Fraude Financiero

El fraude financiero es una preocupación constante para las instituciones bancarias y empresas del sector financiero, ya que genera pérdidas económicas significativas y afecta la confianza de los clientes (Carmona Mora & Londoño Morales, 2021). Con el aumento de transacciones digitales, la detección automática de actividades fraudulentas ha cobrado relevancia, impulsando el uso de algoritmos de aprendizaje automático y minería de datos (Heras Calvo, 2023).

Diversos estudios han demostrado la aplicabilidad de las técnicas de clustering en la detección de fraude, particularmente en el ámbito financiero. La metodología K-means ha sido utilizada en numerosos casos reales debido a su simplicidad y capacidad para segmentar comportamientos anómalos (Hidalgo Cornejo & Solano Barragán, 2019). Por ejemplo, Issa & Vasarhelyi (2011) aplicaron el algoritmo K-means para identificar fraudes financieros, mientras

que Liu, R. et al. (2011) lo emplearon en la detección de actividades relacionadas con el lavado de capital, ambos sobre conjuntos de datos reales. Estas investigaciones evidencian el interés creciente en el uso de K-means como herramienta exploratoria para detectar patrones sospechosos en contextos económicos y fiscales.

### **Algoritmo K-means**

El algoritmo K-means es un método de clustering no supervisado ampliamente utilizado en minería de datos y aprendizaje automático, especialmente para la segmentación y análisis exploratorio de grandes volúmenes de datos (Hidalgo Cornejo & Solano Barragán, 2019; Khan et al., 2024). Esta técnica es fundamental para identificar patrones y comportamientos en conjuntos de datos donde no existen etiquetas previas, como ocurre en el análisis de datos financieros o el comportamiento del cliente en diversos sectores (Barragán Garnica, 2022).

El K-means divide un conjunto de datos en un número predefinido de clústeres (o grupos) basándose en la similitud entre las observaciones, asignando cada punto de datos al clúster cuyo centroide esté más cercano. El centroide de cada clúster es un punto representativo que se calcula como el promedio de las observaciones asignadas a ese clúster, y puede interpretarse como una "representación" de los datos dentro de ese grupo. Este proceso de asignación y recálculo de centroides se repite iterativamente hasta que los puntos de datos ya no cambian de clúster o se cumple un criterio de convergencia (Khan et al., 2024).

Una de las principales fortalezas de K-means es su capacidad para manejar grandes volúmenes de datos y su eficiencia computacional, lo que lo hace adecuado para aplicaciones en sectores donde se generan enormes cantidades de datos, como el sector bancario y la gestión del riesgo de crédito (Heras Calvo, 2023). Al utilizar este algoritmo, las instituciones financieras pueden segmentar a los clientes según sus patrones de gasto o identificar comportamientos

anómalos que podrían ser indicativos de fraude financiero o lavado de dinero (Issa & Vasarhelyi, 2011).

En términos de su funcionamiento, el objetivo de K-means es crear conjuntos de datos homogéneos con propiedades comunes dentro de cada clúster, mientras que mantiene una clara separación entre los clústeres formados, de tal forma que los puntos dentro de un mismo clúster sean más similares entre sí que con los puntos de otros clústeres. Esto permite identificar patrones o segmentos dentro de los datos que podrían no ser evidentes en un análisis superficial (Heras Calvo, 2023; Barragán Garnica, 2022).

### **Outliers o Puntos Atípicos**

En el análisis de datos, los outliers o puntos atípicos se definen como aquellas observaciones que presentan características significativamente diferentes respecto a la mayoría de los datos en un conjunto, pudiendo distorsionar el análisis y la interpretación de resultados (Orellana & Cedillo, 2020; Khan et al., 2024). Su detección y manejo constituye un desafío clave, especialmente en áreas sensibles como las finanzas, donde los valores atípicos pueden estar asociados tanto a errores de registro como a comportamientos anómalos vinculados a fraudes (Issa & Vasarhelyi, 2011; Barragán Garnica, 2022).

Desde el punto de vista estadístico, los outliers se consideran “anomalías” porque no siguen el patrón de comportamiento general de los datos. Su presencia puede influir de manera significativa en técnicas como el clustering, afectando el cálculo de centroides y la definición de grupos homogéneos (Khan et al., 2024). En este sentido, Orellana y Cedillo (2020) explican que los datos atípicos pueden surgir por múltiples factores: errores de medición, errores de ingreso de datos, variabilidad natural o eventos inusuales, siendo estos últimos especialmente relevantes en contextos como la detección de fraudes financieros.

La detección de outliers se realiza a través de dos enfoques principales: global y local. Los enfoques globales evalúan cada observación en el contexto de todo el conjunto de datos, asignando puntuaciones que permiten identificar las instancias más alejadas del comportamiento común. Por otro lado, los enfoques locales comparan las observaciones respecto a sus vecindades más cercanas, siendo particularmente efectivos en conjuntos de datos con densidades variables (Orellana & Cedillo, 2020). Ejemplos de métodos globales incluyen el uso de estadísticas como la distancia de Mahalanobis, mientras que enfoques locales destacan algoritmos como el Local Outlier Factor (LOF), que analiza la densidad local de cada punto (Heras Calvo, 2023).

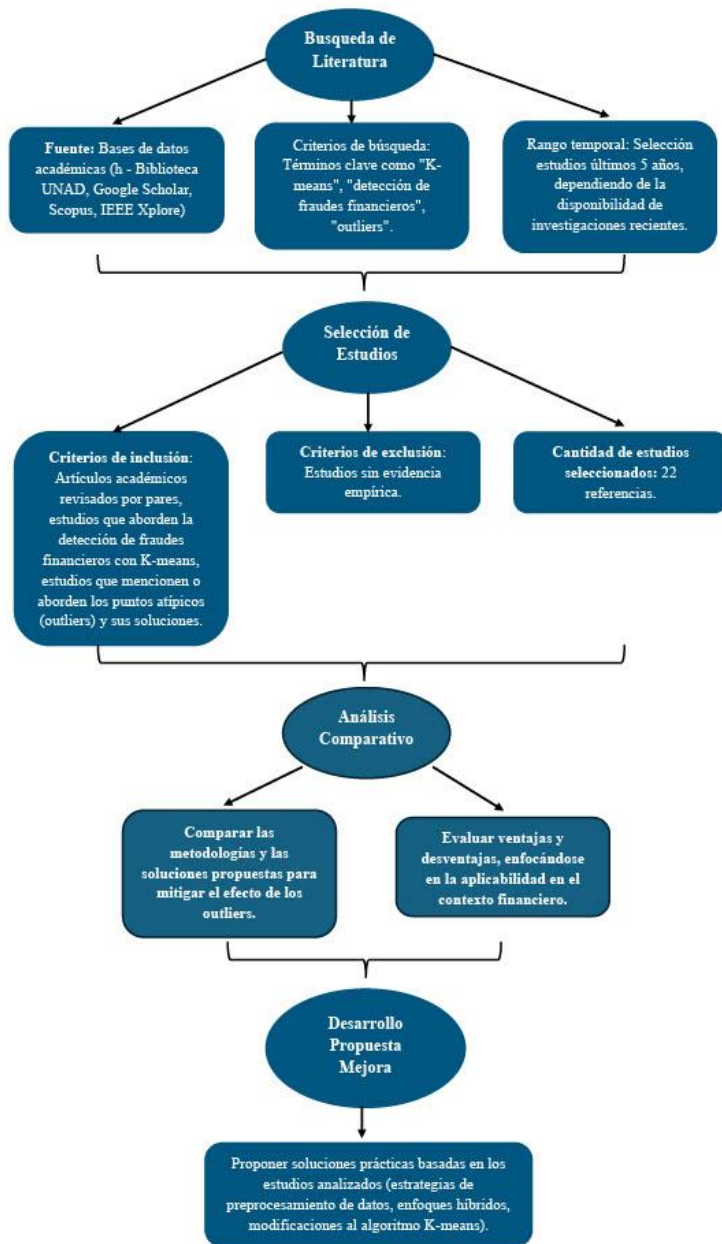
## **Metodología**

### **Método**

El estudio adopta un enfoque cualitativo con un diseño descriptivo y comparativo, centrado en la revisión de la literatura existente sobre la aplicación del algoritmo K-means en la detección de fraudes financieros. Dado el objetivo principal de evaluar cómo los puntos atípicos (outliers) afectan el rendimiento de K-means y las soluciones propuestas para mitigar este problema, se adopta una metodología basada en la revisión bibliográfica sistemática y el análisis comparativo de las metodologías aplicadas. En la figura 1, se describe el proceso metodológico seguido en este proyecto:

Figura 1

## Proceso Metodológico de Revisión



*Nota.* La figura representa la metodología para la revisión sistemática sobre detección de fraudes financieros con K-means y manejo de outliers.

Revisión bibliográfica: Se realizó una búsqueda exhaustiva de estudios académicos, artículos de conferencias y tesis, en total se consultaron 22 referencias de años recientes, que abordan el uso de K-means en la detección de fraudes financieros. La revisión incluyó trabajos que analizan las limitaciones del algoritmo frente a los puntos atípicos, así como soluciones metodológicas para mejorar su robustez. Se recopilaron estudios que presentan enfoques alternativos, como la combinación de K-means con otros algoritmos de detección de anomalías, técnicas de preprocesamiento de datos y la modificación directa del algoritmo para hacerlo más resistente a outliers.

Análisis comparativo: Se llevó a cabo un análisis comparativo de las soluciones propuestas en la literatura para mitigar la influencia de los puntos atípicos sobre K-means. Se identificaron y evaluaron las ventajas y desventajas de cada técnica, con un enfoque particular en la aplicabilidad de estas soluciones al contexto financiero, en donde la presencia de datos atípicos es común debido a la naturaleza impredecible de las transacciones fraudulentas.

Propuestas de mejora: A partir de los hallazgos de la revisión y el análisis comparativo, se desarrollaron recomendaciones prácticas para mejorar la implementación de K-means en la detección de fraudes financieros, específicamente en escenarios con presencia de puntos atípicos. Estas propuestas incluyen estrategias de preprocesamiento de datos, el uso de técnicas híbridas y la aplicación de versiones robustas de K-means que incorporan modificaciones para manejar mejor los outliers.

### **Tipo de Estudio**

Este estudio es descriptivo y exploratorio, con un componente comparativo centrado en la evaluación de soluciones existentes. El análisis tiene como objetivo describir las limitaciones inherentes del algoritmo K-means cuando se enfrenta a la presencia de puntos atípicos y explorar

las metodologías alternativas propuestas para superar dichas limitaciones. Además, se busca comparar la efectividad de las diferentes soluciones en el contexto de la detección de fraudes financieros, analizando cómo mejoran la precisión del modelo.

El estudio sigue un enfoque descriptivo y comparativo, centrado en la evaluación teórica de los enfoques existentes. A través de una revisión exhaustiva de la literatura, se analizan las soluciones propuestas para mejorar el rendimiento de K-means frente a los puntos atípicos, comparando cómo cada solución podría impactar el rendimiento del algoritmo en escenarios financieros reales, basándose en los resultados de estudios previos.

### **Recolección de Datos**

Para este estudio, la recolección de datos se centró en la recopilación de referencias bibliográficas claves de fuentes de datos como h-Biblioteca UNAD, Google Scholar, Scopus, IEEE Xplore, y SpringerLink, estas referencias proporcionan un análisis detallado sobre el uso de K-means en la detección de fraudes financieros. Se incluyó información de fuentes académicas recientes, lo que asegura que los estudios estén actualizados y aborden las metodologías y soluciones más recientes relacionadas con K-means y la detección de fraudes y, que documentan:

**Limitaciones:** Las limitaciones del algoritmo K-means en presencia de puntos atípicos.

**Soluciones:** Las soluciones propuestas para mitigar estos efectos, que incluyen el uso de métodos de preprocesamiento de datos como winsorización, Z-score, y el filtrado por IQR, así como enfoques híbridos que combinan K-means con otros algoritmos de detección de anomalías como DBSCAN, LOF, e Isolation Forest.

**Aplicabilidad:** La aplicabilidad de estas soluciones en el contexto financiero, analizando la efectividad de cada enfoque para mejorar la precisión en la detección de fraudes.

Las fuentes bibliográficas fueron seleccionadas en función de su relevancia, autoría y el impacto de las investigaciones en el campo de la detección de fraudes financieros y el uso de algoritmos de clustering. Se priorizaron las investigaciones que comparan los diferentes métodos de detección de outliers, las que abordan las características específicas del fraude financiero y las que proponen mejoras en la implementación de K-means para hacer frente a los desafíos de los puntos atípicos.

### **Limitaciones del Algoritmo K-means**

Una de las principales limitaciones del algoritmo K-means es su sensibilidad a los puntos atípicos (outliers), ya que estos valores pueden distorsionar el cálculo del centroide y reducir la precisión de la agrupación (Barragán Garnica, 2022), lo cual conlleva a una disminución de la efectividad en la detección de fraudes. Los valores atípicos son puntos de datos que difieren significativamente de la mayoría de las observaciones, es decir, están alejados del comportamiento general de los datos (Orellana & Cedillo, 2020), debido a la variabilidad natural o a errores en la recogida de datos (Khan et al., 2024).

Como lo indica Hidalgo Cornejo & Solano Barragán (2019), existen distintos tipos de algoritmos de clustering y su aplicación varía debido a distintos factores. Los métodos de clustering más utilizados son los jerárquicos y k-means. Sin embargo, con respecto a la sensibilidad al ruido, las técnicas jerárquicas trabajan mejor con bases de datos que contienen valores atípicos en comparación con el k-means.

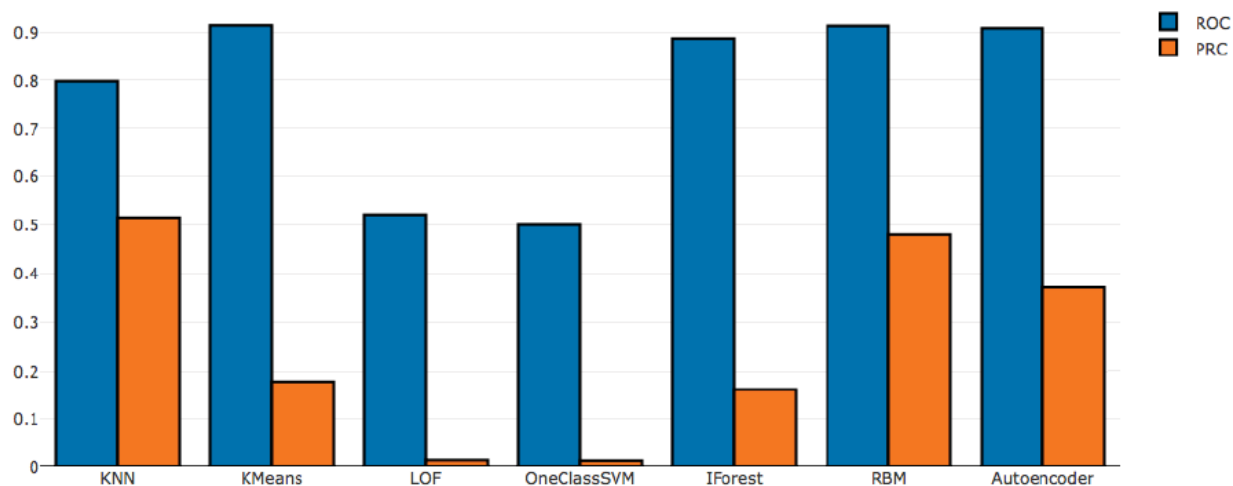
Sahoo & Sahoo (2021) concluye en su trabajo que una limitación importante del algoritmo K-means es su sensibilidad a los valores iniciales de los centroides, lo que puede afectar significativamente la calidad de los resultados, especialmente en presencia de valores atípicos (outliers). Como los valores atípicos pueden distorsionar las medias de los clústeres, el algoritmo puede agrupar incorrectamente los datos, dificultando la identificación precisa de fraudes.

Rodríguez Pérez (2018) demuestra que K-means presenta un rendimiento inferior, especialmente en escenarios donde los datos contienen outliers o están desbalanceados. Las métricas de evaluación (curva ROC y precisión-recall) comparativas evidencian una mayor tasa de falsos positivos y menor efectividad en la detección de fraude en comparación con modelos

como K-nearest neighbors (KNN), Restricted Boltzmann Machine (RBM) y Autoencoder. En la figura 2 se presenta la respectiva comparación de los diferentes modelos.

## Figura 2

### *Comparación de Todos los Algoritmos*



*Nota.* La figura representa la comparación de área bajo las curvas ROC y PR en todos los algoritmos. Tomado de (Rodríguez Pérez, 2018, p. 57).

Además, desde el punto de vista teórico, el trabajo resalta que las técnicas no supervisadas como K-means no son capaces de adaptarse adecuadamente a la superposición de clases ni a los cambios dinámicos en los patrones fraudulentos, lo que agrava su debilidad frente a datos anómalos. Esta limitación es particularmente relevante en el contexto financiero, donde las transacciones fraudulentas tienden a mimetizarse con las legítimas para evitar ser detectadas. De esta manera, el estudio demuestra que si bien K-means puede ser útil en la exploración preliminar de datos, no es adecuado como técnica principal de detección de fraude cuando se presentan valores atípicos o datos complejos.

Ali et al. (2022) presenta una revisión sistemática de la literatura sobre la detección de fraude financiero utilizando técnicas de aprendizaje automático. En este análisis, se evidencia que los métodos de clustering, incluyendo el algoritmo K-means, han sido poco utilizados en comparación con enfoques supervisados, debido a limitaciones inherentes como su sensibilidad a valores atípicos y su bajo desempeño en contextos con datos desbalanceados. Estas características lo hacen poco adecuado para detectar patrones de fraude, que suelen ser irregulares y escasos. El estudio resalta que el uso de K-means puede llevar a una agrupación incorrecta de transacciones, comprometiendo la identificación precisa de actividades fraudulentas.

Bagheri-Gisour Marandyn et al. (2021) muestra la aplicación práctica del algoritmo K-means para la detección de anomalías no supervisadas en un sistema de datos en tiempo real. Sin embargo, se evidenció una caída importante en la precisión cuando las anomalías eran temporales o se encontraban en rangos horarios no frecuentes, lo cual refleja la sensibilidad del algoritmo a valores atípicos.

Huang et al. (2024) evidencia que, aunque K-means permite detectar agrupaciones generales de comportamiento, su rendimiento disminuye considerablemente cuando se enfrenta a fraudes sutiles o distribuidos de manera irregular. Esta debilidad se acentúa debido a su sensibilidad a la inicialización de centroides, la selección del número de clústeres y su incapacidad para modelar adecuadamente los outliers. La estrategia de detección de anomalías mediante la distancia al centroide resulta limitada, ya que puede excluir comportamientos fraudulentos encubiertos o incluir transacciones legítimas como sospechosas. En conjunto, estos hallazgos reafirman que, si bien K-means es una herramienta eficiente para análisis preliminar,

no es suficiente como técnica principal en contextos financieros donde los datos suelen ser desbalanceados y los fraudes se manifiestan de forma atípica.

Nguyen et al. (2020) presentan un análisis comparativo sobre el rendimiento de distintos métodos de aprendizaje profundo frente a técnicas tradicionales en la detección de fraudes con tarjetas de crédito. El trabajo evidencia que los modelos tradicionales como K-means tienden a presentar bajo rendimiento en escenarios con datos desbalanceados y atípicos, lo cual es característico de los fraudes financieros. En particular, destacan que estos algoritmos son sensibles a los cambios en las tácticas de fraude y generan altas tasas de falsos positivos y falsos negativos, debido a su limitada capacidad para capturar patrones complejos o evolutivos en los datos.

Ramírez Mendoza (2022) revela que el algoritmo K-means genera clústeres de tamaño desigual, muchos de ellos con un número reducido o incluso unitario de observaciones, lo que evidencia su sensibilidad a datos extremos. Asimismo, la necesidad de aplicar previamente una reducción de dimensionalidad mediante PCA para lograr agrupaciones más consistentes sugiere que K-means carece de robustez ante datos financieros crudos, que frecuentemente incluyen outliers. En comparación con el modelo aglomerativo, K-means mostró mayor dispersión e inestabilidad en los resultados, reafirmando que, en contextos donde el fraude financiero se manifiesta como anomalías o desviaciones estadísticas, este algoritmo puede resultar poco confiable sin un tratamiento previo adecuado.

Choez Franco (2022) realiza un análisis comparativo de algoritmos de clustering en aprendizaje no supervisado y subraya que K-means no está diseñado para manejar valores atípicos de forma adecuada. Choez destaca que este algoritmo asigna todas las observaciones a algún clúster, sin distinguir si un punto es significativamente diferente del resto, lo cual

representa una limitación crítica en tareas donde los outliers pueden ser precisamente los elementos de interés, como en la detección de fraude. Además, menciona que K-means requiere la definición previa del número de clústeres, dificultando su aplicación en escenarios donde los fraudes emergen de forma no estructurada.

Ponce Del Águila et al. (2019) se enfoca en las buenas prácticas para la gestión del riesgo de fraude interno en instituciones financieras peruanas. Aunque no implementa modelos de clustering como K-means, ofrece un marco conceptual fundamental que resalta la naturaleza atípica, poco frecuente y heterogénea de los fraudes financieros. Esta caracterización pone de manifiesto una limitación inherente del K-means: su tendencia a formar clústeres homogéneos basados en promedios y distancias euclidianas, lo cual compromete su capacidad para identificar conductas fraudulentas aisladas o no recurrentes. Ponce enfatiza que los fraudes internos suelen escapar a patrones fijos, lo que exige técnicas que reconozcan irregularidades estadísticas y no dependan únicamente de estructuras fijas de agrupamiento, como es el caso del K-means tradicional.

Boucher (2020) realiza una evaluación comparativa de seis métodos de detección de outliers aplicados al fraude financiero, incluyendo técnicas de clustering como K-means. En su revisión, clasifica K-means como un método "naive", subrayando que este tipo de algoritmos tiende a producir clústeres demasiado amplios que incorporan tanto datos normales como atípicos, debido a su dependencia de la media aritmética para definir centroides. Esto se traduce en una ineficacia para separar adecuadamente las transacciones fraudulentas del resto del conjunto, especialmente cuando los fraudes no representan patrones claramente diferenciables. Además, Boucher muestra que algoritmos como Logistic Regression o Multi-layer Perceptron obtienen mejores resultados en precisión y F1-score, evidenciando que K-means es menos

adecuado para este tipo de tareas si no se acompaña de preprocesamiento específico para detectar outliers. Los resultados obtenidos se pueden observar en la figura 3.

### Figura 3

#### *Comparativa de Seis Métodos de Detección de Outliers*

	Financial set F-1 score (Class 1)	Credit card set F-1 score (Class 1)	F-1 score (Class 0)	Audit set F-1 score (Class 1)	Accuracy
HDBSCAN		0.00	0.81	0.41	0.71
Statistical method		0.53	0.92	0.86	0.89
Isolation forest	0.10	0.33	0.79	0.76	0.78
Local outlier factor	0.10	0.11	0.89	0.84	0.87
Logistic regression	<b>0.38</b>	<b>0.70</b>	<b>0.96</b>	<b>0.92</b>	<b>0.95</b>
Multi layer perceptron		0.73	1.00	0.99	0.99
Average running time	14min	11min		13s	

*Nota.* Evaluación comparativa de seis métodos de detección de outliers aplicados al fraude financiero. Tomado de (Boucher, É. 2020, p. 43).

Liu, Y. et al. (2022) proponen un modelo mejorado de K-means para contextos de detección de fraude, incorporando técnicas de preservación de privacidad y manejo de datos dispersos. El artículo reconoce explícitamente que el K-means tradicional presenta serias limitaciones cuando se enfrenta a conjuntos de datos escasos o altamente dispersos, como ocurre con perfiles financieros incompletos o normalizados mediante técnicas como one-hot encoding. Además, se identifica que la mayoría de los modelos previos en K-means con preservación de privacidad revelan información sensible durante el proceso, como centroides intermedios o estructuras de los clústeres. Este trabajo evidencia que el modelo clásico no es apto para aplicaciones reales sin ajustes criptográficos y estructurales significativos, particularmente en entornos financieros donde la precisión y confidencialidad son críticas.

Hodge & Austin (2004) presentan una revisión exhaustiva de metodologías para la detección de outliers, incluyendo enfoques como clustering y detección basada en novedad. Los autores identifican que algoritmos como K-means, que buscan minimizar la varianza interna del clúster, tienden a absorber valores atípicos dentro de agrupamientos normales, especialmente si estos no alteran de manera significativa el cálculo de los centroides. Se advierte que este comportamiento distorsiona el modelo de normalidad y puede ocultar casos críticos, como fraudes financieros. Hodge y Austin concluyen que en contextos donde los outliers representan el fenómeno de interés (fraude, sabotaje, intrusión), deben emplearse métodos robustos que permitan diagnosticar y excluir puntos anómalos del modelo central, algo que K-means no ofrece de forma nativa.

En conjunto, la revisión de la literatura evidencia de manera concluyente que el algoritmo K-means presenta limitaciones críticas cuando se enfrenta a valores atípicos en el contexto de la detección de fraudes financieros. Estas limitaciones se manifiestan principalmente en su sensibilidad a la presencia de outliers, su tendencia a formar clústeres homogéneos basados en promedios, su dependencia de la inicialización de centroides y su baja adaptabilidad frente a datos desbalanceados o dispersos. Como consecuencia, K-means tiende a agrupar de manera incorrecta tanto observaciones legítimas como fraudulentas, comprometiendo la precisión del modelo y elevando la tasa de falsos positivos y negativos.

Además, varios estudios han demostrado que la eficacia del algoritmo disminuye considerablemente cuando los fraudes presentan comportamientos sutiles, irregulares o camuflados entre datos normales, lo cual es característico del entorno financiero moderno. En este sentido, aunque K-means puede ser útil en fases exploratorias o como complemento dentro de arquitecturas más robustas, no debe ser considerado como la herramienta principal para la

detección automatizada de fraude si no se le aplican mejoras sustanciales o se combina con técnicas de control de outliers.

A partir de este análisis, se hace necesario examinar las estrategias y soluciones que han sido propuestas o desarrolladas con el fin de superar estas limitaciones. Por ello, en la siguiente sección se procederá a evaluar las técnicas y soluciones implementadas para mejorar la precisión del K-means ante la presencia de outliers en los datos financieros, abordando tanto propuestas metodológicas como adaptaciones algorítmicas desarrolladas en contextos reales y experimentales.

### **Técnicas y Enfoques Propuestos para Mejorar la Precisión del K-means**

Ante las limitaciones ampliamente documentadas del algoritmo K-means frente a valores atípicos tales como: su sensibilidad a centroides iniciales, su dependencia de la forma esférica de los clústeres y su vulnerabilidad ante datos dispersos o desbalanceados, diversos estudios han propuesto mejoras metodológicas, técnicas híbridas o adaptaciones específicas para incrementar su eficacia en la detección de fraudes financieros.

Liu, Y. et al. (2022) propone una versión mejorada del algoritmo K-means diseñada para entornos financieros que exigen no solo eficiencia, sino también preservación de la privacidad y manejo de datos dispersos. La propuesta integra técnicas de criptografía homomórfica y compartición de secretos (secret sharing) para proteger los centroides y las asignaciones de clúster sin revelar información sensible entre instituciones. Además, el algoritmo está optimizado para trabajar con matrices dispersas, lo cual es crucial cuando se analizan grandes volúmenes de datos financieros, donde la presencia de transacciones atípicas o incompletas es común. Esta mejora permite que el algoritmo mantenga su funcionalidad en sistemas de detección de fraude interinstitucional, reduciendo al mismo tiempo la distorsión causada por valores extremos al utilizar mecanismos de tolerancia estructurada frente a outliers.

Vargas et al. (2018) aplica el algoritmo Robust Sparse K-means Clustering (RSKC) para agrupar bancos en Argentina, específicamente con el objetivo de reducir el impacto de valores atípicos. Este enfoque mejora el algoritmo tradicional de K-means al introducir un parámetro de "trimming", que permite excluir del análisis un porcentaje determinado de datos considerados extremos. Esta exclusión ayuda a proteger los centroides del sesgo provocado por outliers, lo que se traduce en clústeres más coherentes y representativos. El enfoque también incluye un componente de penalización para las variables menos relevantes, facilitando la reducción de

ruido en los datos. Los resultados obtenidos por Vargas muestran una mayor estabilidad en los grupos formados y una menor varianza intra-clúster, lo que representa una ventaja significativa en el análisis financiero, donde los datos tienden a presentar alta dispersión.

Aunque el estudio de Nguyen et al. (2020) no implementa directamente el algoritmo K-means, sus observaciones sobre el desbalance de clases, la importancia del preprocesamiento y la presencia de puntos atípicos pueden ser extrapoladas a los desafíos que enfrenta K-means en contextos similares. Dado que K-means tiende a asignar todos los puntos a un clúster, incluso a los valores extremos, resulta ineficiente para detectar fraudes que justamente se comportan como anomalías sin un tratamiento previo adecuado. En este sentido, prácticas como la estandarización de variables o la reducción de dimensionalidad (PCA) podrían mejorar su capacidad para distinguir transacciones fraudulentas en contextos reales.

Huang et al. (2024) utiliza análisis de componentes principales (PCA) para reducir la dimensionalidad del conjunto de datos, conservando las características más relevantes en un espacio más compacto. Esta transformación facilita la aplicación del algoritmo K-means, que se emplea posteriormente para identificar patrones anómalos en un contexto de fraude financiero. El estudio se realiza sobre un conjunto de datos altamente desbalanceado, donde solo el 1.2% de las transacciones corresponden a fraudes. La detección se basa en la distancia de cada punto al centroide de su clúster, clasificando como sospechosas aquellas transacciones significativamente alejadas. El trabajo demuestra que esta combinación de reducción de dimensionalidad y agrupamiento no supervisado permite aislar con efectividad las transacciones fraudulentas y mejorar la precisión del modelo en escenarios reales.

Bagheri-Gisour Marandyn et al. (2021) implementa K-means para detectar anomalías en sistemas de flujo de datos en tiempo real, y destaca que el algoritmo requiere una arquitectura de

preprocesamiento robusta para ofrecer resultados útiles. La solución propuesta incluye varios pasos: identificación y eliminación de datos nulos, normalización, ensamblado de vectores y definición de umbrales de distancia para clasificar valores anómalos. Además, se implementa una validación basada en distancia al centroide, permitiendo marcar registros sospechosos según su proximidad relativa. Aunque esta estrategia no modifica el núcleo del algoritmo, sí configura un entorno en el que K-means puede funcionar de forma más estable, especialmente cuando las anomalías son temporales o de baja frecuencia, como sucede frecuentemente en fraudes financieros.

Ali et al. (2022) En su revisión sistemática de modelos de aprendizaje automático aplicados a la detección de fraude financiero, Ali subraya que los pocos casos exitosos que usan K-means combinan este algoritmo con otros métodos, como PCA, Isolation Forest o LOF. Se señala que K-means rara vez se emplea como técnica principal, y más bien se integra en sistemas híbridos para la reducción de dimensionalidad, agrupamiento exploratorio o segmentación previa. Esta estructura modular permite que K-means actúe en una primera fase de análisis para luego ser complementado por técnicas especializadas en detección de outliers. El artículo resalta que este enfoque puede reducir la tasa de falsos positivos y mejorar la precisión general, especialmente cuando los fraudes presentan patrones no lineales o distribuciones irregulares.

Barragán Garnica (2022) remarca que el algoritmo K-means es altamente sensible a valores extremos, y por lo tanto, debe aplicarse únicamente después de un proceso de estandarización y filtrado estadístico. En su estudio, se utilizó el eliminado de outliers y normalización Z-score antes de implementar K-means, lo cual permitió una segmentación más coherente y útil desde el punto de vista comercial. Aunque el enfoque se centra en marketing, las

implicaciones son trasladables al ámbito financiero, pues refuerza que, sin un tratamiento previo adecuado, K-means puede generar agrupaciones inexactas o engañosas.

Khan et al. (2024) propone la utilización de la winsorización, una técnica estadística que limita el impacto de los outliers reemplazando los valores extremos por un valor umbral definido. Esta técnica conserva el volumen de datos sin eliminar observaciones completas, lo que permite mantener la estructura del conjunto sin distorsionar el cálculo de centroides. Aplicada antes del K-means, la winsorización logra reducir la varianza artificial inducida por valores atípicos extremos, contribuyendo así a una mejor formación de clústeres y a una mayor precisión en la detección de patrones de fraude, especialmente en entornos con alta variabilidad financiera.

Sahoo & Sahoo (2021) proponen el uso del algoritmo K-means como una herramienta preliminar para la detección de fraude contable, destacando su simplicidad, bajo costo computacional y facilidad de implementación. Aunque reconocen que existen técnicas más complejas, el estudio no incorpora modelos adicionales, sino que se enfoca en demostrar que K-means, por sí solo, puede ser eficaz para identificar patrones anómalos en los estados financieros. Como perspectiva futura, los autores sugieren hibridar K-means con variantes como fuzzy c-means o GMM para superar su sensibilidad a los valores iniciales de los centroides.

Torres (2023) implementa K-means para segmentar clientes en el sector financiero y aborda directamente la necesidad de tratamiento previo de outliers. Se utiliza el rango intercuartílico (IQR) para identificar y eliminar valores atípicos antes de aplicar el algoritmo, lo cual mejora la cohesión de los clústeres formados. Además, se evalúan métricas como el índice de Dunn y el coeficiente de silueta para validar la calidad de la agrupación. La estrategia evidencia que la precisión del K-means puede mejorar sustancialmente si los outliers se manejan con criterios estadísticos sólidos antes de la ejecución del algoritmo.

Con base en la revisión detallada de estudios recientes, es posible clasificar las propuestas dirigidas a mejorar el rendimiento del algoritmo K-means frente a la presencia de valores atípicos en tres enfoques principales: 1. estrategias de preprocesamiento de datos, que buscan acondicionar la información antes de aplicar el algoritmo; 2. técnicas híbridas, que combinan K-means con otros algoritmos para compensar sus debilidades; y 3. método alternativo, que modifican la estructura o el funcionamiento interno del algoritmo clásico o proponen uno alternativo. Esta clasificación no solo permite agrupar las soluciones según su naturaleza técnica, sino que facilita la comparación de sus fundamentos metodológicos y su aplicabilidad en entornos financieros. La Tabla 1. resume los aportes más relevantes encontrados en la literatura especializada, destacando el enfoque de mejora utilizado y su descripción:

**Tabla 1**

*Clasificación del Enfoque para Mejorar el Algoritmo K-means*

Estudio	Enfoque de mejora	Descripción
Liu, Y. et al. (2022)	Método alternativo	Se propone una nueva versión del algoritmo K-means que incorpora criptografía y tolerancia a dispersión, lo que modifica la estructura interna del algoritmo.
Vargas et al. (2018)	Método alternativo	El algoritmo RSKC es una modificación directa de K-means con parámetros de trimming y penalización, lo que representa una técnica nueva sobre el algoritmo base.
Nguyen et al. (2020)	Preprocesamiento	Aplicar técnicas de normalización y reducción de dimensionalidad antes de K-means, constituye un preprocesamiento para mejorar la entrada del modelo.
Huang et al. (2024)	Preprocesamiento	Aplicar PCA para reducir la dimensionalidad y después aplicar K-means, haciendo que mejore su precisión en la detección de fraudes financieros.

Estudio	Enfoque de mejora	Descripción
Sahoo & Sahoo (2021)	Técnica híbrida	K-means se emplea como paso preliminar antes de aplicar modelos más robustos, combinando técnicas de diferente complejidad.
Bagheri-Gisour Marandyn et al. (2021)	Preprocesamiento	La propuesta se enfoca en la preparación de datos mediante limpieza, normalización y validación previa, facilitando la operación posterior del K-means.
Ali et al. (2022)	Técnica híbrida	K-means se combina con PCA, LOF e Isolation Forest en una arquitectura mixta, representando una técnica híbrida para compensar debilidades.
Barragán Garnica (2022)	Preprocesamiento	Se realiza filtrado estadístico (Z-score) para eliminar outliers antes de aplicar K-means, mejorando el conjunto de datos de entrada.
Khan et al. (2024)	Preprocesamiento	Se utiliza la winsorización para suavizar valores extremos antes de agrupar, lo que actúa como preprocesamiento estadístico.
Torres (2023)	Preprocesamiento	Se filtran valores atípicos mediante IQR antes de la aplicación de K-means, lo que mejora la calidad de los clústeres mediante preprocesamiento.

*Nota.* Se hace una clasificación de las soluciones presentadas en la literatura y se agrupan de acuerdo con su enfoque para mejorar el algoritmo K-means.

Adicional, la tabla 2 compara las modificaciones realizadas en el algoritmo K-means, destacando las ventajas y limitaciones de cada técnica, así como su contexto de aplicación. La finalidad es proporcionar una visión clara y estructurada de las soluciones más relevantes para superar las limitaciones de K-means en la detección de fraudes financieros.

## **Tabla 2**

*Comparación Soluciones Propuestas para la Mejora del Algoritmo k-means*

Técnica	Ventajas de la Modificación	Limitaciones	Contexto de Aplicación
Winsorización	Reduce la influencia de outliers sin eliminar los puntos atípicos, ayudando a evitar distorsiones significativas en los centroides de K-means.	No elimina los outliers, solo los suaviza. Requiere un umbral bien definido para evitar la pérdida de información importante.	Detección de fraude en entornos con alta variabilidad financiera y presencia de valores extremos.
Z-score	Normaliza los datos antes de aplicar K-means, mejorando la precisión al reducir la influencia de variables con varianzas desproporcionadas.	No aplicable si los datos no son normales (distribución no gaussiana).	Detección de fraudes financieros en clientes o transacciones con características variables.
DBSCAN	K-means se utiliza para la segmentación, mientras que DBSCAN filtra los outliers, mejorando la precisión de los clústeres y reduciendo la distorsión de los centroides.	Requiere un preprocesamiento adecuado. La combinación de ambos algoritmos puede ser costosa computacionalmente.	Detección de fraudes financieros con valores atípicos que son críticos para la identificación.
RSKC (Robust Sparse K-means)	Introduce un proceso de "trimming" para excluir los outliers, lo que ayuda a estabilizar los centroides y mejora la precisión en los clústeres.	La implementación de trimming puede aumentar la complejidad y requerir ajustes cuidadosos.	Agrupación de bancos en sectores financieros donde los datos presentan alta dispersión.

Técnica	Ventajas de la Modificación	Limitaciones	Contexto de Aplicación
PCA (Análisis de Componentes Principales)	Mejora la calidad del input, permitiendo que K-means se enfoque en las variables más relevantes, reduciendo el ruido y las correlaciones innecesarias entre las variables.	Puede perder información importante si la reducción de dimensiones no se realiza correctamente.	Detección de fraude con tarjetas de crédito, donde las transacciones tienen muchas variables.
Normalización y Limpieza de Datos	Prepara los datos para una mejor segmentación, asegurando que las distorsiones generadas por datos nulos o ruidosos no afecten los resultados.	Depende de un preprocesamiento robusto. La calidad de los resultados depende de la calidad del filtrado de los datos.	Detección de anomalías en tiempo real, especialmente cuando las transacciones tienen datos nulos o erróneos.
PCA + Isolation Forest	Mejora la precisión al reducir la dimensionalidad primero, y luego usa Isolation Forest para filtrar los outliers antes de que K-means realice la segmentación.	La combinación de estos métodos aumenta la complejidad computacional y puede ser más costosa en grandes volúmenes de datos.	Detección de fraude en datos financieros donde las transacciones fraudulentas son escasas o distribuidas irregularmente.
IQR (Rango Intercuartílico)	Filtra los puntos atípicos antes de la segmentación, mejorando la cohesión de los clústeres y reduciendo	Solo funciona si los outliers son claramente identificables y la	Segmentación de clientes en instituciones financieras,

Técnica	Ventajas de la Modificación	Limitaciones	Contexto de Aplicación
	la distorsión en los centroides.	estrategia de filtrado es efectiva.	especialmente en el análisis de fraude en productos de inversión.

*Nota.* Se realiza una comparación de las diferentes soluciones propuestas para la mejora del algoritmo k-means y de detalla las principales desventajas de su aplicación.

La clasificación de los enfoques de mejora en preprocesamiento, técnicas híbridas y métodos alternativos permite estructurar de forma clara y sistemática las múltiples estrategias desarrolladas para aumentar la precisión del algoritmo K-means frente a valores atípicos. Esta categorización facilita no solo la comprensión de las soluciones existentes, sino también la elección informada de la técnica más adecuada según el contexto de aplicación. La comparación detallada en la tabla 2 revela que, aunque cada propuesta tiene fortalezas específicas, todas coinciden en la necesidad de adaptar el entorno de aplicación del K-means para mitigar su sensibilidad a datos extremos. Desde estrategias ligeras como la winsorización o el IQR, hasta arquitecturas más complejas que integran PCA e Isolation Forest, los estudios analizados evidencian que la eficacia del algoritmo depende en gran medida del tratamiento previo de los datos y de su integración con otros métodos. En consecuencia, el K-means sigue siendo una herramienta útil en la detección de fraudes financieros, siempre que se emplee dentro de un marco metodológico robusto y contextualizado.

### **Recomendaciones para Mejorar el Uso del Algoritmo K-means**

La revisión de los estudios relacionados con la mejora del algoritmo K-means frente a los outliers revela que, aunque K-means es un algoritmo poderoso y eficiente en la segmentación de datos, presenta limitaciones críticas cuando se enfrenta a puntos atípicos o outliers. Estas limitaciones incluyen sensibilidad a los valores extremos, dependencia de los centroides iniciales, y una baja capacidad de adaptación frente a datos desbalanceados o dispersos (como los que se encuentran en fraudes financieros). Estas deficiencias son especialmente problemáticas cuando se trata de detectar fraudes financieros, donde los comportamientos atípicos pueden estar camuflados entre grandes volúmenes de datos regulares y, si no se manejan correctamente, pueden distorsionar los resultados del análisis.

Dejar el algoritmo K-means intacto no es adecuado para la detección de fraudes financieros en presencia de puntos atípicos. Sin modificaciones, el algoritmo no puede diferenciar adecuadamente entre comportamientos fraudulentos y legítimos cuando los outliers afectan la formación de los clústeres. La sensibilidad a los outliers puede generar falsos positivos y falsos negativos, lo que podría perjudicar la confiabilidad de un sistema de detección de fraude. Además, el algoritmo K-means clásico tiende a suponer que los clústeres tienen una forma esférica y tamaño homogéneo, lo que no siempre se ajusta a la distribución real de los datos financieros, donde las transacciones fraudulentas pueden estar dispersas y no siguen patrones esféricos.

Por lo tanto, modificar K-means es esencial para mejorar la precisión y la robustez del modelo. Las soluciones que incluyen preprocesamiento de los datos, técnicas híbridas o modificaciones internas del algoritmo están diseñadas para mitigar los efectos negativos de los outliers y mejorar la segmentación de los datos. Modificar el algoritmo o complementarlo con

otras técnicas es crucial para mejorar la adaptabilidad de K-means en el contexto de la detección de fraudes financieros. En este contexto, la solución más adecuada dependerá de las características específicas del conjunto de datos y el entorno financiero en el que se aplique. Sin embargo, en términos generales, la combinación de K-means con DBSCAN parece ser la más prometedora, ya que permite que el sistema filtre outliers de manera eficiente antes de aplicar K-means, lo que optimiza tanto la precisión como la velocidad del proceso. Esta solución híbrida puede manejar los outliers de manera más robusta sin incrementar significativamente la complejidad computacional, lo que la hace ideal para entornos financieros reales, donde la velocidad de procesamiento es crucial.

Sin embargo, RSKC también representa una opción sólida, especialmente si se prioriza la estabilidad de los clústeres y se pueden ajustar los parámetros de trimming de manera adecuada para reducir la influencia de los outliers. Este enfoque puede ser útil en escenarios donde se desea mantener la simplicidad del algoritmo sin recurrir a combinaciones complejas de múltiples técnicas.

A partir del análisis de las limitaciones identificadas del algoritmo K-means y la evaluación de diversas técnicas y soluciones implementadas para mitigar su sensibilidad frente a valores atípicos, es posible establecer un conjunto de recomendaciones orientadas a mejorar su aplicación en la detección de fraudes financieros. Estas recomendaciones integran tanto aspectos metodológicos como ajustes técnicos que se desprenden de la literatura revisada y de los casos de estudio analizados.

Aplicar técnicas de preprocesamiento para la gestión de outliers

Los estudios de Barragán Garnica (2022), Khan et al. (2024) y Torres (2023) coinciden en señalar que el rendimiento de K-means mejora significativamente cuando los datos han sido

previamente tratados mediante técnicas de limpieza estadística. Se recomienda, por tanto, implementar procesos de estandarización (Z-score), winsorización o filtrado por rango intercuartílico (IQR) antes de ejecutar el algoritmo, con el fin de reducir la distorsión provocada por valores extremos en el cálculo de centroides.

#### Combinar K-means con métodos de detección de anomalías

Autores como Choez Franco (2022), Ali et al. (2022) y Sahoo & Sahoo (2021) proponen enfoques híbridos en los que K-means se complementa con algoritmos especializados en detección de outliers como DBSCAN, LOF, Isolation Forest o SVM. Estas estrategias permiten a K-means cumplir un rol exploratorio, mientras que los métodos adicionales refinan los resultados e incrementan la sensibilidad del sistema frente a fraudes no evidentes.

#### Implementar variantes robustas del algoritmo

En contextos donde se requiere mantener el uso de K-means como núcleo del proceso analítico, se recomienda optar por versiones adaptadas del algoritmo, como el Robust Sparse K-means Clustering (RSKC) descrito por Vargas et al. (2018), o técnicas de centroides ponderados y trimming. Estas variantes permiten minimizar el impacto de los outliers mediante la exclusión parcial de datos extremos o la penalización de variables irrelevantes.

#### Optimizar la entrada de datos mediante reducción de dimensionalidad

Nguyen et al. (2020) y Bagheri-Gisour Marandyn et al. (2021) destacan la importancia de preparar adecuadamente los datos antes de ejecutar K-means, especialmente cuando se trata de conjuntos financieros con muchas variables. El uso de técnicas como PCA o ensamblaje vectorial ayuda a disminuir el ruido estadístico y a concentrar la información relevante, favoreciendo la correcta identificación de clústeres y la detección de patrones de fraude.

#### Redefinir criterios de agrupamiento con umbrales dinámicos

Huang et al. (2024) propone redefinir la lógica de agrupación mediante el uso de distancias al centroide con umbrales dinámicos. Esta técnica permite clasificar una transacción como sospechosa si su distancia al centroide excede un valor adaptativo, mejorando la sensibilidad del sistema sin depender exclusivamente de la forma del clúster. Se sugiere incorporar métricas como el índice de Dunn y el coeficiente de silueta para evaluar la cohesión y separación de los clústeres generados.

#### Integrar el algoritmo en arquitecturas multinivel

Sahoo & Sahoo (2021) propone una arquitectura de análisis escalonada, en la cual K-means se utiliza como filtro inicial que clasifica grandes volúmenes de datos, y se complementa posteriormente con algoritmos más robustos o supervisados que validan o profundizan el análisis. Este enfoque permite conservar la eficiencia computacional del algoritmo sin comprometer la precisión de la detección.

#### Asegurar la privacidad y adaptabilidad del modelo en entornos dispersos

Finalmente, Liu, Y. et al. (2022) introduce una mejora clave para escenarios interinstitucionales: la implementación de K-means con preservación de privacidad y tolerancia a matrices dispersas. Para instituciones financieras que comparten datos de clientes o transacciones, se recomienda explorar modelos que utilicen criptografía homomórfica y esquemas de compartición de secretos, asegurando la confidencialidad de los datos y la solidez del modelo frente a inconsistencias estructurales.

## Conclusiones

A lo largo de esta monografía se ha analizado el algoritmo K-means y su aplicabilidad en la detección de fraudes financieros, con un enfoque particular en cómo los puntos atípicos (outliers) afectan su rendimiento. El estudio ha permitido identificar tanto las limitaciones inherentes del algoritmo como las soluciones propuestas en la literatura para mitigar sus efectos.

A continuación, se presentan las conclusiones derivadas del desarrollo de los objetivos:

Con respecto al primer objetivo, el cual aborda las limitaciones del algoritmo K-means, se muestra que el algoritmo tiene varias limitaciones inherentes, tales como su sensibilidad a la inicialización de los centroides, su dependencia de la forma esférica de los clústeres y su baja capacidad para manejar datos dispersos o desbalanceados. Estas limitaciones afectan directamente la eficacia del algoritmo en entornos financieros donde los datos pueden ser irregulares o atípicos.

Abordando el segundo objetivo, correspondiente a las soluciones y enfoques alternativos, a través de la revisión de la literatura, se propusieron varias soluciones para mejorar la robustez de K-means frente a los outliers. Estas incluyen el uso de técnicas de preprocesamiento, la combinación con métodos de detección de anomalías como DBSCAN y LOF, y la adopción de versiones robustas del algoritmo, como el Robust Sparse K-means Clustering (RSKC). Estas soluciones permiten mitigar los efectos negativos de los outliers y mejorar la precisión de la detección de fraudes.

Por último, se observa la relevancia del preprocesamiento de datos como la estandarización (Z-score), la winsorización y la eliminación de outliers mediante IQR, es esencial para mejorar el rendimiento de K-means en escenarios financieros. Un manejo adecuado

de los datos antes de aplicar el algoritmo es clave para asegurar que los centroides no sean distorsionados por puntos atípicos.

En conclusión, el análisis realizado en esta monografía ha demostrado que el algoritmo K-means presenta desafíos significativos en su aplicación a la detección de fraudes financieros, principalmente debido a su sensibilidad a los outliers y a sus limitaciones inherentes. Sin embargo, a través de soluciones como el uso de técnicas de preprocesamiento, la combinación con otros métodos de detección de anomalías y la adopción de versiones robustas del algoritmo, es posible mitigar estos efectos y mejorar su rendimiento en contextos financieros. La correcta preparación y tratamiento de los datos es fundamental para asegurar una agrupación precisa, lo que destaca la importancia de un enfoque integral que combine tanto la selección adecuada de técnicas como un manejo riguroso de los datos. Este estudio subraya la necesidad de seguir explorando y adaptando los métodos de aprendizaje automático para mejorar la detección de fraudes, un área crítica en el ámbito financiero.

### **Limitaciones**

El enfoque metodológico fue exclusivamente una revisión bibliográfica, por lo que no se realizó experimento propio ni validaciones empíricas en entornos reales o con conjuntos de datos financieros específicos. Esto implica una dependencia de la calidad, enfoque y resultados de los estudios analizados, que podrían contener sesgos metodológicos propios no controlados en esta monografía.

## Investigaciones Futuras

A continuación, se resaltan futuras investigaciones que pueden surgir después de realizada esta monografía:

Implementación empírica de las técnicas híbridas propuestas

Una línea de investigación clave es la implementación de sistemas de detección que integren K-means con algoritmos como DBSCAN, LOF o Isolation Forest en conjuntos de datos financieros reales, evaluando su desempeño cuantitativamente.

Validación de variantes robustas (como RSKC) en datos desbalanceados

Se recomienda desarrollar estudios empíricos que comparen la eficacia de versiones robustas de K-means frente a algoritmos supervisados en escenarios altamente desbalanceados típicos del fraude financiero.

Desarrollo de marcos automatizados de preprocesamiento de outliers

Dado el impacto del preprocesamiento en la mejora del rendimiento, futuras investigaciones podrían enfocarse en la automatización de procesos de filtrado y normalización estadística antes del clustering.

## Recomendaciones

A partir de los hallazgos y el análisis realizado en esta monografía, se proponen las siguientes recomendaciones para mejorar la efectividad del algoritmo K-means en la detección de fraudes financieros, especialmente en presencia de puntos atípicos. Estas recomendaciones buscan optimizar el uso del algoritmo y proporcionar soluciones prácticas para superar sus limitaciones, incrementando así la precisión de los sistemas de detección.

A continuación, se detallan las recomendaciones clave:

Se recomienda aplicar técnicas de preprocesamiento como la estandarización (Z-score), winsorización o filtrado por IQR antes de ejecutar el algoritmo K-means. Estas técnicas reducirán el impacto de los outliers en el cálculo de los centroides y mejorarán la precisión del modelo.

Es beneficioso integrar K-means con otros algoritmos especializados en la detección de outliers, como DBSCAN, LOF, o Isolation Forest. Esto permitirá utilizar K-means como un primer paso de análisis, seguido de una validación y refinamiento de los resultados con estos métodos adicionales.

Considerar el uso de versiones mejoradas de K-means, como el Robust Sparse K-means Clustering (RSKC), que incorporan mecanismos para excluir o minimizar el impacto de los outliers durante el proceso de agrupación, mejorando la fiabilidad de los resultados.

Para mejorar la efectividad del algoritmo en grandes volúmenes de datos, se recomienda aplicar técnicas de reducción de dimensionalidad como PCA antes de ejecutar K-means. Esto ayudará a reducir el ruido y enfocarse en las variables más relevantes para la detección de fraudes.

Se sugiere implementar una arquitectura escalonada donde K-means se utilice como filtro preliminar para manejar grandes volúmenes de datos, y luego se validen los resultados con algoritmos más robustos y complejos, mejorando tanto la eficiencia como la precisión.

En entornos financieros donde los datos se comparten entre instituciones, es recomendable explorar modelos mejorados de K-means que incorporen técnicas de criptografía homomórfica y compartición de secretos, garantizando la privacidad de los datos y aumentando la tolerancia a los datos dispersos.

## Referencias

- Ali, A., Abd Razak, S., Othman, S. H., Eisa, T. A. E., Al-Dhaqm, A., Nasser, M., Elhassan, T., Elshafie, H., & Saif, A. (2022). Financial Fraud Detection Based on Machine Learning: A Systematic Literature Review. *Applied Sciences*, 12(19), 9637.  
<https://doi.org/10.3390/app12199637>
- Bagheri-Gisour Marandyn, F., Larriva Novo, X. A., & Villagr a Gonz alez, V. A. (2021). *Dise o y evaluaci n de modelos de aprendizaje autom tico no supervisado para la detecci n de anomal as en un sistema Spark*. En: "VI Jornadas Nacionales (JNIC2021 LIVE)", 5/5/2021. <https://oa.upm.es/69003/>
- Barrag n Garnica, D. (2023). Patrones de comportamiento de clientes con tarjetas de cr dito de consumo con deterioro de calificaci n por riesgo utilizando K-means. *ODEON*, (22), 7-37. <https://doi.org/10.18601/17941113.n22.02>
- Boucher,  . (2020). *Outlier detection methods applied to financial fraud* (Tesis de maestr a). Lund University. <http://lup.lub.lu.se/student-papers/record/9030110>
- Carmona Mora, M., & Londo o Morales, L. M. (2021). *Modelos de machine learning para la detecci n de fraude financiero* (Tesis de especializaci n). Universidad de Antioquia. <https://hdl.handle.net/10495/20164>
- Choez Franco,  . S. (2022). *An lisis de las caracter sticas de los tipos de algoritmos de clustering en el aprendizaje no supervisado* (Tesis de grado). Universidad T cnica de Babahoyo. <http://dspace.utb.edu.ec/handle/49000/11594>
- Heras Calvo, D. (2023). *Biblioteca para la evaluaci n sistem tica de algoritmos de clustering*. Trabajo Fin de Grado / Proyecto Fin de Carrera, E.T.S. de Ingenieros Inform ticos (UPM), Boadilla del Monte. <https://oa.upm.es/75555/>

- Hidalgo Cornejo, I. L. M., & Solano Barragán, P. A. (2019). *Evaluación de métodos de clustering para el pronóstico de ventas en empresas productoras y distribuidoras de alimentos procesados* (Tesis de maestría). Universidad Peruana De Ciencias Aplicadas <http://hdl.handle.net/10757/655895>
- Hodge, V. J., & Austin, J. (2004). A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22(2), 85-126. <https://doi.org/10.1023/B:AIRE.0000045502.10941.a9>
- Huang, Z., Zheng, H., Li, C., & Che, C. (2024). Application of machine learning-based K-means clustering for financial fraud detection. *Academic Journal of Science and Technology*, 10(1), 33-39. <https://doi.org/10.54097/74414c90>
- Issa, H., & Vasarhelyi, M. A. (2011). Application of anomaly detection techniques to identify fraudulent refunds. *SSRN*. <https://ssrn.com/abstract=1910468>
- Khan, I. K., Daud, H. B., Zainuddin, N. B., Sokkalingam, R., Abdussamad, A., Museeb, A., & Inayat, A. (2024). Addressing limitations of the K-means clustering algorithm: Outliers, non-spherical data, and optimal cluster selection. *AIMS Mathematics*, 9(9), 25070-25097. <https://doi.org/10.3934/math.20241222>
- Liu, R., Qian, X., Mao, S., & Zhu, S. (2011). Research on anti-money laundering based on core decision tree. In Proceedings of the 2011 Chinese Control and Decision Conference (CCDC) (pp. 4323–4325). IEEE. 10.1109/CCDC.2011.5968986
- Liu, Y., Chen, C., Cui, J., Wang, L., & Wang, L. (2022). Scalable and Sparsity-Aware Privacy-Preserving K-means Clustering with Application to Fraud Detection. *ArXiv*. <https://doi.org/10.48550/arXiv.2208.06093>

- Nguyen, T. T., Tahir, H., Abdelrazek, M., & Babar, A. (2020). Deep Learning Methods for Credit Card Fraud Detection. *ArXiv*. <https://arxiv.org/abs/2012.03754v1>
- Orellana, M., & Cedillo, P. (2020). Detección de valores atípicos con técnicas de minería de datos y métodos estadísticos. *Enfoque UTE*, 11(1), 56-67.  
<https://doi.org/10.29019/enfoque.v11n1.584>
- Ponce Del Águila, A., Pérez Horna, S. Y., Lamadrid Elera, S., & Oliva Melgarejo, S. K. (2019). *Buenas prácticas en la gestión del riesgo de fraude interno: Casos de tres bancos de Lima Metropolitana* (Tesis de maestría). Pontificia Universidad Católica del Perú.  
<http://hdl.handle.net/20.500.12404/15152>
- Ramírez Mendoza, D. Y. (2022). *Métodos de machine learning con algoritmos de clúster no supervisados, una alternativa de segmentación de las pymes colombianas para plantear estrategias de acuerdo con sus condiciones económicas* (Tesis de maestría). Universidad EAFIT. <http://hdl.handle.net/10784/31580>
- Rodríguez Pérez, E. (2018). *Análisis y detección de fraude fiscal mediante técnicas de aprendizaje automático* (Tesis fin de master). Universidad Politécnica de Madrid.  
<https://oa.upm.es/55006/>
- Sahoo, G., & Sahoo, S. S. (2021). Accounting fraud detection using K-means clustering technique. In D. Swain, et al. (Eds.), *Machine Learning and Information Processing* (pp. 171-179). Springer. [https://doi.org/10.1007/978-981-33-4859-2\\_17](https://doi.org/10.1007/978-981-33-4859-2_17)
- Torres, R. (2023). *Clasificación de clientes y predicción de deserciones usando algoritmos K-means y regresión logística* (Trabajo de maestría). Pontificia Universidad Católica del Ecuador.

Vargas, J. M., Díaz, M., & García, F. (2018). Robust Clustering of Banks in Argentina. *Revista de Economía y Estadística*, 56(1), 21-41.

<https://doi.org/10.55444/2451.7321.2018.v56.n1.29385>