



Máquinas de soporte vectorial y árboles de clasificación para la detección de operaciones sospechosas de lavado de activos

Vector support machines and classification trees for detecting suspicious operations of money laundering

Marlon Efraín Gracia Granados*

(Recibido el 02-08-2018. Aprobado el 30-11-2018)

Estilo de citación de artículo:

M. E. Gracia Granados, "Máquinas de soporte vectorial y árboles de clasificación para la detección de operaciones sospechosas de lavado de activos", *Lámpsakos*, (21), pp. 26-38. (enero-junio, 2019). DOI: <https://doi.org/10.21501/21454086.2904>

Resumen.

El lavado de activos es un delito que trae consigo un gran número de consecuencias negativas a la sociedad en general. Para mitigar este problema en las entidades financieras, que es donde principalmente se presenta, se han desarrollado sistemas anti lavado de dinero. Lo anterior origina un nuevo problema: los falsos positivos que se obtienen a partir de dichos sistemas, los cuales representan para las entidades financieras pérdidas de dinero, tiempo y foco, al no tratar las verdaderas operaciones inusuales. Se evalúan los principales métodos de detección de operaciones inusuales de lavado de activos que se encuentran en la literatura, para determinar cuáles técnicas ofrecen los mejores resultados y a partir de estas generar un nuevo modelo que mejore los indicadores registrados. A partir de un proceso de revisión y replicación de metodologías de detección de anomalías encontradas en la literatura, se pudo generar un nuevo modelo que presenta mejores métricas a la hora de clasificar operaciones como normales e inusuales, lo cual puede representar para las entidades financieras una manera de disminuir las tasas de falsos positivos en sus sistemas anti lavado.

Palabras clave: Árboles de clasificación; Soporte vectorial; Métricas; Precisión; Entidades financieras; Falsos positivos; Lavado de activos; Operaciones inusuales; Operaciones sospechosas; Detección.

* Especialista en Estadística, Universidad Nacional sede Medellín, Medellín-Colombia, Correo electrónico: marlon.gracia@udea.edu.co

DOI: <https://doi.org/10.21501/21454086.2904>

Abstract

Money laundering is a crime that brings a large number of negative consequences to society in general. Anti-money laundering systems have been developed to mitigate this problem in financial institutions, which is where it is mainly presented. This causes a new problem: the false positives obtained from these systems, which represent financial losses for the financial entities, as well as time and focus, since they do not deal with the real unusual operations. The main detection methods of unusual operations of money laundering found in the literature are evaluated to determine which techniques offer the best results and from these generate a new model that improves the registered indicators. From a process of review and replication of anomalies detection methodologies found in the literature, a new model that presents better metrics when classifying operations as normal and unusual could be generated, this may represent way to reduce the false positive rates in their anti-money laundering systems in financial institutions.

Keywords: Classification trees, Vector Support, Metrics; Precision; Financial entities; False positives; Money laundering; Unusual operations; Suspicious operations; Detection.

1. INTRODUCCIÓN

El lavado de activos es un delito cuyo objetivo es dar apariencia de legalidad a dineros mal habidos, que son producto de delitos como el narcotráfico, venta de armas, prostitución, extorsión, entre otros [1]. En Colombia, el lavado de activos se encuentra tipificado en el artículo 323 del Código Penal [2], en el que se declaran todas las actividades por las cuales se puede incurrir en esta conducta criminal.

Este delito mueve, según ciertos estudios, entre el 2% y 5% del producto interno bruto mundial [3], [4], [5] y trae consigo efectos negativos sobre la sociedad en general, ya que puede distorsionar los precios de los bienes, alterar los diversos indicadores de un país y adicionalmente, puede ser un multiplicador de actividades ilegales al ayudar a las organizaciones criminales a autofinanciarse y a diversificar sus actividades delictivas [6], [7], [8], [9].

Las entidades financieras son regularmente el punto de entrada de este delito, debido a que los criminales hacen uso de estas para ingresar los dineros mal habidos, y a través de un conjunto de operaciones intentan perder el rastro del origen del dinero, para finalmente usarlo en el sector real sin temor a que se averigüe la procedencia no lícita; lo que constituye los pasos usuales del lavado de activos: colocación, transformación e integración [6].

Debido a lo anterior, en las entidades financieras se implementan sistemas ALD (anti lavado de dinero), con el fin de detectar y prevenir operaciones de lavado de activos realizadas en la entidad. El problema que surge es que estos sistemas suelen presentar un gran número de falsos positivos (operaciones normales que son reportadas por el sistema como sospechosas de lavado de activos) y falsos negativos (operaciones inusuales no detectadas por el sistema), porque por un lado hay operaciones que sí bien son inusuales, no son sospechosas, ejemplo de ello las ganancias ocasionales, y por otro, la astucia de los criminales para camuflar los dineros ilícitos [10].

La producción de falsos positivos y falsos negativos en los sistemas ALD se considera un problema por sus proporciones y por sus implicaciones; para dar una idea de qué tanto se produce este problema, se toma como referencia la encuesta realizada por Dow Jones Risks & Compliance y ACAMS a varios encargados de sistemas ALD en el mundo, en la cual se encontró que más del 44% de los encuestados respondieron que poseen tasas de falsos positivos mayores al 50%, y existen muchos casos en los cuales estos sistemas producen tasas de falsos positivos del 100% [11]. Las implicaciones de esta problemática es la pérdida de tiempo, esfuerzo y dinero en la investigación de operaciones, que al final fueron falsos positivos del sistema, así como la materialización de los efectos negativos del lavado de activos sobre las operaciones que no fueron detectadas por el sistema, es decir, los falsos negativos.

Para mitigar la alta producción de falsos positivos y negativos, varias técnicas de detección de anomalías han sido propuestas en la literatura [12], [13], [14], sin embargo, estas propuestas poseen diferentes métricas, conjuntos de datos y variables, por lo cual los resultados no son fácilmente comparables, y de esta manera es difícil determinar cuál posee mejores resultados.

Por tanto, el objetivo que se plantea en este artículo de investigación es implementar varias de las técnicas más citadas en el ámbito de la detección de operaciones sospechosas de lavado de activos y que posean la información suficiente para ser replicadas, con un conjunto de datos reales suministrado por una entidad financiera, obteniendo métricas que ayuden a determinar cuál de las técnicas tiene un mejor desempeño y usarla como base para generar un nuevo modelo de detección, con mejores características. Todo esto, buscando contribuir a la reducción de falsos positivos y negativos que se producen al detectar operaciones sospechosas de lavado de activos en el sistema financiero, con miras a reducir los efectos negativos que traen consigo.

DOI: <https://doi.org/10.21501/21454086.2904>

El presente artículo se encuentra organizado de la siguiente manera: 1. Se presenta la metodología a utilizar, 2. Se detallan los resultados de las técnicas seleccionadas en la literatura, 3. Se desarrolla la nueva técnica de detección de operaciones sospechosas y se presentan sus resultados, 4. Finalmente se exponen las conclusiones.

2. METODOLOGÍA

La metodología a utilizar tiene como punto de partida el problema planteado, que, en este caso, es la alta producción de falsos positivos por parte de diversos sistemas ALD y el efecto negativo que esta tiene en el sistema financiero y en la sociedad en general.

Para ello, se realizó un estudio del estado del arte de diferentes modelos de detección de lavado de activos, en el que se observó que dichos modelos realizan la clasificación de operaciones entre sospechosas y normales, teniendo en cuenta solo variables de tipo transaccional: montos, tipos y números de operaciones. Con esto en mente, la pregunta de investigación que se desea resolver es: ¿Incluir variables sociodemográficas de quienes realizan las operaciones que deben ser clasificadas entre sospechosas y normales, ayuda a disminuir el porcentaje de falsos positivos?

La metodología que se seguirá para resolver la pregunta de la investigación es la siguiente:

En una primera parte, se preparan los diferentes insumos para realizar la implementación y comparación de los diferentes modelos de detección; uno de ellos es la selección de las métricas que se utilizarán para determinar las bondades de cada método y compararlos, selección de las variables a utilizar, obtención de los datos y separación de estos en las poblaciones de entrenamiento y validación. Es muy importante que en los datos obtenidos se encuentre una marcación que indique operaciones sospechosas reales de lavado de activos, con lo cual se puede tener una mayor certeza de los resultados.

La segunda parte de la investigación consistirá en replicar los diferentes métodos seleccionados de la literatura y aplicárselos al conjunto de validación a utilizar, el cual es el mismo para todos, incluido el que se implementará; sobre los resultados al aplicar los métodos al conjunto de validación, se obtendrán las métricas seleccionadas y se hará un análisis de bondades y desventajas.

La última parte consistirá en crear un nuevo método de detección de anomalías, que se aproveche de la información sociodemográfica que fue agregada, y el cual busca superar las bondades de los métodos replicados de la literatura.

Métricas a utilizar

Sobre las métricas a utilizar para determinar las cualidades de detección de cada uno de los modelos a analizar, son muchas las propuestas que existen en la literatura [15], [16], [17]. Las que serán utilizadas en este artículo son unas métricas básicas, pero efectivas, que se obtienen de una tabla de contingencia como la que se muestra en la Tabla 1.

Tabla 1. Tabla de contingencia para clasificar las operaciones obtenidas

Clasificación planteada	Verdadera clasificación	
	Normal	Inusual
Normal	Verdadero negativo	Falso negativo
Inusual	Falso positivo	Verdadero positivo

Teniendo en cuenta la siguiente convención:

VP: Verdaderos positivos
VN: Verdaderos negativos
FP: Falsos positivos
FN: Falsos negativos

Se obtienen las siguientes métricas:

Error: métrica que mide el porcentaje de malas clasificaciones, se calcula según se indica en la Fórmula 1.

$$\text{Error} = 1 - \frac{VP+VN}{VP+VN+FP+FN} \quad (1)$$

Tasa de detección: métrica que mide el porcentaje de operaciones inusuales identificadas correctamente, se calcula según se indica en la Fórmula 2.

$$\text{Tasa de detección: } \frac{VP}{VP+FN} \quad (2)$$

Tasa de falsos positivos: métrica que mide el porcentaje de operaciones erróneamente clasificadas como inusuales, se calcula según se indica en la Fórmula 3.

$$\text{Tasa de falsos positivos: } \frac{FP}{FP+VN} \quad (3)$$

Precisión: métrica que mide el porcentaje de aciertos en las transacciones reportadas como inusuales, se calcula según se indica en la Fórmula 4.

$$\text{Precisión} = \frac{VP}{VP+FP} \quad (4)$$

Cada una de las métricas anteriores ayudan a caracterizar el modelo de clasificación; la pretensión es evaluar cada una de las anteriores métricas con respecto a cada modelo y, a partir de estas, determinar el modelo con las mejores características.

Variables a obtener

Para obtener el listado de variables a solicitar a la entidad financiera, lo primero que se requiere es determinar cuáles modelos de identificación de operaciones sospechosas de lavado de activos se intentará replicar.

Luego del análisis de la literatura, los modelos escogidos fueron: Redes bayesianas dinámicas [18], Redes neuronales de base radial [19], Máquinas de soporte vectorial [20], Clúster de dos fases [21], Maximización de la esperanza [22] y Sequence Matching [23].

Las variables utilizadas por estos modelos son: montos de entrada y salida a las cuentas de los clientes, frecuencia de entrada y salida a las cuentas de los clientes, tiempos entre transacciones, tipos de transacciones (efectivo, electrónica), periodo de la operación (inicio del mes, mitad del mes, final del mes) y una marcación que determina si la persona es o no sospechosa de lavado de activos en la entidad financiera.

Las variables sociodemográficas escogidas que se cree pueden aportar a caracterizar mejor a quien realiza las operaciones y clasificarlas son: edad, ingresos declarados, egresos declarados, patrimonio, ventas, segmento de cliente (clasificación que le da la entidad al cliente), tipo de persona (natural, jurídica), ocupación, actividad económica, montos enviados y recibidos internacionalmente, montos desembolsados en créditos.

Obtención de datos

Las variables enunciadas en la etapa anterior fueron solicitadas a una entidad financiera local, de la cual se obtuvo la información completamente anónima y correspondiente a un año.

Dado que algunos de los modelos elegidos requieren información consolidada, lo obtenido fue separado en dos bases independientes, una que contiene el detalle, registro a registro y otra que se encuentra consolidada a nivel de persona.

El número de registros de las dos bases son diferentes, dada la anterior condición. La base de datos consolidada, como va a nivel de persona, indicando la actividad de cada una en todo un año, es apenas de 7 millones de registros; mientras que la base detallada posee 156 millones de registros, esto es algo a tener en cuenta en la división de la población que será llevada a cabo en la siguiente etapa.

DOI: <https://doi.org/10.21501/21454086.2904>

Población de entrenamiento y de validación

Las dos bases obtenidas en el anterior paso requieren ser subdivididas en dos conjuntos independientes, uno de entrenamiento a partir del cual se entrenarán todos los modelos, aunque algunos de ellos requieran todos los datos o solo un subconjunto, y otro como población de validación, con el cual todos los modelos serán ejecutados y, a partir de allí, se tomarán las métricas con las que se determinará el modelo con las mejores características.

Con la división del universo de registros en estas dos bases, la de entrenamiento y validación, hay tres temas a tener en cuenta; el primero es definir cuál será el porcentaje de datos destinado a cada fin. Para esto existen recomendaciones de cómo realizar la subdivisión [24], aunque esto exige cuidado cuando el conjunto de datos es muy limitado; en el caso de esta investigación se poseen dos bases con un gran número de registros, así existe mayor flexibilidad a la hora de escoger los porcentajes; y debido a que es importante tener un buen número de registros para validar los modelos replicados, se elige particionar los conjuntos de datos en 40% para la población de entrenamiento y 60% para la población de validación.

El segundo tema a tener en cuenta es que la población que se está tratando no está balanceada, es decir, existen muchos más registros normales que inusuales. Esto implica que, al realizar la subdivisión en población de entrenamiento y validación, se debe cuidar que también la distribución de operaciones normales e inusuales se mantenga en los dos subconjuntos, de lo contrario podría tener impacto en los resultados [25]; por tanto, si en el universo de datos un 98% es normal y un 2% es inusual, esta proporción se debe intentar mantener en las poblaciones de entrenamiento y validación.

Finalmente, el último asunto a considerar es que si bien la población está desbalanceada, algunos estudios sugieren que, a nivel del entrenamiento del modelo, los dos conjuntos de datos a clasificar posean un número similar de datos [26], [27]; esto aplica para los modelos que requieren tomar muestras de la población de entrenamiento.

La manera para realizar este balanceo de datos, para las muestras de entrenamiento, será el balanceo por debajo [28], que consiste en tomar la totalidad de inusualidades de la población de entrenamiento, y ponerlo en la muestra balanceada; seguido se elige aleatoriamente un número igual de registros no inusuales, de tal manera que la muestra final para ser utilizada en el entrenamiento de los modelos que requieran muestras más pequeñas, posean un número igual de operaciones normales e inusuales.

3. IMPLEMENTACIÓN DE LOS MODELOS DE DETECCIÓN

Con los pasos previos realizados, ya se cuenta con bases detalladas y consolidadas, tanto de entrenamiento como de validación, además tienen todas las variables requeridas por los modelos estudiados, adicional a ello se obtuvieron las métricas a calcular con las cuales se determinarán los pro y contras de los modelos de detección, solo resta implementarlos con la información obtenida y discutir los resultados que estos ofrecen.

En las siguientes secciones se resumen las técnicas utilizadas, los modelos matemáticos en los que se sustentan y los resultados obtenidos al replicar dichas técnicas, partiendo de los datos de validación que se poseen; sin embargo, no se dará el detalle de las mismas, ya que se relacionan los artículos originales expuestos por sus autores, por sí se desea mayor profundización.

Redes bayesianas dinámicas

La primera técnica es la implementada por Raza y Haider, que consiste en la utilización de una red bayesiana dinámica para la obtención de datos anómalos en el ámbito del lavado de activos [18].

La implementación de esta técnica está separada en tres etapas; la primera, de segmentación, en la que cada uno de los individuos es asignado a un grupo con

el que comparte características comunes, con el que se realizan fácilmente comparaciones. La segunda etapa corresponde a la identificación de la Red Bayesiana Dinámica de cada uno de los grupos en los que se segmentó la población. La tercera y última etapa es la de identificación de anomalías, que se hace con base en la Red Bayesiana identificada en la etapa anterior y con una métrica especial creada por los autores de esta técnica de detección, denominada AIRE (índice de anomalía utilizando rango y entropía).

Los resultados de esta técnica, luego de ser replicados con los datos y variables que se poseen, son los presentados en la Tabla 2.

Tabla 2. Resultados implementación Red bayesiana dinámica

Modelo	Red Bayesiana Dinámica
Número de muestras	2.957.028
Error	0,89%
Verdaderos positivos	197
Verdaderos negativos	2.930.627
Falsos Positivos	19.034
Falsos negativos	7.170
Tasa de detección	2,67%
Tasa de falsos positivos	0,65%
Precisión	1,02%

Se puede concluir de las mediciones anteriores, lo siguiente:

- Los puntos positivos de la técnica radican en las bajas tasas de error y de falsos positivos, pero los puntos negativos están en la poca precisión y la muy baja tasa de detección, lo que vuelve la técnica poco aplicable como modelo principal de detección.
- Por otra parte, este modelo aún se podría considerar un modelo de filtrado, con el que se eliminarían muchos registros de la muestra que no son anómalos, y que pueden ser utilizados por modelos más precisos para obtener mejores resultados.

Redes neuronales de base radial

Este modelo de detección fue generado utilizando la técnica de redes neuronales de base radial, para obtener las transacciones anómalas dentro de un conjunto de datos [19].

El modelo se divide en tres partes, la primera es la obtención de las variables de entrada, que antes de ser analizadas por el modelo, son normalizadas de manera que sus valores estén en el rango 0 a 1. El segundo paso corresponde a la conformación de las n capas ocultas, que posee la red neuronal, a partir de los datos de entrada y utilizando un algoritmo de APCIII Clúster que resulta en un número n de grupos que se deberían conformar; y de estos grupos se obtiene el ancho gaussiano de los mismos, indispensable para ejecutar el modelo. Finalmente, el último paso consiste en encontrar los pesos con los cuales se puedan mapear las capas ocultas a las de salida, siendo dos (probabilidad de ser usual y probabilidad de ser inusual), para esto se utiliza un algoritmo de RLS con el que se completan los datos solicitados por el modelo.

Dado lo anterior, se ejecuta el modelo con la población de entrenamiento y posteriormente se envía la población de validación para obtener los resultados de clasificación y realizar las mediciones planteadas; los resultados están presentados en la Tabla 3.

Tabla 3. Resultados modelo de Red neuronal de base radial

Modelo	Red neuronal de base radial
Número de muestras	4.514.994
Error	52,47%
Verdaderos positivos	18.117
Verdaderos negativos	3
Falsos positivos	2.127.981
Falsos negativos	5.779
Tasa de detección	75,82%
Tasa de falsos positivos	52,62%
Precisión	0,76%

Fuente: elaboración propia.

DOI: <https://doi.org/10.21501/21454086.2904>

La conclusión sobre este modelo fue:

- Se puede observar que este modelo posee una buena tasa de detección, no obstante, la tasa de error es demasiado alta, al igual que la tasa de falsos positivos; mientras que la precisión se mantiene muy baja, por lo cual es un modelo de poca aplicabilidad en un ambiente real.

Máquinas de soporte vectorial

Este modelo utiliza la técnica de máquinas de soporte vectorial, para realizar hiperplanos que dividan la población entre usual e inusual, a partir de la población de entrenamiento, con el cual se logra clasificar la población de validación [20].

La técnica consiste básicamente en dos pasos; en el primero se determinan las variables a ingresar y se les realiza un pequeño proceso de transformación, que busca estandarizar el rango de cada una de estas; el segundo consiste en ingresar los datos que recibe la técnica, que son el factor de control y el castigo por mala clasificación; teniendo esto presente basta con entrenar el modelo con individuos de la población de entrenamiento y posteriormente realizar la clasificación de la población de validación. Una vez realizado esto, se obtienen las métricas que se consignan en la Tabla 4.

Tabla 4. Resultados máquinas de soporte vectorial

Modelo	Máquina de soporte vectorial
Número de muestras	4.514.994
Error	2,61%
Verdaderos positivos	13.860
Verdaderos negativos	4.383.436
Falsos positivos	107.662
Falsos negativos	10.036
Tasa de detección	58,00%
Tasa de falsos positivos	2,40%
Precisión	11,41%

Las conclusiones sobre este modelo son las siguientes:

- Este es el modelo más interesante hasta el momento, ya que posee una tasa de error baja, una tasa de detección significativa, además una tasa de falsos positivos baja y su precisión es alta comparada con las métricas anteriores.
- El problema que podría surgir de esta técnica es el alto volumen de datos clasificados como inusuales, en total 121.522 datos, que representa un gran número de inusualidades para ser atendidas por un equipo de Compliance.

Clúster de dos fases

El Clúster de dos fases se basa en la modificación de una técnica de clasificación en conjunto con un modelo de escogencia de grupos anómalos [21].

Como su nombre lo dice, esta posee dos fases; en la primera fase propone que sobre el grupo de datos se realice una técnica de clasificación por medio de un algoritmo de K medias, pero este algoritmo ha sido modificado de tal manera que no le importe que el número de individuos de los grupos que entregue quede desbalanceado.

La segunda fase toma los grupos arrojados en la primera fase y calcula sobre el promedio de distancia de estos, un árbol de expansión mínimo sobre el cual se elimina la arista más grande para obtener así dos poblaciones separadas y, mediante un conteo de elementos en estas poblaciones, se determina el conjunto inusual.

Al comparar la clasificación realizada con la población de validación, se obtienen las muestras que se presentan en la Tabla 5.

Tabla 5. Resultados Clúster de dos fases

Modelo	Clúster de dos fases
Número de muestras	4.514.994
Error	0,530%
Verdaderos positivos	7
Verdaderos negativos	4.491.045
Falsos positivos	53
Falsos negativos	23.889
Tasa de detección	0,029%
Tasa de falsos positivos	0,001%
Precisión	11,67%

Las conclusiones sobre este modelo son las siguientes:

- Se observa en las métricas obtenidas que el error y la tasa de falsos positivos es muy bajo, lo cual es positivo, pero la tasa de detección también lo es. Por otra parte, la precisión comparada con las anteriores es alta.
- Este modelo, aunque tenga una tasa de detección tan baja, aún puede ser aplicable en el mundo real, gracias a que los resultados arrojados por esta técnica son pocos y aunque esto suceda entrega datos inusuales, asimismo, el esfuerzo que llevaría estudiar todos los casos arrojados no es tan elevado y se está obteniendo buena recompensa al reportar varias operaciones que en efecto son inusuales.

Maximización de la esperanza

Este modelo busca por medio del algoritmo de maximización de la esperanza, encontrar de manera iterativa el conjunto de datos que maximiza la probabilidad de ser inusuales [22].

Esta técnica es no supervisada, por lo que no requiere conocer si los individuos son inusuales o no en la etapa de entrenamiento, pues en un número de iteraciones realiza la agrupación de los individuos en varios grupos, en los cuales mide ciertas métricas para determinar la correctitud de estas clasificaciones, cuando llega al final del número de iteraciones parametrizadas o no supera la tolerancia indicada, arroja cuál de los grupos formados es el inusual.

Al llegar a los resultados de transacciones anómalas, se obtienen las métricas definidas, las cuales se presentan en la Tabla 6.

Tabla 6. Resultados Maximización de la esperanza

Modelo	Maximización de la esperanza
Número de muestras	4.514.994
Error	0,530%
Verdaderos positivos	10
Verdaderos negativos	4.491.015
Falsos positivos	83
Falsos negativos	23.886
Tasa de detección	0,04%
Tasa de falsos positivos	0,002%
Precisión	10,75%

Las conclusiones sobre este modelo son las siguientes:

- Este modelo también presenta una tasa de error baja, tasa de falsos positivos baja, pero una de detección baja. La precisión es razonable vistos los modelos anteriores.
- Al igual que el modelo de Clúster de dos fases, los resultados son muy pocos, pero dentro de estos hay operaciones sospechosas, lo cual le da aplicabilidad en el mundo real evaluando el esfuerzo/ganancia.

Sequence Matching

Sequence Matching tiene un enfoque diferente y en vez de observar una transacción, observa una secuencia de esta comparándola simultáneamente con secuencias de referencia para determinar si son inusuales o no [23].

Esta técnica se realiza en cuatro partes; una en la que para cada variable se obtienen umbrales, para determinar a partir de qué punto se puede considerar el valor como "muy grande"; en la segunda se seleccionan las transacciones en las que por lo menos una variable de una operación supere el umbral definido anteriormente; en la tercera se seleccionan transac-

DOI: <https://doi.org/10.21501/21454086.2904>

ciones de referencia; y en la cuarta se comparan las transacciones elegidas en la segunda con las elegidas en la tercera y las que definitivamente tengan un comportamiento diferente a dichas operaciones de referencia son clasificadas como inusuales.

Al implementar esta técnica con la población de validación se obtienen las métricas que se presentan en la Tabla 7.

Tabla 7. Resultados Sequence Matching

Modelo	Sequence Matching
Número de muestras	3.437.420
Error	1,6%
Verdaderos positivos	383
Verdaderos negativos	3.381.921
Falsos positivos	47.257
Falsos negativos	7.859
Tasa de detección	4,64%
Tasa de falsos positivos	1,38%
Precisión	0,8%

Las conclusiones sobre este modelo son las siguientes:

- De esta técnica hay poco que rescatar, si bien la tasa de errores y de falsos positivos es relativamente baja, la tasa de detección y la precisión también lo es, lo cual le resta mucha aplicabilidad en el sector real.

4. IMPLEMENTACIÓN DEL NUEVO MODELO

A partir del conocimiento adquirido al implementar los anteriores modelos, obtener las métricas de cada uno y poder observar sus fortalezas y debilidades, se formula un nuevo modelo con el cual se espera conseguir mejores indicadores.

De las técnicas revisadas, la que presenta métricas más robustas en cuanto a que maneja errores bajos, con una tasa de falsos positivos también baja, pero una tasa de detección y una precisión decente, es la técnica de máquina de soporte vectorial; entre sus

virtudes se encuentra la capacidad de filtrar un gran número de operaciones no inusuales manteniendo un grupo de inusualidades considerable, no obstante la desventaja radica que en un ambiente real esa cantidad de registros puede llegar a ser difícil de atender en profundidad.

Por otro lado, se tiene una técnica sencilla de implementar, pues sus resultados siguen siendo decentes y, sobre todo, es aplicable en el sector real debido a que con muy poco esfuerzo se pueden obtener operaciones inusuales en un gran conjunto de datos; esa técnica es la de clúster de dos fases. Otro aspecto importante de esta es que toma varias técnicas y las aplica consecutivamente a un conjunto de datos para ir refinando los resultados encontrados.

En ese orden de ideas, tomando esos dos aprendizajes obtenidos luego de la replicación de los diferentes modelos de detección estudiados, el que la máquina de soporte vectorial sirve como buena técnica para filtrar datos no inusuales y que se pueden encadenar métodos diferentes a un grupo de datos para obtener mejores resultados, se forma el nuevo modelo de detección que se plantea en este artículo.

El modelo que se plantea en la presente investigación posee cuatro fases; en la primera de ellas se determina qué variables serán utilizadas en el modelo, para ello se utilizará la población de entrenamiento y se realizará un proceso de selección de variables Backward, Forward y Stepwise sobre el total de variables contra la variable respuesta, si la persona tuvo operaciones inusuales o no. Luego de ello, se tomarán todas las variables que por lo menos hayan sido seleccionadas por uno de los métodos anteriores.

En la segunda fase se aplicará la técnica de máquina de soporte vectorial sobre las variables seleccionadas en el anterior paso; esto en la población de entrenamiento. El objetivo de la técnica en este punto será filtrar operaciones no inusuales tratando de disminuir la población y manteniendo un número considerable de operaciones inusuales.

En la tercera fase, sobre las operaciones que superaron el anterior filtro se construirá un árbol de clasificación, con el objetivo de obtener un conjunto de datos no tan amplio como la población de entrenamiento inicial y que tenga las características de transacciones capaces de superar el filtro de las máquinas de soporte vectorial; pero como siguen siendo operaciones de entrenamiento, mantienen la marcación de si la operación es inusual o no, a partir de lo cual se crea el entrenamiento del árbol de clasificación.

La última fase consiste en aplicar en la población de validación las fases dos y tres, con lo cual sus datos serán filtrados gracias a la técnica de máquina de soporte vectorial entrenada en la fase dos y con las variables obtenidas en la fase 1; posteriormente será aplicado, sobre los datos restantes, el árbol de clasificación entrenado en la fase tres; y los que sean clasificados por este como inusuales, serán los reportados como tal por la nueva técnica.

Al aplicar estas cuatro fases sobre la población de validación se obtienen las métricas observadas en la Tabla 8.

Tabla 8. Resultados nuevo modelo

Modelo	Nuevo Modelo
Número de muestras	4.514.994
Error	0,66%
Verdaderos positivos	1.756
Verdaderos negativos	4.483.282
Falsos positivos	7.816
Falsos negativos	22.140
Tasa de detección	7,35%
Tasa de falsos positivos	0,17%
Precisión	18,34%

Como se puede ver, comparativamente son muchas las bondades de este nuevo modelo sobre los modelos estudiados y contrastados anteriormente, la precisión obtenida del 18.34% mejora en un 57.16% más la mejor precisión obtenida hasta ahora que era la del clúster de dos fases; las tasas de error y de falsos positivos siguen siendo significativamente muy bajas, del 0.66% y 0.17%, respectivamente; lo más importante del mé-

todo es la cantidad de registros que al final reporta como sospechosos, tan solo 9.572 individuos, teniendo en cuenta el gran número de individuos del que se parte originalmente (4.514.994 individuos), el que se tenga un método que clasifique a una cantidad no tan significativa como sospechosa teniendo una precisión bastante buena según los indicadores de referencia resulta provechoso, porque este número de registros brinda la posibilidad de realizar una investigación minuciosa en cada uno de los individuos marcados como sospechosos por el modelo; así gran parte del esfuerzo realizado no se perderá y sube considerablemente la cantidad de reportes realizados a los entes superiores.

5. TRABAJOS FUTUROS

Los resultados de la implementación del modelo creado fueron positivos en comparación a lo revisado en la literatura; no obstante, hay modelos más robustos que pueden ser empleados con el fin de refinar mucho más los resultados, porque si bien lo arrojado es aplicable en la vida real con un resultado aceptable, continúa siendo alto el umbral de operaciones no detectadas.

Esta investigación otorga un punto de partida sobre la utilización de variables no transaccionales en la detección de operaciones sospechosas de lavado de activos, y cómo éstas pueden ayudar a la obtención de mejores resultados.

6. CONCLUSIONES

Este artículo da cuenta de una breve introducción al problema del lavado de activos, las consecuencias de este en la sociedad y de cómo las entidades financieras trabajan para mitigar su impacto implementando sistemas de detección de lavado de activos, con lo cual se introduce un nuevo problema que radica en poder detectar efectivamente el lavado de activos sin producir altos índices de falsos positivos, pues esto ocasiona pérdida de recursos como tiempo de personas encargadas de investigar estas operaciones y di-

DOI: <https://doi.org/10.21501/21454086.2904>

nero de las entidades financieras, además que la poca eficiencia de los sistemas de detección da cabida a que se materialicen los riesgos del lavado de activos.

Al implementar las principales metodologías de lavado de activos encontradas en la literatura, las cuales solo hacen uso de variables transaccionales y no tienen en cuenta una caracterización de quién realiza la operación, se compararon con unas métricas definidas en el artículo, estableciendo que la técnica que más bondades ofrece es la máquina de soporte vectorial. Este conocimiento y la experiencia obtenida al implementar las otras técnicas permitió crear una nueva, la cual, sí hace uso de variables sociodemográficas de quien realiza la operación, y fue evaluada con las mismas métricas, así se estableció que ofrece mejores resultados.

Finalmente, la nueva técnica implementada consistió en la utilización del mejor modelo estudiado, que es la máquina de soporte vectorial, así los datos arrojados por este fueron posteriormente utilizados por un árbol de clasificación consiguiendo ser refinados y se obtuvieron mejores resultados.

REFERENCIAS

- [1] B. Buchanan, "Money laundering a global obstacle", *Res. Int. Bus. Finance*, vol. 18, no. 1, pp. 115-127, April. 2004. DOI: <https://doi.org/10.1016/j.ribaf.2004.02.001>
- [2] Congreso de Colombia, "Ley 599 de 2000, Art. 323", *Diario Oficial* no. 44.097, julio 24 de 2000. Recuperado de <https://www.wipo.int/edocs/lexdocs/laws/es/co/co028es.pdf>
- [3] B. Unger et al., *The amounts and effects of money laundering*, Ministry of Finance, 2006.
- [4] UNODC, "Estimating illicit financial flows resulting from drug trafficking and other transnational organized crimes". United Nations Office on Drugs and Crime, 2011.
- [5] J. Walker and B. Unger, "Measuring global money laundering: the Walker gravity model", *Review of Law & Economics*, vol. 5, no. 2, pp. 821-853, 2009.
- [6] M. Levi, "Money Laundering and Its Regulation", *Annals of the American Academy of Political and Social Science*, vol. 582, no 1, pp. 181-194, January. 2002. DOI: <https://doi.org/10.1177/000271620258200113>
- [7] F. A. T. Force, "The forty recommendations", 2003.
- [8] B. Bartlett, "The negative effects of money laundering on economic development", *The Asian Development Bank Regional Technical Assistance Project No. 5967*, 2002.
- [9] P. J. Quirk, "Money laundering: muddying the macroeconomy", *Finance Dev.*, vol. 34, pp. 7-9, 1997. Retrieved from <https://www.imf.org/external/pubs/ft/fandd/1997/03/pdf/quirk.pdf>
- [10] R. Menon and S. Kuman, «Understanding the role of technology in anti money laundering compliance», *Infosys Technol. Ltd*, vol. 1, pp. 2-4, 2005.
- [11] DOWJONES, *AML Survey Results from Dow Jones Risk & Compliance & ACAMS*, 2011.
- [12] Z. Gao and M. Ye, "A framework for data mining based anti money laundering research", *Journal of Money Laundering Control*, vol. 10, no. 2, pp. 170-179, 2007. DOI: <https://doi.org/10.1108/13685200710746875>
- [13] N. A. L. Khac and M. T. Kechadi, "Application of Data Mining for Anti money Laundering Detection: A Case Study", in *2010 IEEE International Conference on Data Mining Workshops*, 2010, pp. 577-584.
- [14] K. D. Rohit and D. B. Patel, "Review On Detection of Suspicious Transaction In Anti Money Laundering Using Data Mining Framework", *International Journal for Innovative Research in Science & Technology*, vol. 1, no. 8, pp. 129-133, 2015.

- [15] P. Baldi, S. Brunak, Y. Chauvin, C. A. Andersen, and H. Nielsen, "Assessing the accuracy of prediction algorithms for classification: an overview", *Bioinformatics*, vol. 16, no. 5, pp. 412-424, 2000. DOI: <https://academic.oup.com/bioinformatics/article/16/5/412/192336>
- [16] T. Fawcett, "An introduction to ROC analysis", *Pattern Recognit. Lett.*, vol. 27, no. 8, pp. 861-874, 2006. DOI: <https://doi.org/10.1016/j.patrec.2005.10.010>
- [17] C. Liu, P. M. Berry, T. P. Dawson, and R. G. Pearson, "Selecting thresholds of occurrence in the prediction of species distributions", *Ecography*, vol. 28, no. 3, pp. 385-393, 2005. DOI: <https://doi.org/10.1111/j.0906-7590.2005.03957.x>
- [18] S. Raza and S. Haider, "Suspicious activity reporting using dynamic bayesian networks", *Procedia Computer Science*, vol. 3, pp. 987-991, 2011. DOI: <https://doi.org/10.1016/j.procs.2010.12.162>
- [19] L.T. Lv, N. Ji, and J.L. Zhang, "A RBF neural network model for anti money laundering", *International Conference on Wavelet Analysis and Pattern Recognition*, 2008. ICWAPR '08, 2008, vol. 1, pp. 209-215. DOI: 10.1109/ICWAPR.2008.4635778
- [20] J. Tang and J. Yin, "Developing an intelligent data discriminating system of anti money laundering based on SVM", in *Proceedings of 2005 International Conference on Machine Learning and Cybernetics*, 2005, vol. 6, pp. 3453-3457. DOI: 10.1109/ICMLC.2005.1527539
- [21] M. F. Jiang, S. S. Tseng, and C. M. Su, "Two phase clustering process for outliers detection", *Pattern Recognition Letters*, vol. 22, no. 6-7, pp. 691-700, may 2001. DOI: [https://doi.org/10.1016/S0167-8655\(00\)00131-8](https://doi.org/10.1016/S0167-8655(00)00131-8)
- [22] Z. Chen, L. D. V. Khoa, A. Nazir, E. N. Teoh, and E. K. Karupiah, "Exploration of the effectiveness of expectation maximization algorithm for suspicious transaction detection in anti money laundering", in *IEEE Conference on Open Systems (ICOS)*, 2014, pp. 145-149. DOI: 10.1109/UEMCON.2016.7777919
- [23] X. Liu, P. Zhang, and D. Zeng, "Sequence matching for suspicious activity detection in anti money laundering", in *Intelligence and Security Informatics, Springer*, 2008, pp. 50-61. DOI: 10.1007/978-3-540-69304-8_6
- [24] J. Friedman, T. Hastie, and R. Tibshirani, *The elements of statistical learning*, vol. 1. Springer series in statistics Springer, Berlin, 2001.
- [25] Q. Wei and R. L. Dunbrack Jr, "The role of balanced training and testing data sets for binary classifiers in bioinformatics", *PLoS One*, vol. 8, no. 7, p. e67863, 2013. DOI: <https://doi.org/10.1371/journal.pone.0067863>
- [26] N. V. Chawla, "C4. 5 and imbalanced data sets: investigating the effect of sampling method, probabilistic estimate, and decision tree structure", in *Proceedings of the ICML*, 2003, DOI: https://doi.org/10.1007/978-3-540-69304-8_6
- [27] N. Japkowicz, "Learning from imbalanced data sets: a comparison of various strategies", in *AAA/ workshop on learning from imbalanced data sets*, 2000, vol. 68, pp. 10-15. Retrieved from <https://pdfs.semanticscholar.org/1af9/6acae07b1e141f98f3df973eaf9e0a9226fb.pdf>
- [28] C. Drummond and R. C. Holte, "C4. 5, class imbalance, and cost sensitivity: why under sampling beats over sampling", in *Workshop on learning from imbalanced datasets II*, 2003, vol. 11. Retrieved from <https://pdfs.semanticscholar.org/144b/bbaf2f0876c23295019b6e380c9fe4feda3.pdf>