

Las técnicas de aprendizaje automático para tareas de ingeniería de software relacionadas con la ciberseguridad son cada vez más populares.

En este survey se estudian los diferentes enfoques basados en el aprendizaje automático, existentes y se muestra que tipo de representaciones se utilizaron en tareas para ciberseguridad y lenguajes de programación.

- Se ha descubierto que las representaciones de grafos son la categoría más popular, y que los tokenizadores y los árboles de sintaxis abstracta (AST) son las dos representaciones más populares en general.
- La tarea de detección de vulnerabilidades es la más popular, y que el lenguaje más usado es C para la mayoría de técnicas.
- Los modelos basados en secuencias son la categoría de modelos más popular, y que las máquinas de vectores de soporte son el modelo más popular en general.

1. Introducción:

Las vulnerabilidades de software son defectos que afectan a las propiedades de seguridad previstas por un sistema de software y que permiten a los atacantes hacer acciones maliciosas.

A medida que los humanos tienden a volver su vida más dependiente de la tecnología los proveedores de software se ven forzados a diseñar software cada vez más seguros. Existen varias prácticas para abordar este tipo de problemas, y se pueden hacer en cada fase del ciclo de vida del software, aunque estas prácticas ayudan de alguna manera u otra para evaluar y mejorar estos problemas de seguridad, toman mucho tiempo, llegan a ser costosas y pueden ser propensas a errores. Entonces es aquí donde entra la IA y el ML que se aplican a tareas de ciberseguridad, como la detección de vulnerabilidades, detección de malware, etc. Entonces con esto nos damos cuenta que estas técnicas son valiosas, ya que ayudan a mejorar la calidad del software y ahorramos tiempo, que de otro modo serían propensas a errores y llevarían mucho tiempo.

Un ejemplo claro sería: un modelo que detecte vulnerabilidades antes de la implementación, esto ahorraría tiempo y dinero, y a la par aumentaría la seguridad del sistema en conjunto.

Los modelos de ML no son capaces de procesar el código fuente como tal, por ende se busca la manera de pasarlo a datos numéricos para ajustar las ponderaciones necesarias para el aprendizaje, esto quiere decir que debemos buscar la manera que el código fuente sea entendible de manera eficiente tanto en su información estructural como semánticamente, esto hace que el modelo aprenda y comprenda de los patrones subyacentes en el código. Por esto el código fuente es una parte crucial durante el desarrollo de técnicas de ML, un código se puede representar como: AST (ÁRBOL DE SINTAXIS ABSTRACTA), CFG (gráfico de flujo de control), o tokenizado.

En este artículo veremos un análisis sistemático de la literatura, para comprender el estado actual de técnica de representación del código fuente para técnicas basadas en ML para tareas de ciberseguridad, Investigaremos los lenguajes con los cuales se hacen

estas técnicas, la popularidad de ciertas representaciones y los típicos modelos de ML. El objetivo del estudio es permitir a los investigadores comprender las lagunas en este ámbito, en particular si hay ciertas tareas de ciberseguridad o lenguajes que son descuidados por las técnicas existentes.

2. Background: En esta sección se aborda la terminología para que el público entienda de manera correcta el survey.

2.1) ML para ing. de software segura

Los ing de software se encargan de generar código nuevo de calidad, y son responsables de ver que errores hay, que vulnerabilidades existen y cualquier otro problema del código fuente, estas tareas pueden llegar a ser tediosas, propensas a errores, y lentas. Una tarea de ciberseguridad es aquella que tiene como objetivo frustrar a los posibles intrusos, por lo tanto las tareas que hacen los ing de soft. son aquellas que se encargan de resolver e identificar vulnerabilidades en el código que podría permitir al atacante aprovecharse del sistema. Es ahí donde entra el ML, la IA podría ayudar a mejorar en tiempo y eficiencia estos trabajos y así mitigar los problemas que suceden en estos casos y mejorar la calidad del código fuente que se pone a disposición del público.

2.2) Representación del código fuente e incrustaciones de código.

Una representación de código fuente captura sintaxis y la semántica del código fuente, de modo que el modelo de ML puede aprender características clave, para poder representar un código fuente hay muchas formas, algunas son: representación de grafos, técnicas de procesamiento de lenguaje natural, etc. Estas representaciones ofrecen información diferente sobre el código y de esta manera hacen que el modelo de ML aprenda más cosas sobre él. Por ejemplo las técnicas de NLP se basan más en la información semántica del código, por lo que el modelo aprenderá relaciones semánticas, más no relaciones estructuradas del código.

Los modelos aprenden a base de incrustaciones vectoriales, estas son una forma de baja dimensión para representar datos de alta dimensión, para este caso las incrustaciones se crean a partir de la representación del código fuente. La representación del código fuente es el primer nivel de abstracción del código fuente original, a esto se le llama fase extracción de características, los datos originales se transforman en datos de menor dimensión, esto tiene como objetivo preservar la mayor cantidad de información de los posibles datos originales. El segundo nivel de abstracción son las incrustaciones vectoriales y son las que permiten aplicar técnicas de ML a las representaciones transformándolas en representación de forma numérica, que es la forma en la cual la máquina es capaz de entender. Nos centraremos en el primer nivel de abstracción, aunque hay varias técnicas para generar incrustaciones como word2vec, doc2vec y GraphCode2vec.

3. Related Work: Según el conocimiento del creador de este survey, menciona que existen muchos más como este survey pero la diferencia radica en que los otros artículos se centran principalmente en los modelos de ML utilizados, pero pocos mencionan o describen detalladamente las representaciones utilizadas en estos trabajos. El artículo se centra principalmente en el tema de las representaciones utilizadas para el ML, en

tareas relacionadas a la seguridad. Se hace una comparativa del survey de los autores con otros surveys que ya existen, el enfoque de este survey es algo netamente referente a las representaciones de código fuente que se usan para técnicas de ML para tareas de ciberseguridad.

- Comparación con Surveys Generales sobre ML en la ing de soft
- Comparación con estudios específicos en detección de vulnerabilidades.
- Comparación con revisiones sobre técnicas de PLN.
- Comparación con un estudio terciario.
- Comparación con artículos que crean taxonomía
- Comparación con surveys sobre ciberseguridad en general

4. Metodología: Para llevar a cabo nuestra revisión sistemática, que implica 3 actividades:

1. Planificar: Definimos las preguntas de investigación de este estudio y la consulta de búsqueda utilizada para los artículos.
2. Realizar: Durante esta fase buscamos 3 fuentes bibliográficas y descargamos todos los artículos en formato CSV. Luego de esto, se aplican criterios de inclusión y exclusión en tres fases para eliminar artículos hasta llegar al grupo final de artículos que incluyen en este estudio. Dos revisores leyeron cada artículo de forma independiente e hicieron un análisis extrayendo info. relevante para las preguntas de investigación que se desarrollaron en la fase de planificación. Se revisa y se resuelven las discrepancias para obtener el análisis final, se calcula el kappa de Cohen para evaluar la fiabilidad de nuestra evaluación. La puntuación que se obtiene es de 0.97 y esto indica que obtuvimos una concordancia casi perfecta en nuestro análisis.

3. Presentar: En esta fase, analizamos los datos y los organizamos para poder responder a las preguntas de investigación que planteamos.

4.1) Research Questions: A través de la revisión sistemática, pretendemos responder a las siguientes preguntas:

- ¿Cuáles son las representaciones de código más utilizadas?
- ¿Ciertas tareas de ciberseguridad utilizan solo o principalmente un tipo de representación del código fuente?
- ¿Qué tareas de ciberseguridad cubren las técnicas que se han creado?
- ¿Qué lenguajes de programación son los principales objetivos de las técnicas basadas en el aprendizaje automático para las tareas de ciberseguridad?
- ¿Qué modelos se utilizan habitualmente con diferentes representaciones?

4.2) Search Method: Se utilizó la siguiente cadena de búsqueda para encontrar estudios primarios relacionados con la representación o representaciones del código fuente para tareas de ciberseguridad basadas en ML.

- ("machine learning" OR "deep learning" OR "artificial intelligence") AND ("security" OR "vulnerability") AND ("code")

Aunque esta cadena de búsqueda suele ser muy general, arrojó un resultado de 67 512 artículos, lo que se decidió es, en lugar de utilizar una cadena específica que pudiera omitir una categoría de tareas o representaciones de seguridad de software, creamos una cadena general y elimináramos

manualmente cualquier artículo que no cumplieran nuestros criterios de evaluación (inclusión o exclusión)

Se busco en 3 base de datos:

- ACM library
- IEEE xplora
- Springer Link
- Tambien se busco en nueve conferencias de ML y NLP

4.3) Criterio de Inclusion y Exclusion: La tabla de abajo enumera los criterios de inclusion y exclusion aplicados a los articulos en multiples etapas para eliminar aquellos que no eran relevantes para el estudio.

Table 1. Inclusion and Exclusion Criteria

Inclusion Criteria	Exclusion Criteria
I1 Written between 2012 and May 2023	E1 Duplicated studies
I2 A full paper	E2 Books, reference work entries, reference works
I3 Focused on ML for cybersecurity tasks	E3 Position papers, short papers, tool demo papers, keynotes, reviews, tutorials, and panel discussions
I4 Contains information regarding the source code's representation	E4 Studies not in English
	E5 Survey/comparative studies.

4.4) Paper Selection:

- Comenzamos con un totla de 67512 artículos
- Excluimos los duplicados y los que no se encontraban dentro del rango de años comprendido entre el 2012 y mayo del 2023, asi como los articulos que no eran completos. Esto dio un resultado de 52 836 articulos
- Posteriormente revisamos titulo, palabras clave, y el resumen de cada articulo para incluirlos o excluirlos en funcion de si se ajustaban a nuestros criterios, tras esta busqueda nos quedamos con 520 articulos.
- A continuación aplicamos los mismo criterios esta vez leyendo el articulo completo y fruto de este resultado no quedo un total de 141 articulos que se incluyen en este survey.



Fig. 1. Overview of the Three Stages of our Search Process

4.5) Extraccion de Data: El objetivo principal era responder a las 5 preguntas de investigacion, para ello, de cada uno de los articulos, extrajeron sistematicamente 4 datos fundamentales:

- La representacion utilizada
- La tarea de ciberseguridad
- Los lenguajes de programacion
- El tipo de modelo de IA

¿Como manejan los detalles de la "representacion"?

Aqui los investigadores establecieron una regla muy importante para ser consistentes:

- Regla de "Un paper, multiples representaciones"

- Regla de la "representacion final"

¿Que hicieron ademas de responder las preguntas?

Los investigadores no se detuvieron ahí, dieron un paso mas allá para ofrecer un analisis mas completo:

- Extraccion de Metricas de Rendimiento

- El proposito de extraccion extra: (agregar y comparar)

5. RQ1 Results: What Are the Commonly Used Source Code Representations?:

Los autores identificaron 39 representaciones distintas, las cuales agruparon en 4 categorias principales:

- Basadas en arboles:
 - Árbol de Sintaxis Abstracta (AST - Abstract Syntax Tree):
 - Árbol de Análisis Sintáctico (Parse Tree):
 - AST+
- Basadas en grafos:
 - Grafo de Flujo de Control (CFG - Control Flow Graph):
 - Grafo de Dependencia de Programa (PDG - Program Dependence Gxraph)
 - Grafo de Flujo de Datos (DFG - Data Flow Graph)
 - Grafo de Llamadas (Call Graph)
 - Grafo de Propiedades de Código (CPG - Code Property Graph)
 - Otras representaciones de grafos
- Representaciones Lexicas:
 - Tokenizador (Tokenizer):
 - iSeVC y sSyVC
 - Fragmento de Contrato (Contract Snippet)
- Representaciones Miscelaneas
 - Imagen (Image)
 - Métricas de Código (Code Metrics)
 - Gadgets de Código (Code Gadgets)
 - Secuencias de Códigos de Operación (Opcode Sequences):
 - Expresión Regular (Regular Expression)
 - Información de la Aplicación (Application Information) y Llamadas a API (API Calls)

6. RQ2: Do Certain Tasks Only Use or Mostly Use One Type of Representation? Este pregunta analiza la relacion entre las diferentes representaciones de codigo fuente y las tareas especificas de ciberseguridad para las que se utilizan. La conclusión principal es que, si bien existen preferencias claras para ciertas tareas, la elección de una representación a menudo implica un equilibrio entre la robustez de la información que proporciona y los recursos computacionales necesarios para generarla y procesarla.

- Deteccion de vulnerabilidades: La tareas mas común, dado que es la tarea mas investigada, es tambien la que utiliza la mayor variedad de representaciones.

- Tokenizadores y Árboles de Sintaxis Abstracta (ASTs)
- Representaciones Basadas en Grafos
- Representaciones Específicas para Vulnerabilidades:
- Detección de malware:
 - Grafo de Flujo de Control Interprocedimental (ICFG), Grafo de Llamadas y CADG (Grafo de Dependencia de API Contextual):
- El equilibrio: Necesidades vs. Recursos
 - Un **Grafo de Propiedades de Código (CPG)** es extremadamente robusto, ya que combina AST, CFG y PDG, pero es muy complejo y costoso de generar para miles de muestras de código.
 - En contraste, si el objetivo es detectar un ataque de inyección, un **Grafo de Flujo de Datos (DFG)** podría ser suficiente, ya que proporciona la información necesaria sobre el flujo de datos sin ser excesivamente complejo o intensivo en recursos.

Hallazgos clave de la RQ2:

- La detección de vulnerabilidades es la tarea más popular
- Ciertas representaciones son altamente especializadas
- Hay una clara preferencia para la detección de malware

7. RQ3: What Cybersecurity Tasks Are Covered by the ML-Based Techniques? Para organizar y comprender el alcance de las tareas de ciberseguridad abordadas por las técnicas ML, los autores clasificaron las distintas tareas dentro del marco del ciclo de Proceso Unificado Racional, un proceso del desarrollo del software, de las nueve disciplinas del ciclo RUP, se observó que las tareas de ciberseguridad analizadas encajan en 5 categorías. La detección de vulnerabilidades es la tarea más popular, siendo esta con 75 artículos de los 141 en total

- Analisis y Diseño: Tareas que hay en esta fase.
 - Detección, predicción y clasificación de malware.
 - Analisis de seguridad, aplicado especialmente a contratos inteligentes
- Implementación
 - Reparación de vulnerabilidades
 - Predicción de vulnerabilidades
- Pruebas (Testing)
 - Detección de vulnerabilidades
 - Detección de mal uso de criptografía
 - Detección, clasificación, filtrado, desofuscación y localización de código malicioso
 - Detección de comportamiento malicioso
 - Detección de fugas de contraseña y ataques de inyección
 - Detección de ataques de reentrada
 - Analisis, clasificación, localización, y extrapolación de vulnerabilidades
 - Detección de clones de código vulnerable
 - Descubrimiento de vulnerabilidades API no protegidas

- Coincidencia de código fuente binario para tareas como detección de software.
 - Configuración y Gestión del cambio
 - Clasificación de commits de seguridad
 - Identificación de commits que introducen o corrigen vulnerabilidades
 - Detección de parches de seguridad
 - Entorno
 - Detección de paquetes maliciosos
- Importante:
- Detección vs. Localización: Las tareas de *detección* se centran en identificar de forma general que existe una característica (por ejemplo, un fragmento de código vulnerable), mientras que la *localización* se enfoca en señalar exactamente dónde ocurre esa característica (por ejemplo, la línea específica de la vulnerabilidad).
 - **Predicción vs. Detección:** Las técnicas de *predicción* buscan encontrar problemas antes de que se manifiesten en un sistema lanzado oficialmente. En cambio, las técnicas de *detección* se enfocan en identificar problemas que ya existen en un sistema en funcionamiento.

8. RQ4: What Programming Languages Are Predominantly Targeted by the ML-Based Techniques for Cybersecurity Tasks?

Table 3. Languages Supported by Existing Techniques

Lang.	#Papers	Lang.	#Papers	Lang.	#Papers	Lang.	#Papers	Lang.	#Papers	Lang.	#Papers
C	81 (57.4%)	JS	12 (8.5%)	Python	6 (4.3%)	C#	2 (1.4%)	Gecko	1 (0.7%)	Powershell	1 (0.7%)
C++	50 (35.5%)	Solidity	12 (8.5%)	CSS	3 (2.1%)	SQL	2 (1.4%)	Go	1 (0.7%)	Ruby	1 (0.7%)
Java	36 (25.5%)	PHP	8 (5.7%)	Rust	3 (2.1%)	TS	2 (1.4%)	HTML	1 (0.7%)	Smali	1 (0.7%)
										XUL	1 (0.7%)

JS = JavaScript; TS = TypeScript; SM = SpiderMonkey.

Esta sección analiza que LP son los más estudiados por las técnicas de ciberseguridad basadas en ML

- C y C++ son los más populares: C es el lenguaje más popular con un total de 81 artículos cubiertos (57.4%), seguido de C++ cubierto por 50 (35.5%). La popularidad de ambos lenguajes se debe a dos factores principales: la gran disponibilidad de conjunto de datos para tareas de ciberseguridad y el mayor riesgo inherente de vulnerabilidades relacionadas con la memoria en estos lenguajes.
- Enfoque en Java y Solidity: Se observa un notable enfoque en tareas de seguridad para aplicaciones de Android, donde los estudios se centraron exclusivamente en aplicaciones escritas en **Java**. Además, debido a la creciente popularidad de los contratos inteligentes, varias técnicas se desarrollaron para Solidity, el lenguaje utilizado para escribirlos.
- Una brecha sorprendente en Python y JS: A pesar de ser lenguajes extremadamente populares entre los desarrolladores, hay sorprendentemente pocas técnicas de seguridad para **Python** (solo 6) y **JavaScript** (solo 12). Los autores sugieren que esto podría deberse a la falta de conjuntos de datos disponibles para que los investigadores entrenen y prueben sus modelos. Esta es considerada una brecha que debería ser abordada en futuros trabajos
- Ausencia de herramientas agnósticas al lenguaje: Ninguna de las técnicas analizadas

es "agnóstica al lenguaje", lo que significa que cada herramienta fue creada para uno o, como mucho, unos pocos lenguajes específicos. Crear una herramienta que funcione para cualquier lenguaje es muy difícil debido a la gran variedad de paradigmas de programación existentes.

9. RQ5: What Models Are Commonly Used with Different Representations?
10. Threats to Validity
11. Final Considerations