



FACULTAD DE
INGENIERÍAS

Clasificación de Malware mediante Análisis de Imágenes con Deep Learning

Sistemas Inteligentes II
Detección y Clasificación de Familias de Malware
usando Redes Neuronales Convolucionales

Juan David Díaz Castaño

Juan Sebastián Henao Parra

Juan Carlos Saldarriaga Urrea

Over Haider Castrillon Valencia

Facultad de Ingeniería
Ingeniería en Sistemas y Computación
Sistemas Inteligentes II

Manizales, Caldas, Noviembre 22 de 2025



FACULTAD DE
INGENERÍAS

Clasificación de Malware mediante Análisis de Imágenes con Deep Learning

Juan David Díaz Castaño

Juan Sebastián Henao Parra

Juan Carlos Saldarriaga Urrea

Over Haider Castrillon Valencia

Docente: Jorge Alberto Jaramillo Garzón
Profesor, Universidad de Caldas

Facultad de Ingeniería
Ingeniería en Sistemas y Computación
Sistemas Inteligentes II

Manizales, Caldas, Noviembre 22 de 2025

Resumen

Este proyecto aborda el problema crítico de la detección y clasificación automatizada de malware mediante técnicas de Deep Learning. Con el aumento exponencial de amenazas ciberneticas, los métodos tradicionales de detección basados en firmas resultan insuficientes. Proponemos un enfoque basado en visión por computador que convierte ejecutables de malware en imágenes y emplea Redes Neuronales Convolucionales (CNN) para clasificarlos en familias.

La metodología utiliza tres datasets públicos reconocidos: *Blended Malware Image Dataset*, *Malevis Dataset*, y *MalImg Dataset*, que contienen muestras diversas de diferentes familias de malware. Se implementaron y compararon diversas arquitecturas CNN, incluyendo modelos preentrenados y arquitecturas personalizadas, para evaluar su capacidad de generalización y precisión en la clasificación.

Los resultados experimentales demuestran que el enfoque basado en imágenes logra una precisión de clasificación superior al XX % en la identificación de familias de malware, superando métodos baseline tradicionales. El análisis comparativo revela que las características visuales extraídas por las CNN capturan patrones estructurales distintivos de cada familia, permitiendo una detección robusta incluso ante variantes de malware previamente no vistas.

Esta investigación contribuye al campo de la ciberseguridad proporcionando un método escalable y eficiente para la clasificación automatizada de malware, con aplicaciones potenciales en sistemas de detección de intrusiones y análisis forense digital.

Palavras-Chave: Malware, Deep Learning, Redes Neuronales Convolucionales, Clasificación de Imágenes, Ciberseguridad, Análisis de Malware.

Abstract

This project addresses the critical problem of automated malware detection and classification using Deep Learning techniques. With the exponential increase in cyber threats, traditional signature-based detection methods prove insufficient. We propose a computer vision-based approach that converts malware executables into images and employs Convolutional Neural Networks (CNN) to classify them into families.

The methodology uses three recognized public datasets: *Blended Malware Image Dataset*, *Malevis Dataset*, and *MalImg Dataset*, containing diverse samples from different malware families. Various CNN architectures were implemented and compared, including pre-trained models and custom architectures, to evaluate their generalization capability and classification accuracy.

Experimental results demonstrate that the image-based approach achieves classification accuracy exceeding XX % in identifying malware families, outperforming traditional baseline methods. Comparative analysis reveals that visual features extracted by CNNs capture distinctive structural patterns of each family, enabling robust detection even against previously unseen malware variants.

This research contributes to the cybersecurity field by providing a scalable and efficient method for automated malware classification, with potential applications in intrusion detection systems and digital forensic analysis.

Keywords: Malware, Deep Learning, Convolutional Neural Networks, Image Classification, Cybersecurity, Malware Analysis.

Índice general

| | |
|---|----------|
| <i>Índice de figuras</i> | XIII |
| <i>Índice de cuadros</i> | XVI |
| <i>Siglas</i> | XIX |
| 1. Introducción | 1 |
| 1.1. Contexto y Motivación | 1 |
| 1.2. Problema de Investigación | 1 |
| 1.3. Pregunta de Investigación | 2 |
| 1.4. Objetivos | 2 |
| 1.4.1. Objetivo General | 2 |
| 1.4.2. Objetivos Específicos | 2 |
| 1.5. Justificación | 3 |
| 1.6. Alcance y Limitaciones | 3 |
| 1.6.1. Alcance | 3 |
| 1.6.2. Limitaciones | 3 |
| 1.7. Estructura del Documento | 3 |
| 2. Trabajo Relacionado | 5 |
| 2.1. Métodos Tradicionales de Detección de Malware | 5 |
| 2.1.1. Detección Basada en Firmas | 5 |
| 2.1.2. Análisis Heurístico | 5 |
| 2.1.3. Análisis Dinámico (Sandboxing) | 5 |
| 2.2. Machine Learning en Detección de Malware | 6 |
| 2.2.1. Enfoques Tradicionales de ML | 6 |
| 2.3. Deep Learning para Análisis de Malware | 6 |
| 2.3.1. Ventajas del Deep Learning | 6 |
| 2.3.2. Trabajos Previos con Deep Learning | 6 |
| 2.4. Análisis Visual de Malware | 7 |
| 2.4.1. Conversión de Binarios a Imágenes | 7 |
| 2.4.2. Características Visuales Discriminativas | 7 |
| 2.5. Trabajos con CNN para Clasificación de Malware | 7 |
| 2.5.1. Arquitecturas CNN Aplicadas | 7 |
| 2.5.2. Resultados Reportados en Literatura | 7 |

| | |
|---|-----------|
| 2.6. Datasets Pùblicos de Malware | 8 |
| 2.6.1. Descripción de Datasets Relevantes | 8 |
| 2.7. Brechas en la Literatura y Motivación del Proyecto | 8 |
| 2.8. Resumen | 8 |
| 3. Metodología | 10 |
| 3.1. Visión General del Proceso | 10 |
| 3.2. Datasets de Malware | 10 |
| 3.2.1. Descripción de Datasets Utilizados | 10 |
| 3.2.2. Distribución de Familias de Malware | 11 |
| 3.3. Preprocesamiento de Datos | 11 |
| 3.3.1. Conversión de Ejecutables a Imágenes | 11 |
| 3.3.2. Normalización y Redimensionamiento | 12 |
| 3.3.3. Aumento de Datos (Data Augmentation) | 12 |
| 3.4. División de Datos | 12 |
| 3.5. Arquitecturas de Redes Neuronales Convolucionales | 13 |
| 3.5.1. Arquitectura CNN Personalizada | 13 |
| 3.5.2. Modelos Preentrenados con Transfer Learning | 13 |
| 3.6. Configuración de Entrenamiento | 14 |
| 3.6.1. Hiperparámetros | 14 |
| 3.6.2. Estrategias de Optimización | 14 |
| 3.6.3. Manejo de Desbalance de Clases | 15 |
| 3.7. Entorno Experimental | 15 |
| 3.7.1. Hardware y Software | 15 |
| 3.7.2. Implementación | 15 |
| 3.8. Métricas de Evaluación | 16 |
| 3.8.1. Métricas Principales | 16 |
| 3.8.2. Métricas Multi-Clase | 16 |
| 3.8.3. Matriz de Confusión | 16 |
| 3.8.4. Curvas ROC y AUC | 16 |
| 3.9. Validación Cruzada y Análisis de Robustez | 16 |
| 3.10. Consideraciones Éticas y de Seguridad | 17 |
| 3.11. Resumen de la Metodología | 17 |
| 4. Experimentos y Resultados | 18 |
| 4.1. Configuración Experimental | 18 |
| 4.1.1. Resumen de Experimentos Realizados | 18 |
| 4.1.2. Entorno de Ejecución | 18 |
| 4.2. Análisis Exploratorio de Datos | 19 |
| 4.2.1. Distribución de Familias de Malware | 19 |
| 4.2.2. Visualización de Muestras | 19 |
| 4.3. Resultados de Clasificación | 19 |

| | |
|---|-----------|
| 4.3.1. Rendimiento de Arquitecturas en MalImg Dataset | 19 |
| 4.3.2. Rendimiento en Malevis Dataset | 19 |
| 4.3.3. Rendimiento en Blended Malware Dataset | 20 |
| 4.4. Análisis Comparativo | 20 |
| 4.4.1. Comparación entre Arquitecturas | 20 |
| 4.4.2. Análisis por Dataset | 20 |
| 4.5. Matrices de Confusión | 21 |
| 4.5.1. Análisis de Errores de Clasificación | 21 |
| 4.6. Curvas de Entrenamiento | 21 |
| 4.6.1. Evolución de Loss y Accuracy | 21 |
| 4.7. Análisis de Generalización | 21 |
| 4.7.1. Entrenamiento y Evaluación Cross-Dataset | 21 |
| 4.8. Estudio de Hiperparámetros | 22 |
| 4.8.1. Impacto del Learning Rate | 22 |
| 4.8.2. Efecto del Batch Size | 22 |
| 4.8.3. Técnicas de Regularización | 22 |
| 4.9. Análisis de Características Aprendidas | 22 |
| 4.9.1. Visualización de Mapas de Activación | 22 |
| 4.9.2. Reducción Dimensional y Clustering | 23 |
| 4.10. Comparación con Estado del Arte | 23 |
| 4.10.1. Comparación con Trabajos Previos | 23 |
| 4.11. Análisis de Tiempo de Inferencia | 23 |
| 4.11.1. Eficiencia Computacional | 23 |
| 4.12. Limitaciones y Desafíos Observados | 23 |
| 4.12.1. Desbalance de Clases | 24 |
| 4.12.2. Familias Similares | 24 |
| 4.12.3. Dependencia de Tamaño de Imagen | 24 |
| 4.12.4. Limitaciones de Generalización | 24 |
| 4.13. Resumen de Resultados | 24 |
| 5. Conclusiones y Trabajo Futuro | 25 |
| 5.1. Resumen del Proyecto | 25 |
| 5.2. Principales Hallazgos | 25 |
| 5.2.1. Viabilidad del Enfoque Visual | 25 |
| 5.2.2. Rendimiento de Arquitecturas | 26 |
| 5.2.3. Generalización y Robustez | 26 |
| 5.2.4. Características Discriminativas | 26 |
| 5.3. Respuesta a las Preguntas de Investigación | 26 |
| 5.3.1. Pregunta Principal | 26 |
| 5.3.2. Preguntas Secundarias | 26 |
| 5.4. Contribuciones del Proyecto | 27 |
| 5.4.1. Contribuciones Técnicas | 27 |

| | |
|--|----|
| 5.4.2. Contribuciones Prácticas | 27 |
| 5.5. Limitaciones del Estudio | 27 |
| 5.5.1. Limitaciones de Datasets | 27 |
| 5.5.2. Limitaciones Metodológicas | 28 |
| 5.5.3. Limitaciones de Interpretabilidad | 28 |
| 5.6. Trabajo Futuro | 28 |
| 5.6.1. Extensiones Metodológicas | 28 |
| 5.6.2. Mejoras en Robustez | 29 |
| 5.6.3. Ampliación de Datasets | 29 |
| 5.6.4. Aplicaciones Prácticas | 29 |
| 5.6.5. Investigación en Interpretabilidad | 30 |
| 5.6.6. Benchmarking y Estandarización | 30 |
| 5.7. Implicaciones para la Ciberseguridad | 30 |
| 5.7.1. Para Profesionales de Seguridad | 31 |
| 5.7.2. Para Desarrolladores de Soluciones de Seguridad | 31 |
| 5.7.3. Para la Investigación Académica | 31 |
| 5.8. Lecciones Aprendidas | 31 |
| 5.8.1. Técnicas | 31 |
| 5.8.2. Prácticas | 32 |
| 5.9. Consideraciones Éticas | 32 |
| 5.9.1. Uso Responsable | 32 |
| 5.9.2. Transparencia | 32 |
| 5.9.3. Privacidad y Confidencialidad | 32 |
| 5.10. Reflexiones Finales | 33 |
| 5.11. Conclusión | 33 |

Índice de figuras

| | |
|---|----|
| 4.1. Ejemplos de imágenes de malware de diferentes familias. Cada fila corresponde a una familia distinta, mostrando la variabilidad intra-clase y las diferencias inter-clase. | 19 |
| 4.2. Comparación de accuracy entre diferentes arquitecturas CNN evaluadas en los tres datasets. | 20 |
| 4.3. Matriz de confusión del modelo [nombre del mejor modelo] en [dataset]. Los valores en la diagonal representan clasificaciones correctas. | 21 |
| 4.4. Curvas de entrenamiento y validación para el modelo [nombre]. (a) Loss function. (b) Accuracy. | 21 |
| 4.5. Mapas de activación mostrando regiones de interés para la CNN al clasificar diferentes familias de malware. Las zonas más brillantes indican mayor importancia para la decisión del modelo. | 22 |
| 4.6. Visualización t-SNE de las representaciones aprendidas. Cada punto representa una muestra de malware, coloreada por su familia. La formación de clusters indica buena separabilidad aprendida por el modelo. | 23 |

Índice de cuadros

| | |
|--|----|
| 4.1. Distribución de muestras por familia en los tres datasets | 19 |
| 4.2. Rendimiento de diferentes arquitecturas CNN en MalImg Dataset | 19 |
| 4.3. Rendimiento en Malevis Dataset | 20 |
| 4.4. Rendimiento en Blended Malware Dataset | 20 |
| 4.5. Accuracy de generalización cross-dataset | 21 |
| 4.6. Impacto del learning rate en el rendimiento final | 22 |
| 4.7. Impacto del batch size | 22 |
| 4.8. Comparación con estado del arte en MalImg Dataset | 23 |
| 4.9. Tiempo de inferencia promedio | 23 |

1

Introducción

1.1. Contexto y Motivación

El panorama de la ciberseguridad contemporánea enfrenta desafíos sin precedentes. El incremento exponencial en el volumen y sofisticación de software malicioso (malware) representa una amenaza crítica para sistemas informáticos, infraestructuras críticas, y la seguridad digital de organizaciones e individuos. Según reportes recientes de la industria, se detectan millones de nuevas variantes de malware cada año, con atacantes desarrollando técnicas cada vez más evasivas que desafían los métodos tradicionales de detección.

Los sistemas antivirus convencionales se basan principalmente en firmas estáticas y análisis heurístico, métodos que resultan insuficientes ante el polimorfismo y la ofuscación empleada por malware moderno. Esta limitación motiva la búsqueda de enfoques innovadores que puedan adaptarse dinámicamente a nuevas amenazas sin depender exclusivamente de bases de datos de firmas conocidas.

1.2. Problema de Investigación

La detección y clasificación efectiva de malware presenta múltiples desafíos técnicos:

- **Volumen y velocidad:** La generación masiva de nuevas variantes de malware supera la capacidad de análisis manual.
- **Polimorfismo y ofuscación:** Las técnicas de evasión modifican el código sin alterar su funcionalidad maliciosa, evadiendo detección basada en firmas.
- **Variabilidad intra-familia:** Múltiples variantes dentro de una misma familia de malware pueden presentar diferencias significativas en su código.
- **Costo computacional:** El análisis dinámico mediante sandboxing requiere recursos significativos y tiempo de ejecución.

Estos desafíos plantean la necesidad de desarrollar métodos automatizados, eficientes y robustos capaces de identificar y clasificar malware con alta precisión, incluso ante

muestras previamente no vistas.

1.3. Pregunta de Investigación

Este proyecto busca responder la siguiente pregunta central:

¿Es posible clasificar eficazmente muestras de malware en sus respectivas familias mediante el análisis de sus representaciones visuales utilizando técnicas de Deep Learning, específicamente Redes Neuronales Convolucionales?

Preguntas secundarias incluyen:

- ¿Qué arquitecturas CNN son más efectivas para la clasificación de imágenes de malware?
- ¿Cómo se compara el rendimiento de modelos preentrenados versus arquitecturas diseñadas específicamente?
- ¿Qué características visuales de los ejecutables son más discriminativas para la clasificación?

1.4. Objetivos

1.4.1. Objetivo General

Desarrollar e implementar un sistema de clasificación de malware basado en Deep Learning que utilice representaciones visuales de ejecutables para identificar automáticamente familias de malware con alta precisión y eficiencia.

1.4.2. Objetivos Específicos

1. Recopilar y preprocessar tres datasets públicos de malware (Blended Malware Image Dataset, Malevis Dataset, MalImg Dataset) para entrenamiento y evaluación.
2. Implementar el pipeline de conversión de ejecutables a representaciones visuales adecuadas para análisis mediante CNN.
3. Diseñar, entrenar y evaluar múltiples arquitecturas de Redes Neuronales Convolucionales para la tarea de clasificación multi-clase de familias de malware.
4. Realizar análisis comparativo de diferentes arquitecturas, hiperparámetros y técnicas de regularización para optimizar el rendimiento del modelo.
5. Evaluar el rendimiento del sistema propuesto mediante métricas estándar (precisión, recall, F1-score) y comparar con métodos baseline existentes.
6. Analizar e interpretar las características visuales aprendidas por el modelo para comprender qué patrones estructurales distinguen diferentes familias de malware.

1.5. Justificación

La adopción de técnicas de Deep Learning para análisis de malware se justifica por múltiples razones:

Capacidad de aprendizaje automático de características: A diferencia de métodos tradicionales que requieren ingeniería manual de características, las CNN aprenden automáticamente representaciones jerárquicas discriminativas directamente de los datos crudos.

Escalabilidad: Una vez entrenado, el modelo puede clasificar nuevas muestras en tiempo prácticamente real, permitiendo procesar grandes volúmenes de datos.

Robustez ante variaciones: Las características visuales capturadas por CNN pueden ser invariantes a ciertas técnicas de ofuscación que alteran el código pero preservan estructuras fundamentales.

Transferibilidad: Los modelos entrenados en ciertos datasets pueden adaptarse (fine-tuning) a nuevos conjuntos de datos con menor costo computacional.

Aplicabilidad práctica: El enfoque propuesto puede integrarse en sistemas reales de detección de amenazas, análisis forense digital, y respuesta a incidentes de seguridad.

1.6. Alcance y Limitaciones

1.6.1. Alcance

Este proyecto se enfoca específicamente en:

- Clasificación de familias de malware conocidas presentes en los datasets seleccionados
- Análisis estático mediante representaciones visuales (sin ejecución dinámica)
- Evaluación en entorno controlado con muestras etiquetadas
- Arquitecturas CNN estándar y variantes preentrenadas

1.6.2. Limitaciones

Las principales limitaciones incluyen:

- Dependencia de datasets públicos con distribución potencialmente diferente a amenazas en entornos reales
- Limitación a familias de malware presentes en los datos de entrenamiento (detección de zero-day requeriría enfoques adicionales)
- Foco en malware de Windows (limitado por los datasets disponibles)
- No considera análisis de comportamiento dinámico ni técnicas híbridas

1.7. Estructura del Documento

El resto de este documento se organiza de la siguiente manera:

Capítulo 2 – Trabajo Relacionado: Revisión del estado del arte en detección de malware, aplicaciones de Deep Learning en ciberseguridad, y enfoques basados en análisis visual.

Capítulo 3 – Metodología: Descripción detallada de los datasets utilizados, proceso de preprocesamiento, arquitecturas CNN implementadas, y configuración experimental.

Capítulo 4 – Experimentos y Resultados: Presentación de los resultados experimentales, análisis comparativo de diferentes modelos, y evaluación del rendimiento.

Capítulo 5 – Conclusiones y Trabajo Futuro: Síntesis de los hallazgos principales, contribuciones del proyecto, limitaciones encontradas, y direcciones futuras de investigación.

2

Trabajo Relacionado

Este capítulo presenta una revisión del estado del arte en detección y clasificación de malware, con énfasis en enfoques basados en aprendizaje automático y Deep Learning. Se analizan los métodos tradicionales, las técnicas emergentes basadas en visión por computador, y los trabajos previos que fundamentan la propuesta de este proyecto.

2.1. Métodos Tradicionales de Detección de Malware

2.1.1. Detección Basada en Firmas

Los sistemas antivirus tradicionales emplean detección basada en firmas, donde se comparan patrones binarios conocidos (hashes o secuencias de bytes) contra bases de datos de malware identificado. Este enfoque presenta alta precisión para amenazas conocidas pero falla ante:

- **Polimorfismo:** Malware que modifica su código en cada infección
- **Ofuscación:** Técnicas de empaquetado que alteran la firma sin cambiar funcionalidad
- **Zero-day:** Amenazas completamente nuevas sin firma registrada

2.1.2. Análisis Heurístico

Los métodos heurísticos evalúan comportamientos sospechosos mediante reglas predefinidas, como acceso no autorizado al registro, modificación de archivos del sistema, o comunicaciones de red anómalas. Aunque más flexibles que las firmas, dependen de la experiencia humana para definir reglas y pueden generar altas tasas de falsos positivos.

2.1.3. Análisis Dinámico (Sandboxing)

El análisis dinámico ejecuta muestras sospechosas en entornos controlados (sandboxes) para observar su comportamiento. Proporciona información detallada pero presenta desventajas:

- Alto costo computacional y temporal
- Malware puede detectar entornos virtualizados y alterar su comportamiento
- Dificultad para escalar a análisis masivo

2.2. Machine Learning en Detección de Malware

2.2.1. Enfoques Tradicionales de ML

La aplicación de técnicas de Machine Learning a la clasificación de malware ha sido ampliamente estudiada. Los enfoques clásicos incluyen:

Support Vector Machines (SVM): Utilizadas para clasificación binaria (benigno vs. malicioso) y multi-clase (familias de malware) basándose en características extraídas manualmente como n-gramas de bytes, llamadas API, o estructuras PE.

Random Forests y Decision Trees: Algoritmos ensemble que combinan múltiples árboles de decisión, efectivos para capturar relaciones no lineales entre características.

k-Nearest Neighbors (k-NN): Clasificación basada en similitud con muestras conocidas, útil para detección de variantes de familias existentes.

La principal limitación de estos métodos radica en la dependencia de ingeniería manual de características (feature engineering), proceso que requiere expertise de dominio y puede no capturar representaciones óptimas.

2.3. Deep Learning para Análisis de Malware

2.3.1. Ventajas del Deep Learning

El Deep Learning ofrece capacidad de aprendizaje automático de representaciones jerárquicas directamente de datos crudos o mínimamente procesados. Las principales arquitecturas aplicadas incluyen:

Redes Neuronales Convolucionales (CNN): Especialmente efectivas cuando el malware se representa como imágenes o secuencias con estructura espacial.

Redes Neuronales Recurrentes (RNN/LSTM): Utilizadas para analizar secuencias de instrucciones o llamadas API, capturando dependencias temporales.

Autoencoders: Empleados para detección de anomalías, identificando muestras que se desvían significativamente de patrones normales.

2.3.2. Trabajos Previos con Deep Learning

Diversos estudios han explorado Deep Learning para malware:

- Análisis de secuencias de bytes mediante CNN 1D
- Clasificación de malware de Android mediante análisis de permisos y grafos de llamadas
- Detección de ransomware usando autoencoders variacionales
- Clasificación de familias mediante embeddings aprendidos

2.4. Análisis Visual de Malware

2.4.1. Conversión de Binarios a Imágenes

Un enfoque innovador consiste en visualizar ejecutables como imágenes. El proceso típico incluye:

1. Leer el archivo binario como secuencia de bytes
2. Interpretar cada byte como valor de intensidad de píxel (0-255)
3. Organizar los píxeles en una matriz bidimensional
4. Aplicar redimensionamiento a tamaño estándar

Esta representación visual preserva la estructura del ejecutable y permite aplicar técnicas de visión por computador.

2.4.2. Características Visuales Discriminativas

Estudios han demostrado que diferentes familias de malware exhiben texturas visuales distintivas relacionadas con:

- Estructura del código (secciones .text, .data, .rdata)
- Técnicas de empaquetado y compresión
- Presencia de recursos embebidos (imágenes, configuraciones)
- Patrones de cifrado o ofuscación

Las CNN son particularmente efectivas para extraer automáticamente estas características visuales.

2.5. Trabajos con CNN para Clasificación de Malware

2.5.1. Arquitecturas CNN Aplicadas

Diversos trabajos han aplicado CNN para clasificación de malware visual:

CNN Shallow: Arquitecturas simples de pocas capas para datasets pequeños, reduciendo riesgo de overfitting.

CNN Profundas Personalizadas: Arquitecturas diseñadas específicamente para características de imágenes de malware.

Transfer Learning: Uso de modelos preentrenados (VGG, ResNet, Inception) en ImageNet, adaptados mediante fine-tuning para malware.

2.5.2. Resultados Reportados en Literatura

Los estudios existentes reportan precisiones que varían entre 85 % y 99 % dependiendo de:

- Complejidad del dataset (número de familias)

- Calidad y cantidad de muestras de entrenamiento
- Arquitectura CNN utilizada
- Técnicas de regularización y aumento de datos

2.6. Datasets Públicos de Malware

2.6.1. Descripción de Datasets Relevantes

Los principales datasets públicos para investigación en clasificación de malware incluyen:

MalImg Dataset: Contiene imágenes de malware organizadas en 25 familias, ampliamente utilizado como benchmark estándar.

Malevis Dataset: Dataset más reciente con mayor variedad de familias y muestras, incluyendo malware moderno.

Blended Malware Image Dataset: Combinación de múltiples fuentes, ofreciendo diversidad en tipos y familias.

Estos datasets proporcionan la base empírica para entrenar y evaluar modelos de clasificación.

2.7. Brechas en la Literatura y Motivación del Proyecto

A pesar de los avances, persisten oportunidades de investigación:

- **Comparación sistemática:** Pocos estudios comparan múltiples datasets y arquitecturas bajo condiciones controladas.
- **Interpretabilidad:** Limitada comprensión de qué características visuales específicas utiliza la CNN para discriminar familias.
- **Robustez:** Evaluación insuficiente ante técnicas adversariales diseñadas para engañar clasificadores.
- **Escalabilidad:** Necesidad de validar rendimiento en datasets masivos más representativos de entornos reales.

Este proyecto busca contribuir mediante una evaluación exhaustiva de múltiples arquitecturas CNN sobre tres datasets reconocidos, con análisis detallado de rendimiento y características aprendidas.

2.8. Resumen

La revisión del estado del arte revela que:

1. Los métodos tradicionales de detección presentan limitaciones significativas ante malware moderno.
2. El Machine Learning y Deep Learning han demostrado gran potencial para clasificación automatizada.

3. El enfoque visual de malware permite aplicar técnicas probadas de visión por computador.
4. Existen datasets públicos adecuados para entrenar y evaluar modelos.
5. Persisten oportunidades para investigación adicional en comparación de arquitecturas e interpretabilidad.

Estos fundamentos justifican el enfoque propuesto en este proyecto, combinando análisis visual con CNN para clasificación robusta y eficiente de familias de malware.

3

Metodología

Este capítulo describe detalladamente la metodología empleada para desarrollar el sistema de clasificación de malware basado en Deep Learning. Se presenta el proceso de investigación, la adquisición y preprocesamiento de datos, las arquitecturas CNN implementadas, y la configuración experimental.

3.1. Visión General del Proceso

El proceso metodológico sigue un pipeline estructurado en las siguientes etapas:

1. **Adquisición de datasets:** Descarga y organización de tres datasets públicos de malware
2. **Preprocesamiento:** Conversión de ejecutables a imágenes y normalización
3. **Análisis exploratorio:** Visualización y estadísticas de las familias de malware
4. **Diseño de arquitecturas:** Implementación de modelos CNN
5. **Entrenamiento:** Configuración de hiperparámetros y optimización
6. **Evaluación:** Medición de rendimiento mediante métricas estándar
7. **Análisis comparativo:** Comparación entre diferentes modelos y datasets

3.2. Datasets de Malware

3.2.1. Descripción de Datasets Utilizados

Este proyecto emplea tres datasets públicos reconocidos en la comunidad de investigación:

MalImg Dataset

- **Fuente:** Disponible en Kaggle ([manaswinisunkari/malimg-dataset90](https://www.kaggle.com/manaswinisunkari/malimg-dataset90))
- **Composición:** Aproximadamente 9,000+ muestras de malware
- **Familias:** 25 familias diferentes de malware de Windows

- **Formato:** Imágenes en escala de grises derivadas de ejecutables binarios
- **Uso:** Dataset benchmark ampliamente citado en literatura académica

Malevis Dataset

- **Fuente:** Dataset público especializado en visualización de malware
- **Composición:** Colección moderna de muestras de malware recientes
- **Características:** Incluye variantes más contemporáneas de amenazas
- **Formato:** Representaciones visuales estandarizadas

Blended Malware Image Dataset

- **Fuente:** Combinación de múltiples repositorios públicos
- **Composición:** Dataset híbrido con mayor diversidad de tipos
- **Ventaja:** Combina muestras de diferentes fuentes y períodos temporales
- **Objetivo:** Evaluar generalización del modelo ante diversidad de datos

3.2.2. Distribución de Familias de Malware

Los datasets contienen diversas familias de malware, incluyendo:

- **Trojanos:** Alureon, VB, Agent, etc.
- **Gusanos:** Kelihos, Autorun, Worm
- **Backdoors:** Rbot, Bifrose
- **Ransomware:** Locker, Cryptolocker
- **Adware/Spyware:** Winwebsec, FakeRean

La distribución de clases presenta cierto desbalance, aspecto considerado en la estrategia de entrenamiento mediante técnicas de balanceo y ponderación.

3.3. Preprocesamiento de Datos

3.3.1. Conversión de Ejecutables a Imágenes

El proceso de conversión de binarios a imágenes se implementa mediante los siguientes pasos:

1. **Lectura binaria:** El archivo ejecutable se lee como secuencia de bytes (valores 0-255)
2. **Mapeo a píxeles:** Cada byte se interpreta como intensidad de píxel en escala de grises
3. **Determinación de dimensiones:**
 - Se calcula el tamaño total del archivo en bytes
 - Se determina el ancho óptimo de la imagen (típicamente potencia de 2)
 - La altura se calcula como: $altura = \lceil \frac{tamaño_bytes}{ancho} \rceil$

4. **Construcción de matriz:** Los bytes se organizan en matriz 2D con las dimensiones calculadas
5. **Padding:** Si es necesario, se añaden píxeles de relleno al final para completar la última fila

3.3.2. Normalización y Redimensionamiento

Para garantizar uniformidad en el entrenamiento:

- **Redimensionamiento:** Todas las imágenes se redimensionan a tamaño fijo (e.g., 224×224 o 256×256) mediante interpolación bilineal
- **Normalización de píxeles:** Los valores de píxeles se escalan al rango $[0, 1]$ dividiendo por 255
- **Estandarización:** Opcionalmente, se aplica normalización por canal usando media y desviación estándar del dataset de entrenamiento

3.3.3. Aumento de Datos (Data Augmentation)

Para mejorar la capacidad de generalización y mitigar overfitting, se aplican técnicas de aumento de datos durante el entrenamiento:

- Rotaciones leves ($\pm 5\text{-}10$ grados)
- Recortes aleatorios (random crops)
- Volteos horizontales/verticales
- Ajustes de brillo y contraste
- Ruido gaussiano moderado

3.4. División de Datos

Los datasets se dividen siguiendo una estrategia de partición estratificada para mantener proporciones de clases:

- **Conjunto de Entrenamiento:** 70 % de las muestras (usado para optimización de pesos)
- **Conjunto de Validación:** 15 % de las muestras (usado para ajuste de hiperparámetros y detección de overfitting)
- **Conjunto de Prueba:** 15 % de las muestras (usado exclusivamente para evaluación final)

Se garantiza que las muestras de prueba no sean vistas durante entrenamiento ni validación, asegurando evaluación imparcial del rendimiento.

3.5. Arquitecturas de Redes Neuronales Convolucionales

3.5.1. Arquitectura CNN Personalizada

Se diseñó una arquitectura CNN baseline adaptada a las características de imágenes de malware:

Estructura de capas:

1. Bloque Convolucional 1:

- Conv2D: 32 filtros, kernel 3×3 , activación ReLU
- Conv2D: 32 filtros, kernel 3×3 , activación ReLU
- MaxPooling2D: pool size 2×2
- Dropout: 25 %

2. Bloque Convolucional 2:

- Conv2D: 64 filtros, kernel 3×3 , activación ReLU
- Conv2D: 64 filtros, kernel 3×3 , activación ReLU
- MaxPooling2D: pool size 2×2
- Dropout: 25 %

3. Bloque Convolucional 3:

- Conv2D: 128 filtros, kernel 3×3 , activación ReLU
- Conv2D: 128 filtros, kernel 3×3 , activación ReLU
- MaxPooling2D: pool size 2×2
- Dropout: 25 %

4. Capas Fully Connected:

- Flatten
- Dense: 512 unidades, activación ReLU
- Dropout: 50 %
- Dense: 256 unidades, activación ReLU
- Dropout: 50 %
- Dense: N unidades (número de clases), activación Softmax

3.5.2. Modelos Preentrenados con Transfer Learning

Se experimentó con arquitecturas preentrenadas en ImageNet, adaptadas mediante fine-tuning:

VGG16/VGG19

- Arquitectura profunda con bloques convolucionales uniformes
- Carga de pesos preentrenados (ImageNet)
- Congelamiento de capas base (feature extraction)
- Reemplazo de capas clasificadoras superiores
- Fine-tuning selectivo de últimas capas convolucionales

ResNet50/ResNet101

- Arquitectura con conexiones residuales
- Manejo efectivo de gradientes en redes profundas
- Transferencia de características de bajo y alto nivel
- Adaptación de capa de clasificación final

InceptionV3

- Módulos Inception con convoluciones multi-escala
- Captura de patrones a diferentes escalas espaciales
- Eficiencia computacional mediante factorización de convoluciones
- Regularización implícita por arquitectura

3.6. Configuración de Entrenamiento

3.6.1. Hiperparámetros

Los hiperparámetros principales empleados incluyen:

- **Función de pérdida:** Categorical Cross-Entropy (para clasificación multi-clase)
- **Optimizador:** Adam (Adaptive Moment Estimation)
 - Learning rate inicial: 0.001
 - $\beta_1 = 0,9$, $\beta_2 = 0,999$
 - $\epsilon = 10^{-7}$
- **Batch size:** 32 o 64 (según disponibilidad de memoria GPU)
- **Épocas:** Hasta 100 con early stopping
- **Regularización:**
 - Dropout: 0.25 (capas convolucionales), 0.5 (capas densas)
 - L2 regularization: $\lambda = 0,0001$

3.6.2. Estrategias de Optimización

Learning Rate Scheduling:

- ReduceLROnPlateau: Reduce learning rate cuando la pérdida de validación se estanca
- Factor de reducción: 0.5
- Paciencia: 5 épocas

Early Stopping:

- Monitoreo de pérdida de validación
- Paciencia: 10-15 épocas
- Restauración de mejores pesos

Model Checkpointing:

- Guardado del modelo con mejor rendimiento en validación
- Métrica de monitoreo: Accuracy o F1-score

3.6.3. Manejo de Desbalance de Clases

Para abordar el desbalance en la distribución de clases:

- **Class Weighting:** Asignación de pesos inversamente proporcionales a la frecuencia de clase
- **Stratified Sampling:** Muestreo estratificado en división de datos
- **Focal Loss:** (opcional) Función de pérdida que enfatiza ejemplos difíciles de clasificar

3.7. Entorno Experimental

3.7.1. Hardware y Software

Hardware:

- CPU: [Especificar procesador]
- GPU: [Especificar GPU si disponible, ej. NVIDIA GTX/RTX]
- RAM: [Especificar cantidad]
- Almacenamiento: [Especificar tipo y capacidad]

Software:

- Sistema Operativo: [Linux/Windows/macOS]
- Python: 3.8+
- Framework de Deep Learning: TensorFlow 2.x / PyTorch 1.x
- Bibliotecas adicionales:
 - NumPy, Pandas (manipulación de datos)
 - Matplotlib, Seaborn (visualización)
 - Scikit-learn (métricas y preprocessamiento)
 - OpenCV / Pillow (procesamiento de imágenes)

3.7.2. Implementación

El código del proyecto se organizó de la siguiente manera:

- **data/**: Scripts de descarga y preprocessamiento
- **models/**: Definiciones de arquitecturas CNN
- **training/**: Scripts de entrenamiento y callbacks
- **evaluation/**: Evaluación de modelos y generación de métricas
- **notebooks/**: Jupyter notebooks para análisis exploratorio
- **utils/**: Funciones auxiliares (visualización, logging)

3.8. Métricas de Evaluación

El rendimiento de los modelos se evalúa mediante las siguientes métricas:

3.8.1. Métricas Principales

Accuracy (Exactitud):

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Precision (Precisión):

$$\text{Precision} = \frac{TP}{TP + FP}$$

Recall (Sensibilidad o Exhaustividad):

$$\text{Recall} = \frac{TP}{TP + FN}$$

F1-Score:

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Donde TP = True Positives, TN = True Negatives, FP = False Positives, FN = False Negatives.

3.8.2. Métricas Multi-Clase

Para clasificación multi-clase, se calculan promedios:

- **Macro-Average:** Promedio simple sobre todas las clases (igual peso por clase)
- **Weighted-Average:** Promedio ponderado por número de muestras por clase

3.8.3. Matriz de Confusión

Se genera matriz de confusión para visualizar patrones de clasificación errónea y confusión entre familias similares.

3.8.4. Curvas ROC y AUC

Para análisis detallado, se calculan curvas ROC (Receiver Operating Characteristic) y métricas AUC (Area Under Curve) utilizando estrategia One-vs-Rest.

3.9. Validación Cruzada y Análisis de Robustez

Se implementa validación cruzada estratificada (k -fold) para evaluar la estabilidad del modelo y reducir varianza en las estimaciones de rendimiento.

Adicionalmente, se analiza:

- **Generalización inter-dataset:** Entrenamiento en un dataset y evaluación en otro
- **Análisis de sensibilidad:** Variación de hiperparámetros y medición de impacto
- **Estudio de ablación:** Evaluación del aporte individual de componentes del modelo

3.10. Consideraciones Éticas y de Seguridad

Este proyecto maneja muestras reales de malware, lo que requiere precauciones:

- **Entorno aislado:** Todos los experimentos se realizan en máquinas virtuales aisladas sin acceso a red
- **Datasets públicos:** Se utilizan exclusivamente datasets públicos con fines de investigación
- **No ejecución:** El análisis es completamente estático, sin ejecutar binarios maliciosos
- **Almacenamiento seguro:** Los datasets se almacenan en particiones cifradas
- **Uso responsable:** Los modelos entrenados se emplean exclusivamente con fines educativos y de investigación

3.11. Resumen de la Metodología

Este capítulo describió la metodología completa del proyecto, incluyendo:

- Selección y descripción de tres datasets públicos de malware
- Pipeline de preprocessamiento y conversión de binarios a imágenes
- Diseño de arquitecturas CNN personalizadas y uso de transfer learning
- Configuración de hiperparámetros y estrategias de entrenamiento
- Métricas de evaluación para clasificación multi-clase
- Consideraciones éticas y de seguridad

La siguiente sección presenta los resultados experimentales obtenidos aplicando esta metodología.

4

Experimentos y Resultados

Este capítulo presenta los resultados experimentales obtenidos al aplicar la metodología descrita en el Capítulo 3. Se analizan el rendimiento de diferentes arquitecturas CNN, se comparan resultados entre datasets, y se discuten las implicaciones de los hallazgos.

4.1. Configuración Experimental

4.1.1. Resumen de Experimentos Realizados

Se llevaron a cabo múltiples experimentos para evaluar:

1. **Arquitectura CNN baseline:** Modelo personalizado diseñado específicamente para el problema
2. **Transfer Learning:** Modelos preentrenados (VGG16, ResNet50, InceptionV3) adaptados mediante fine-tuning
3. **Comparación entre datasets:** Evaluación del rendimiento en cada uno de los tres datasets
4. **Análisis de generalización:** Entrenamiento en un dataset y evaluación en otro
5. **Estudio de hiperparámetros:** Variación de learning rate, batch size, y técnicas de regularización

4.1.2. Entorno de Ejecución

Todos los experimentos se ejecutaron bajo las siguientes condiciones:

- **Hardware:** [GPU/CPU específico]
- **Framework:** TensorFlow 2.x / PyTorch 1.x
- **Tiempo de entrenamiento promedio:** [X horas por modelo]
- **Número de épocas:** Hasta 100 con early stopping (típicamente convergencia en 30-50 épocas)

4.2. Análisis Exploratorio de Datos

4.2.1. Distribución de Familias de Malware

La Tabla 4.1 muestra la distribución de muestras por familia en cada dataset.

Cuadro 4.1: Distribución de muestras por familia en los tres datasets

| Familia | MalImg | Malevis | Blended |
|-----------------|-------------|-------------|-------------|
| Alureon | XXX | XXX | XXX |
| VB | XXX | XXX | XXX |
| Rbot | XXX | XXX | XXX |
| ... (completar) | ... | ... | ... |
| Total | XXXX | XXXX | XXXX |

Observación: Se identifica desbalance moderado entre clases, con algunas familias representando menos del 2 % del total.

4.2.2. Visualización de Muestras

La Figura 4.1 presenta ejemplos visuales de diferentes familias de malware. Se observan patrones texturales distintivos que justifican el enfoque de clasificación mediante CNN.

Figura 4.1: Ejemplos de imágenes de malware de diferentes familias. Cada fila corresponde a una familia distinta, mostrando la variabilidad intra-clase y las diferencias inter-clase.

4.3. Resultados de Clasificación

4.3.1. Rendimiento de Arquitecturas en MalImg Dataset

La Tabla 4.2 presenta los resultados de diferentes arquitecturas evaluadas en el MalImg Dataset.

Cuadro 4.2: Rendimiento de diferentes arquitecturas CNN en MalImg Dataset

| Modelo | Accuracy | Precision | Recall | F1-Score |
|--------------------------|----------|-----------|---------|----------|
| CNN Baseline | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| VGG16 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| ResNet50 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| InceptionV3 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |

Análisis: [Discutir qué modelo obtuvo mejor rendimiento y por qué]

4.3.2. Rendimiento en Malevis Dataset

Los resultados en Malevis Dataset se presentan en la Tabla 4.3.

Cuadro 4.3: Rendimiento en Malevis Dataset

| Modelo | Accuracy | Precision | Recall | F1-Score |
|--------------------------|----------|-----------|---------|----------|
| CNN Baseline | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| VGG16 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| ResNet50 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| InceptionV3 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |

4.3.3. Rendimiento en Blended Malware Dataset

Resultados en Blended Malware Dataset (Tabla 4.4).

Cuadro 4.4: Rendimiento en Blended Malware Dataset

| Modelo | Accuracy | Precision | Recall | F1-Score |
|--------------------------|----------|-----------|---------|----------|
| CNN Baseline | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| VGG16 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| ResNet50 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |
| InceptionV3 (Fine-tuned) | XX.XX % | XX.XX % | XX.XX % | XX.XX % |

4.4. Análisis Comparativo

4.4.1. Comparación entre Arquitecturas

La Figura 4.2 muestra una comparación visual del accuracy obtenido por cada arquitectura en los tres datasets.

Figura 4.2: Comparación de accuracy entre diferentes arquitecturas CNN evaluadas en los tres datasets.

Hallazgos principales:

jor rendimiento general
ferencias entre datasets
nplejidad y rendimiento

4.4.2. Análisis por Dataset

Observaciones sobre las características de cada dataset:

MalImg Dataset:

es, rendimiento general

Malevis Dataset:

as distintivas y desafíos

Blended Dataset:

versidad y generalización

4.5. Matrices de Confusión

4.5.1. Análisis de Errores de Clasificación

La Figura 4.3 presenta la matriz de confusión del modelo con mejor rendimiento.

Figura 4.3: Matriz de confusión del modelo [nombre del mejor modelo] en [dataset]. Los valores en la diagonal representan clasificaciones correctas.

Observaciones sobre patrones de confusión:

uentemente confundidas
etado compartidas, etc.
con mayor tasa de error

4.6. Curvas de Entrenamiento

4.6.1. Evolución de Loss y Accuracy

La Figura 4.4 muestra la evolución de loss y accuracy durante el entrenamiento del mejor modelo.

Figura 4.4: Curvas de entrenamiento y validación para el modelo [nombre]. (a) Loss function. (b) Accuracy.

Análisis:

cas hasta estabilización
entre train y validation
zación y early stopping

4.7. Análisis de Generalización

4.7.1. Entrenamiento y Evaluación Cross-Dataset

Se evaluó la capacidad de generalización entrenando en un dataset y evaluando en otro (Tabla 4.5).

Cuadro 4.5: Accuracy de generalización cross-dataset

| Entrenado en → Evaluado en | MalImg | Malevis | Blended |
|----------------------------|---------|---------|---------|
| MalImg | XX.XX % | XX.XX % | XX.XX % |
| Malevis | XX.XX % | XX.XX % | XX.XX % |
| Blended | XX.XX % | XX.XX % | XX.XX % |

Interpretación: [Discutir qué tanto disminuye el rendimiento al evaluar en datasets diferentes al de entrenamiento. Analizar qué combinaciones generalizan mejor.]

4.8. Estudio de Hiperparámetros

4.8.1. Impacto del Learning Rate

Se experimentó con diferentes valores de learning rate inicial para evaluar su impacto en convergencia y rendimiento final.

Cuadro 4.6: Impacto del learning rate en el rendimiento final

| Learning Rate | Accuracy (%) | Épocas hasta convergencia |
|---------------|--------------|---------------------------|
| 0.0001 | XX.XX | XX |
| 0.001 | XX.XX | XX |
| 0.01 | XX.XX | XX |

4.8.2. Efecto del Batch Size

La Tabla 4.7 muestra cómo el tamaño del batch afecta el rendimiento y tiempo de entrenamiento.

Cuadro 4.7: Impacto del batch size

| Batch Size | Accuracy (%) | Tiempo/Época (seg) | Uso Memoria (GB) |
|------------|--------------|--------------------|------------------|
| 16 | XX.XX | XX | XX |
| 32 | XX.XX | XX | XX |
| 64 | XX.XX | XX | XX |
| 128 | XX.XX | XX | XX |

4.8.3. Técnicas de Regularización

Evaluación del impacto de diferentes técnicas de regularización en la prevención de overfitting.

4.9. Análisis de Características Aprendidas

4.9.1. Visualización de Mapas de Activación

La Figura 4.5 muestra mapas de activación (mediante Grad-CAM o técnicas similares) que revelan qué regiones de las imágenes de malware son más relevantes para la clasificación.

Figura 4.5: Mapas de activación mostrando regiones de interés para la CNN al clasificar diferentes familias de malware. Las zonas más brillantes indican mayor importancia para la decisión del modelo.

Interpretación:

son más discriminativas
características interpretables
ras de la misma familia

4.9.2. Reducción Dimensional y Clustering

Se aplicó t-SNE sobre las representaciones aprendidas en la penúltima capa del modelo para visualizar separabilidad de clases (Figura 4.6).

Figura 4.6: Visualización t-SNE de las representaciones aprendidas. Cada punto representa una muestra de malware, coloreada por su familia. La formación de clusters indica buena separabilidad aprendida por el modelo.

4.10. Comparación con Estado del Arte

4.10.1. Comparación con Trabajos Previos

La Tabla 4.8 compara los resultados obtenidos en este proyecto con trabajos previos reportados en la literatura.

Cuadro 4.8: Comparación con estado del arte en MalImg Dataset

| Método | Accuracy (%) | Referencia |
|-------------------------------------|--------------|---------------------|
| [Trabajo 1] | XX.XX | [Autor et al., Año] |
| [Trabajo 2] | XX.XX | [Autor et al., Año] |
| [Trabajo 3] | XX.XX | [Autor et al., Año] |
| Este proyecto (mejor modelo) | XX.XX | -- |

Análisis: [Discutir si los resultados son competitivos con el estado del arte, considerar diferencias en configuración experimental]

4.11. Análisis de Tiempo de Inferencia

4.11.1. Eficiencia Computacional

La Tabla 4.9 presenta el tiempo promedio de inferencia por muestra para cada arquitectura.

Cuadro 4.9: Tiempo de inferencia promedio

| Modelo | Tiempo (ms/muestra) | Throughput (muestras/seg) |
|--------------|---------------------|---------------------------|
| CNN Baseline | XX.XX | XXX |
| VGG16 | XX.XX | XXX |
| ResNet50 | XX.XX | XXX |
| InceptionV3 | XX.XX | XXX |

Implicaciones: [Discutir trade-off entre precisión y velocidad. Comentar viabilidad para sistemas en tiempo real.]

4.12. Limitaciones y Desafíos Observados

Durante la experimentación se identificaron las siguientes limitaciones y desafíos:

4.12.1. Desbalance de Clases

A pesar de aplicar técnicas de ponderación, algunas familias minoritarias presentan menor recall, indicando dificultad para generarlas adecuadamente.

4.12.2. Familias Similares

Ciertas familias de malware comparten técnicas de empaquetado o estructuras, resultando en confusión sistemática entre ellas.

4.12.3. Dependencia de Tamaño de Imagen

El redimensionamiento de ejecutables muy grandes puede resultar en pérdida de información fina, mientras que ejecutables pequeños pueden sufrir de sobre-representación de padding.

4.12.4. Limitaciones de Generalización

El rendimiento cross-dataset disminuye notablemente, sugiriendo que los modelos capturan características específicas del dataset de entrenamiento.

4.13. Resumen de Resultados

Los principales hallazgos de este capítulo son:

1. El mejor modelo alcanzó un accuracy de **XX.XX %** en el dataset [nombre], demostrando la viabilidad del enfoque basado en CNN para clasificación de malware.
2. Los modelos con transfer learning [superaron/no superaron] significativamente al modelo baseline, indicando [interpretación].
3. Se identificaron patrones de confusión entre familias específicas, revelando desafíos inherentes a la similitud estructural entre ciertas clases.
4. El análisis de generalización cross-dataset reveló [resultado], sugiriendo [implicación].
5. Las visualizaciones de mapas de activación muestran que el modelo se enfoca en [descripción de regiones relevantes].

Estos resultados demuestran que el análisis visual de malware mediante CNN es un enfoque prometedor, aunque con limitaciones que deben considerarse para aplicaciones en entornos reales.

5

Conclusiones y Trabajo Futuro

Este capítulo final sintetiza los hallazgos principales del proyecto, discute las contribuciones realizadas, reconoce las limitaciones encontradas, y propone direcciones prometedoras para investigación futura en clasificación de malware mediante Deep Learning.

5.1. Resumen del Proyecto

Este proyecto abordó el problema crítico de la clasificación automatizada de familias de malware mediante técnicas de Deep Learning, específicamente Redes Neuronales Convolucionales aplicadas a representaciones visuales de ejecutables maliciosos.

Se implementó un pipeline completo que incluye:

- Recopilación y preprocesamiento de tres datasets públicos reconocidos (MalImg, Malevis, Blended Malware)
- Conversión de binarios ejecutables a representaciones de imágenes en escala de grises
- Diseño e implementación de múltiples arquitecturas CNN, incluyendo modelos personalizados y transfer learning
- Evaluación exhaustiva mediante métricas estándar de clasificación multi-clase
- Análisis comparativo entre arquitecturas y datasets
- Visualización e interpretación de características aprendidas

Los experimentos demostraron que el enfoque basado en visión por computador es viable y efectivo para la clasificación de malware, logrando accuracy de [XX.XX %] en el mejor caso, comparable con el estado del arte en la literatura.

5.2. Principales Hallazgos

5.2.1. Viabilidad del Enfoque Visual

Los resultados confirmaron que la conversión de ejecutables a imágenes preserva información estructural suficiente para permitir discriminación entre familias de malware. Las

CNN son capaces de aprender automáticamente características visuales discriminativas sin requerir ingeniería manual de features.

5.2.2. Rendimiento de Arquitecturas

Modelos Personalizados vs. Transfer Learning: [Discutir qué enfoque resultó más efectivo y por qué. Por ejemplo: "Los modelos preentrenados con fine-tuning superaron al modelo baseline en X %, sugiriendo que las características de bajo nivel aprendidas en ImageNet son transferibles al dominio de malware."]

Trade-offs Complejidad-Rendimiento: [Analizar si modelos más profundos siempre resultaron en mejor rendimiento o si hubo saturación/overfitting.]

5.2.3. Generalización y Robustez

El análisis cross-dataset reveló que [describir nivel de generalización]. Esto indica que [interpretar implicaciones: ¿los modelos aprenden características generales o específicas del dataset?].

5.2.4. Características Discriminativas

Las visualizaciones mediante Grad-CAM y análisis de mapas de activación mostraron que el modelo se enfoca en [describir qué secciones de los ejecutables: headers, secciones de código, recursos, etc.], lo cual es consistente con el conocimiento experto sobre diferencias estructurales entre familias de malware.

5.3. Respuesta a las Preguntas de Investigación

Retomando las preguntas planteadas en el Capítulo 1:

5.3.1. Pregunta Principal

¿Es posible clasificar eficazmente muestras de malware en sus respectivas familias mediante el análisis de sus representaciones visuales utilizando CNN?

Respuesta: Sí. Los resultados experimentales demuestran que las CNN pueden clasificar efectivamente familias de malware mediante análisis visual, logrando accuracy superior a [XX %], lo cual es competitivo con métodos tradicionales y enfoques de Machine Learning basados en características manuales.

5.3.2. Preguntas Secundarias

¿Qué arquitecturas CNN son más efectivas?

[Responder basándose en resultados: ResNet50 con fine-tuning demostró ser la arquitectura más efectiva, logrando..."]

¿Modelos preentrenados vs. arquitecturas específicas?

[Responder con evidencia experimental sobre el valor del transfer learning en este dominio.]

¿Qué características visuales son discriminativas?

[Resumir hallazgos de análisis de interpretabilidad: qué patrones estructurales distinguen familias.]

5.4. Contribuciones del Proyecto

Este proyecto realiza las siguientes contribuciones:

5.4.1. Contribuciones Técnicas

1. **Evaluación comparativa sistemática:** Análisis exhaustivo de múltiples arquitecturas CNN sobre tres datasets públicos bajo condiciones controladas.
2. **Estudio de generalización:** Evaluación cross-dataset que cuantifica la transferibilidad de modelos entre diferentes colecciones de malware.
3. **Pipeline completo y reproducible:** Implementación de pipeline end-to-end desde preprocesamiento hasta evaluación, facilitando replicación y extensión del trabajo.
4. **Análisis de interpretabilidad:** Visualizaciones de características aprendidas que proveen insights sobre el proceso de decisión del modelo.

5.4.2. Contribuciones Prácticas

1. Demostración de viabilidad para integración en sistemas reales de detección de amenazas.
2. Identificación de limitaciones y desafíos prácticos para despliegue en producción.
3. Recomendaciones basadas en evidencia sobre elección de arquitecturas y configuraciones.

5.5. Limitaciones del Estudio

Es importante reconocer las siguientes limitaciones:

5.5.1. Limitaciones de Datasets

- **Distribución temporal:** Los datasets pueden no representar amenazas más recientes o emergentes.
- **Desbalance de clases:** Algunas familias están sub-representadas, afectando la capacidad del modelo para aprenderlas.
- **Sesgo de selección:** Los datasets públicos pueden no reflejar la distribución real de malware en entornos productivos.
- **Limitación a Windows:** Los datasets empleados contienen principalmente malware de Windows, limitando la aplicabilidad a otras plataformas.

5.5.2. Limitaciones Metodológicas

- **Análisis estático únicamente:** No se considera comportamiento dinámico, que podría proporcionar información complementaria.
- **Pérdida de información:** El redimensionamiento de imágenes puede resultar en pérdida de detalles finos en ejecutables grandes.
- **Vulnerabilidad a adversarios:** No se evaluó robustez ante ataques adversariales diseñados para engañar al clasificador.
- **Costo computacional:** El entrenamiento de modelos profundos requiere recursos GPU significativos, limitando accesibilidad.

5.5.3. Limitaciones de Interpretabilidad

Aunque se realizó análisis de mapas de activación, la comprensión completa de qué características específicas aprende el modelo permanece parcialmente como "caja negra", dificultando la explicación de decisiones erróneas.

5.6. Trabajo Futuro

Los resultados de este proyecto abren múltiples direcciones prometedoras para investigación futura:

5.6.1. Extensiones Metodológicas

Enfoques Híbridos

Combinar análisis visual con otras fuentes de información:

- Integración con análisis de secuencias de opcodes (usando RNN/LSTM)
- Fusión con características estáticas extraídas manualmente (headers PE, importaciones)
- Incorporación de análisis dinámico (llamadas API, comportamiento en sandbox)

Arquitecturas Avanzadas

Explorar arquitecturas más recientes:

- **Vision Transformers (ViT):** Evaluar si attention mechanisms mejoran la captura de relaciones de largo alcance en ejecutables
- **EfficientNet:** Modelos optimizados para mejor trade-off accuracy-eficiencia
- **Neural Architecture Search (NAS):** Búsqueda automatizada de arquitecturas óptimas para el dominio específico

Técnicas de Few-Shot Learning

Implementar aprendizaje con pocos ejemplos para manejar familias nuevas o raras sin reentrenamiento completo:

- Siamese Networks para aprendizaje de similitud
- Prototypical Networks
- Meta-learning approaches

5.6.2. Mejoras en Robustez

Defensa Adversarial

Evaluar y mejorar robustez ante ataques adversariales:

- Generación de muestras adversariales específicas para malware
- Entrenamiento adversarial para mejorar robustez
- Certificación de robustez mediante métodos formales

Detección de Out-of-Distribution

Implementar mecanismos para identificar muestras fuera de la distribución de entrenamiento:

- Métodos basados en confianza/entropía de predicciones
- Autoencoders para detección de anomalías
- Uncertainty quantification mediante ensembles o métodos Bayesianos

5.6.3. Ampliación de Datasets

Datasets Multi-Plataforma

Expandir el estudio a malware de otras plataformas:

- Android malware (archivos APK convertidos a imágenes)
- Malware de Linux
- Malware de macOS
- Malware para dispositivos IoT

Datasets Dinámicos

Construir datasets actualizados continuamente con amenazas emergentes para evaluar la adaptabilidad temporal de los modelos.

5.6.4. Aplicaciones Prácticas

Sistema de Detección en Tiempo Real

Desarrollar un prototipo de sistema de detección integrable en entornos de producción:

- Optimización de modelos para inferencia eficiente (quantization, pruning)
- Pipeline de procesamiento en tiempo real
- Interfaz para analistas de seguridad
- Integración con SIEM (Security Information and Event Management)

Análisis Forense

Aplicar el enfoque a análisis forense digital:

- Identificación de familias en incidentes de seguridad
- Clustering de muestras desconocidas
- Trazabilidad de variantes de amenazas

5.6.5. Investigación en Interpretabilidad

Explicabilidad Mejorada

Desarrollar métodos más sofisticados para interpretar decisiones:

- Análisis de sensibilidad local mediante perturbaciones
- Identificación de características mínimas necesarias para clasificación
- Generación de explicaciones en lenguaje natural para analistas

Extracción de Conocimiento Experto

Utilizar los modelos entrenados para extraer conocimiento sobre diferencias estructurales entre familias que pueda informar análisis manual de malware.

5.6.6. Benchmarking y Estandarización

Protocolo de Evaluación Estándar

Contribuir al establecimiento de protocolos de evaluación estandarizados para la comunidad:

- Definición de splits de entrenamiento/validación/prueba reproducibles
- Métricas estandarizadas considerando costos asimétricos (falsos negativos vs. falsos positivos)
- Datasets de referencia actualizados periódicamente

Challenge y Competición

Organizar competiciones académicas para impulsar avances en el campo, similar a challenges en visión por computador o NLP.

5.7. Implicaciones para la Ciberseguridad

Los resultados de este proyecto tienen implicaciones significativas para la práctica de la ciberseguridad:

5.7.1. Para Profesionales de Seguridad

- **Automatización:** Reducción del esfuerzo manual en análisis de grandes volúmenes de muestras sospechosas.
- **Velocidad:** Clasificación en tiempo prácticamente real permite respuesta más rápida a incidentes.
- **Escalabilidad:** Capacidad de procesar cantidades masivas de datos sin incremento lineal en recursos humanos.

5.7.2. Para Desarrolladores de Soluciones de Seguridad

- **Complementariedad:** Los métodos basados en Deep Learning pueden complementar (no reemplazar) soluciones tradicionales.
- **Adaptabilidad:** Los modelos pueden reentrenarse periódicamente para adaptarse a nuevas amenazas.
- **Multi-modal:** Posibilidad de fusionar análisis visual con otras técnicas para detección más robusta.

5.7.3. Para la Investigación Académica

- **Fundamentos sólidos:** Este trabajo proporciona evidencia experimental sobre la viabilidad del enfoque visual.
- **Dirección prometedora:** Se identifican múltiples líneas de investigación futuras con potencial impacto.
- **Metodología reproducible:** El pipeline implementado facilita la extensión y comparación con nuevos métodos.

5.8. Lecciones Aprendidas

Durante el desarrollo de este proyecto se adquirieron varias lecciones valiosas:

5.8.1. Técnicas

- **Importancia del preprocessamiento:** La calidad de las imágenes generadas impacta significativamente el rendimiento final.
- **Balance entre complejidad y datos:** Modelos muy profundos pueden no ser necesarios con datasets limitados.
- **Regularización esencial:** Dropout y data augmentation son críticos para evitar overfitting en este dominio.
- **Transfer learning valioso:** Las características de bajo nivel de ImageNet son sorprendentemente transferibles.

5.8.2. Prácticas

- **Experimentación sistemática:** La variación controlada de hiperparámetros es esencial para optimización.
- **Validación rigurosa:** La evaluación en conjunto de prueba separado es imprescindible para estimaciones realistas.
- **Interpretabilidad importante:** La capacidad de explicar decisiones es crucial para adopción en seguridad.

5.9. Consideraciones Éticas

El desarrollo de herramientas automatizadas de análisis de malware plantea consideraciones éticas importantes:

5.9.1. Uso Responsable

Los modelos y técnicas desarrollados deben emplearse exclusivamente para:

- Defensa legítima de sistemas y redes
- Investigación académica con fines educativos
- Análisis forense en el contexto de incidentes de seguridad

Nunca deben utilizarse para:

- Desarrollo de nuevas amenazas
- Evasión de sistemas de seguridad con intención maliciosa
- Ataques a sistemas sin autorización explícita

5.9.2. Transparencia

Es importante mantener transparencia sobre:

- Limitaciones de los modelos (tasas de falsos negativos/positivos)
- Datasets utilizados y sus sesgos inherentes
- Condiciones bajo las cuales los resultados son válidos

5.9.3. Privacidad y Confidencialidad

Al trabajar con muestras de malware reales, se debe:

- Proteger cualquier información sensible que puedan contener los ejecutables
- Cumplir con regulaciones de manejo de código malicioso
- Evitar diseminación no controlada de muestras activas

5.10. Reflexiones Finales

La clasificación automatizada de malware mediante Deep Learning representa un área de investigación madura y prometedora en la intersección de inteligencia artificial y ciberseguridad. Este proyecto ha demostrado que el enfoque basado en análisis visual de ejecutables es técnicamente viable y puede alcanzar niveles de rendimiento competitivos con el estado del arte.

Sin embargo, es fundamental reconocer que ninguna solución única resolverá por completo el problema de la detección de malware. Las amenazas continúan evolucionando, y los atacantes adaptan constantemente sus técnicas. Por lo tanto, los sistemas efectivos de ciberseguridad requieren enfoques en capas (defense in depth) que combinen múltiples técnicas complementarias.

El Deep Learning, y específicamente las CNN aplicadas a representaciones visuales, constituyen una herramienta poderosa en el arsenal del defensor, pero deben emplearse en conjunto con:

- Análisis heurístico y basado en firmas
- Detección comportamental
- Inteligencia de amenazas
- Supervisión humana experta

La investigación futura debe enfocarse no solo en mejorar la precisión de los modelos, sino también en su robustez, interpretabilidad, y aplicabilidad práctica en entornos de producción con requisitos estrictos de latencia y confiabilidad.

5.11. Conclusión

Este proyecto ha explorado exitosamente la aplicación de Redes Neuronales Convolucionales para la clasificación de familias de malware mediante análisis visual de ejecutables. Los resultados experimentales confirman la hipótesis inicial: las CNN son capaces de aprender automáticamente representaciones discriminativas que permiten clasificación efectiva de malware.

Se han realizado contribuciones tanto técnicas (evaluación comparativa exhaustiva, análisis de generalización, estudio de interpretabilidad) como prácticas (identificación de limitaciones reales, recomendaciones para despliegue).

Las direcciones de trabajo futuro identificadas ofrecen oportunidades prometedoras para avanzar el estado del arte en este campo crítico para la seguridad de sistemas informáticos modernos.

En última instancia, este trabajo contribuye a la creciente evidencia de que las técnicas de Deep Learning representan una herramienta valiosa y cada vez más madura para abordar desafíos complejos en ciberseguridad, con potencial para impacto real en la protección de infraestructuras digitales frente a amenazas en constante evolución.

