

**НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ
«КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ
імені ІГОРЯ СІКОРСЬКОГО»
НАВЧАЛЬНО-НАУКОВИЙ ФІЗИКО-ТЕХНІЧНИЙ ІНСТИТУТ
Кафедра інформаційної безпеки**

До захисту допущено
Завідувач кафедри

_____ Дмитро ЛАНДЕ
(підпис)

« _____ » _____ 2024 р.

**Дипломна робота
на здобуття ступеня бакалавра
за освітньо-професійною програмою «Системи, технології та математичні
методи кібербезпеки»
спеціальності 125 «Кібербезпека»**

на тему: Метод класифікації цифрових доказів на основі аналізу метаданих методами
машинного навчання

Виконала: здобувач вищої освіти **IV** курсу, групи ФБ-02

Лета Яна Василівна

Керівник к.т.н., доцент Барановський Олексій Миколайович



Рецензент к.ф.-м.н., доцент Хмельницький Микола Олексійович



Засвідчую, що у цій дипломній роботі немає
запозичень з праць інших авторів без відповідних
посилань.

Здобувач вищої освіти _____
(підпис)

Київ – 2024 року

НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ
«КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ
імені ІГОРЯ СІКОРСЬКОГО»
НАВЧАЛЬНО-НАУКОВИЙ ФІЗИКО-ТЕХНІЧНИЙ ІНСТИТУТ
Кафедра інформаційної безпеки

Рівень вищої освіти – перший (бакалаврський)

Спеціальність – 125 «Кібербезпека»

Освітньо-професійна програма «Системи, технології та математичні методи кібербезпеки»

ЗАТВЕРДЖУЮ

Завідувач кафедри

_____ Дмитро ЛАНДЕ

(підпис)

«__» _____ 2024 р.

ЗАВДАННЯ
на дипломну роботу здобувачу вищої освіти

Лета Яна Василівна

1. Тема роботи: «Метод класифікації цифрових доказів на основі аналізу метаданих методами машинного навчання»,
керівник роботи: Барановський Олександр Миколайович, к.т.н., доцент каф. ІБ,
затверджені наказом по університету від 31 травня 2024 р. № 2251-с
2. Термін подання здобувачем вищої освіти роботи 10 червня 2024 р.
3. Вихідні дані до роботи: існуючі методи аналізу цифрових доказів методами машинного навчання.
4. Зміст роботи: Провести аналіз різних літературних джерел для обґрунтування теми. Провести дослідження над існуючими методами аналізу цифрових доказів. Розробити метод автоматичного визначення пріоритетних підозрілих файлів. Зібрати та підготувати набір даних для навчання і тестування моделі. Здійснити навчання та тестування моделі машинного навчання. Оцінити результати та проаналізувати можливі подальші зміни для підвищення ефективності та вдосконалення створеного методу класифікації.
5. Перелік ілюстративного матеріалу (із зазначенням плакатів, презентацій тощо): презентація.
6. Дата видачі завдання: 5 лютого 2024 р.

Календарний план

№ з/п	Назва етапів виконання дипломної роботи	Термін виконання етапів дипломної роботи	Примітка
1	Вибір тематики роботи	10.02.2024	Виконано
2	Дослідження різних напрямів обраної тематики	20.02.2024	Виконано
3	Визначення теми роботи	01.03.2024	Виконано
4	Дослідження існуючих методів аналізу цифрових доказів	15.03.2024	Виконано
5	Опис методу визначення пріоритету підозрілих цифрових доказів	18.03.2024 - 08.04.2024	Виконано
6	Збір і підготовка набору даних для створення методу класифікації	10.04.2024 - 19.04.2024	Виконано
7	Програмна реалізація моделі	20.04.2024 - 18.05.2024	Виконано
8	Створення змісту дипломної роботи	20.05.2024	Виконано
9	Написання теоретичної частини роботи	21.05.2024 - 25.05.2024	Виконано
10	Оформлення практичної частини	26.05.2024 - 30.05.2024	Виконано
11	Оцінка результатів роботи	31.05.2024 - 03.06.2024	Виконано
12	Аналіз майбутніх дослідження для вдосконалення методу класифікації	04.06.2024 - 05.06.2024	Виконано
13	Створення презентації для передзахисту	07.06.2024 - 09.06.2024	Виконано
14	Підготовка до захисту	12.06.2024	Виконано

Здобувач вищої освіти

Керівник роботи



Яна ЛЕТА

(Власне ім'я, ПРІЗВИЩЕ)

Олексій БАРАНОВСЬКИЙ

(Власне ім'я, ПРІЗВИЩЕ)

РЕФЕРАТ

Обсяг роботи 67 сторінок, 7 ілюстрацій, 2 таблиці, 2 додатки і 39 бібліографічних найменувань за переліком посилань.

Метою роботи є розробка методу, що використовує алгоритми машинного навчання для автоматичної класифікації цифрових доказів на визначення підозрілості на основі аналізу їх метаданих, щоб пришвидшити процесу відбору релевантних файлів під час цифрових криміналістичних розслідувань.

Об'єктом дослідження є метадані цифрових файлів, наявні програмні засоби і методи для аналізу цифрових доказів.

Предметом дослідження є методи машинного навчання та обробка метаданих, що використовуються для автоматичної класифікації цифрових доказів.

У даній роботі були застосовані такі методи дослідження, як аналіз літературних джерел для визначення сучасних методів аналізу цифрових доказів; порівняльний аналіз існуючих алгоритмів для виявлення найефективніших підходів; збір та видобування метаданих для створення набору даних; використання методів машинного навчання для навчання і тестування моделей.

Створений метод класифікації цифрових доказів на основі аналізу метаданих рекомендовано в подальшому інтегрувати з системами автоматизованого збору метаданих, що відкриває можливість широкого використання розробленого рішення в цифрових криміналістичних розслідуваннях, дозволяючи значно полегшити та прискорити процес аналізу цифрових доказів.

Результати роботи були представлені на XXII Всеукраїнській науково-практичній конференції студентів, аспірантів та молодих вчених «Теоретичні і прикладні проблеми фізики, математики та інформатики» 2024 року.

Ключові слова: цифрова криміналістика, цифрові докази, метадані, класифікація, машинне навчання.

ABSTRACT

The volume of the work is 67 pages, 7 illustrations, 2 tables, 2 appendices and 39 bibliographic names according to the list of references.

The aim of the work is to develop a method that uses machine learning algorithms to automatically classify digital evidence for suspiciousness based on the analysis of their metadata, in order to speed up the process of selecting relevant files during digital forensic investigations.

The object of the study is the metadata of digital files, available software tools and methods for analyzing digital evidence.

The subject of research is machine learning methods and metadata processing used for automatic classification of digital evidence.

This work used such research methods as the analysis of literary sources to determine modern methods of digital evidence analysis; comparative analysis of existing algorithms to identify the most effective approaches; collecting and extracting metadata to create a dataset; the use of machine learning methods for training and testing models.

The created method of digital evidence classification based on metadata analysis is recommended to be further integrated with automated metadata collection systems, which opens up the possibility of wide use of the developed solution in digital forensic investigations, allowing to significantly facilitate and speed up the process of digital evidence analysis.

The results of the work were presented at the XXII All-Ukrainian scientific and practical conference of students, postgraduates and young scientists "Theoretical and Applied Problems of Physics, Mathematics and Informatics" in 2024.

Keywords: digital forensics, digital evidence, metadata, classification, machine learning.

ЗМІСТ

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ, СИМВОЛІВ, ОДИНИЦЬ, СКОРОЧЕНЬ І ТЕРМІНІВ	9
ВСТУП.....	11
1 ЦИФРОВА КРИМІНАЛІСТИКА І МЕТАДАНІ ЦИФРОВИХ ФАЙЛІВ	14
1.1 Цифрова криміналістика.....	15
1.2 Процес розслідування у цифровій криміналістиці.....	17
1.3 Історія метаданих цифрових файлів	19
1.4 Проблема класифікації метаданих.....	21
Висновки до розділу 1	24
2 КЛАСИФІКАЦІЯ ЦИФРОВИХ ДОКАЗІВ ІЗ ВИКОРИСТАННЯМ МАШИННОГО НАВЧАННЯ.....	25
2.1 Класифікація цифрових доказів	25
2.2 Існуючі методи аналізу цифрових доказів на основі машинного навчання...30	30
2.3 Ефективність існуючих методів аналізу цифрових доказів.....	34
2.4 Визначення пріоритетних підозрілих файлів.....	38
2.5 Приклади застосування машинного навчання у цифровій криміналістиці ...39	39
2.6 Обмеження машинного навчання у цифровій криміналістиці.....	44
Висновки до розділу 2.....	46
3 РОЗРОБКА МЕТОДУ КЛАСИФІКАЦІЇ ЦИФРОВИХ ДОКАЗІВ НА ОСНОВІ МЕТАДАНИХ.....	47
3.1 Збір даних	47
3.2 Обробка даних	50
3.3 Створення моделі класифікації та аналіз результатів	53
3.4 Майбутні дослідження	56
Висновки до розділу 3	57
ВИСНОВКИ.....	58
ПЕРЕЛІК ДЖЕРЕЛ ПОСИЛАНЬ	60
ДОДАТОК А.....	64
ДОДАТОК Б	66

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ, СИМВОЛІВ, ОДИНИЦЬ, СКОРОЧЕНЬ І ТЕРМІНІВ

AB — Adaptive Boosting, Алгоритм адаптивного підсилення.

ANN — Artificial Neural Networks, Штучні нейронні мережі.

CSAM — Child Sexual Abuse Media, медіа матеріали сексуального насильства над дітьми.

CSV — Comma-Separated Values, текстовий формат для представлення табличних даних.

DDoS — Distributed Denial of Service, атака на комп'ютерну систему або мережу, в якій атакуючий намагається перевантажити її шляхом надлишкового надходження запитів чи трафіку з багатьох джерел одночасно.

DT — Decision Tree, Дерево рішень.

FAT — File Allocation Table, сімейство файлових систем розроблених компанією Microsoft.

GNB — Gaussian Naive Bayes, Гаусів Наївний Баєсів класифікатор.

IBM OS/360 — IBM Operating System/360, сімейство операційних систем компанії IBM.

IDS — Intrusion Detection System, система виявлення вторгнень.

IP — Internet Protocol, міжмережевий протокол.

K-Means — Метод k-середніх.

KNN — K-Nearest Neighbors, Метод k-найближчих сусідів.

LR — Logistic Regression, Логістична регресія.

MapReduce — програмна модель для обробки та генерації великих обсягів даних паралельним та розподіленим способом, розроблена компанією Google.

MS-DOS — Microsoft Disk Operating System, операційна система компанії Microsoft.

NB — Naive Bayes, Наївний Баєсів класифікатор.

NTFS — New Technology File System, стандартна файлова система для сімейства операційних систем Microsoft Windows NT.

RAW — тип збереження даних без будь-яких додаткових обробок або кодувань.

RF — Random Forest, Випадковий ліс.

ROC-крива — Receiver Operating Characteristic curve, відображає відношення між чутливістю (True Positive Rate, TPR) і специфічністю (1 - False Positive Rate, FPR) класифікатора при різних порогових значеннях.

SMF — System Management Facilities, компонент IBM z/OS для мейнфреймів.

SVM — Support Vector Machine, Метод опорних векторів.

ОС — операційна система.

ПЗ — програмне забезпечення.

Влучність — Precision, метрика оцінки класифікатора, яка вимірює відсоток правильно виявлених позитивних елементів серед всіх елементів, які були виявлені як позитивні.

Повнота — Recall, метрика оцінки класифікатора, яка вимірює відсоток правильно виявлених позитивних елементів серед всіх дійсних позитивних елементів.

Точність — Accuracy, метрика оцінки класифікатора, яка визначає відсоток правильно класифікованих елементів з усіх елементів, які були класифіковані.

F-міра — F-score, гармонічне середнє між повнотою (Recall) та влучністю (Precision).

ВСТУП

У сучасному світі де технології постійно розвиваються, обсяг даних, які збираються й аналізуються під час цифрових судових розслідувань, зростає експоненційно. Цей процес стає дедалі складнішим через велику кількість інформації, яка часто не має відношення до справи. В результаті цього, вручну знайти релевантні файли стає надзвичайно важко і вимагає значних часових і ресурсних витрат, що нагадує пошук голки в копиці сіна. З розвитком методів машинного навчання з'являються нові можливості для автоматизації та оптимізації процесів аналізу цифрових доказів.

У даному дослідженні пропонується метод, що використовує машинне навчання для класифікації цифрових доказів на основі аналізу метаданих. Цей підхід спрямований на те, щоб допомогти фахівцям у галузі цифрової криміналістики ефективніше визначати пріоритетність файлів, які можуть бути підозрілими та мати відношення до розслідування.

«Великі цифрові криміналістичні дані» є проблемою, з якою стикаються правоохоронні органи в усьому світі, оскільки поширеність цифрових пристроїв постійно зростає [1]. Як швидко виявити відповідні розслідуванню файли, це проблема, яку потрібно вирішити. Даний метод визначає потенційну підозрілість цифрового файлу, а не видає кінцевий результат, бо для нього потрібно рішення експерта. Використовується підхід керованого машинного навчання, що базується на результатах попередніх оброблених випадків.

Робота є актуальною через значні виклики для фахівців, які займаються аналізом цифрових доказів. Використання методів машинного навчання для процесу класифікації цифрових доказів є перспективним напрямком, який дозволяє значно зменшити витрати часу та зусиль. Впровадження новітніх технологій у правоохоронні органи є актуальним завданням, розробка

ефективного інструменту для аналізу цифрових доказів може значно підвищити ефективність розслідувань та сприяти боротьбі з кіберзлочинністю, тому подальше вдосконалення такого методу для автоматизації задачі класифікації може принести велику користь.

Метою даного дослідження є розробка методу, який використовує алгоритми машинного навчання для автоматичної класифікації цифрових доказів на визначення підозрілості на основі аналізу їх метаданих, щоб пришвидшити процесу відбору релевантних файлів під час цифрових криміналістичних розслідувань.

Завдання дослідження:

- 1) Провести дослідження над існуючими методами аналізу цифрових доказів.
- 2) Розробити метод автоматичного визначення пріоритетних підозрілих файлів.
- 3) Зібрати та підготувати набір даних для навчання і тестування моделі.
- 4) Здійснити навчання та тестування моделі машинного навчання.
- 5) Оцінити результати та проаналізувати можливі подальші зміни для підвищення ефективності та вдосконалення створеного методу класифікації.

Об'єктом дослідження є метадані цифрових файлів, наявні програмні засоби і методи для аналізу цифрових доказів.

Предметом дослідження є методи машинного навчання та обробка метаданих, що використовуються для автоматичної класифікації цифрових доказів.

Методи дослідження: аналіз літературних джерел для визначення сучасних методів аналізу цифрових доказів; порівняльний аналіз існуючих алгоритмів для виявлення найефективніших підходів; збір та видобування метаданих для

створення набору даних; використання методів машинного навчання для навчання і тестування моделей.

Було застосовано методи машинного навчання для задачі класифікації цифрових файлів на основі їх метаданих, адже у відкритому доступі не було знайдено відповідного програмного рішення.

Практичне значення одержаних результатів полягає у створенні програмного засобу для автоматичної класифікації цифрових доказів на основі аналізу метаданих з подальшою можливістю інтеграції з системами автоматизованого збору метаданих. Це відкриває можливість широкого використання розробленого рішення в цифрових криміналістичних розслідуваннях, дозволяючи значно полегшити та прискорити процес аналізу цифрових доказів.

Апробація результатів роботи: робота була представлена на XXII Всеукраїнській науково-практичній конференції студентів, аспірантів та молодих вчених «Теоретичні і прикладні проблеми фізики, математики та інформатики» із 13 по 17 травня 2024 року в місті Київ, Україна [2].

1 ЦИФРОВА КРИМІНАЛІСТИКА І МЕТАДАНІ ЦИФРОВИХ ФАЙЛІВ

Часто зневажають, що за звичайними документами, зображеннями та іншими цифровими файлами ховається прихована інформація, що міститься у метаданих. Ці дані неможливо побачити під час перегляду файлу, але вони можуть бути надзвичайно цінними для цифрових судових розслідувань.

Різні типи файлів містять різні типи метаданих. Наприклад, документ Word може містити інформацію про автора, дату створення, кількість разів редагування та навіть назви комп'ютерів, на яких він відкривався. Зображення можуть містити інформацію про камеру, яка його зробила, дату та час зйомки, а також GPS-координати, якщо фото було зроблено з GPS-пристроєм.

Метадані можуть надати величезну кількість інформації цифровому слідчому, а також їх можна використовувати для звинувачення або виправдання автора, рецензента, власника або видавця документа чи файлу.

Метадані можуть надати слідчому велику кількість інформації про файли, які досліджуються. Крім того, судовий слідчий може використовувати метадані для отримання інформації, наприклад, автора файлу, дату та час створення, кількість разів, коли файл було змінено, у тому числі, коли відбулася зміна [3].

Метадані файлової системи, включені до категорії метаданих файлової системи, — це інформація про розмір файлу, конкретні виділені одиниці даних і доступ до міток дати чи часу у файлі [3].

Метадані можуть бути корисними при вирішенні правових спорів, оскільки їх можна використовувати як доказ для підтвердження або спростування інших доказів, наданих у судовій справі. Крім того, метадані можна зберігати на різних

сайтах у файлах. Дослідники можуть отримати інформацію з цих метаданих, яка корелює з правом власності та потенційними власниками [4][5].

Метадані, та прихована інформація, що міститься в цифрових файлах, може стати ключовим елементом у розслідуванні злочину. Розкриваючи сліди змін або маніпуляцій, метадані дають слідчим та юристам цінну інформацію для формування висновків щодо справи.

Однак, успіх розслідування залежить не лише від наявності доказів, а й від їх належного збору та аналізу. Неправильне використання інструментів, системні помилки, приховування виправдувальних фактів, спотворення інформації, некомпетентність або свідоме фальсифікування доказів – все це може призвести до невірних висновків та помилкових звинувачень. Таким чином, практикуючий юрист повинен розуміти, як збираються цифрові докази; і зв'язок між процесом збору та підтвердженням потенційних доказів [6].

Отже, метадані можуть створювати докази, які не завжди є очевидними, і слідчий повинен забезпечити збір таких даних і перевірку їх для представлення в суді. Важливо, щоб юристи, які займаються справами, пов'язаними з метаданими, знали про природу та цінність метаданих для допомоги в процесі судового переслідування.

1.1 Цифрова криміналістика

Криміналістику визначають як характеристику доказів, яка підтверджує їхню придатність для визнання фактом і здатність переконувати на основі доказів або високої статистичної достовірності [7]. Якщо застосувати це визначення до цифрової криміналістики, стає зрозуміло, що ця сфера охоплює представників з різних дисциплін і з різним досвідом. Окрім інформатики, цифрова криміналістика є важливою для спеціалістів у таких галузях, як право, правоохоронні органи, політика та стандартизація. Таким чином, література з

цифрової криміналістики охоплює широкий спектр матеріалів і повинна бути орієнтована на різноманітну аудиторію, включаючи експертів із цифрових доказів, юристів, адвокатів, суддів, політиків, розробників і дослідників [8].

Цифрова криміналістика – це процес аналізу та представлення цифрових доказів, зібраних із різних джерел, таких як бази даних, комп'ютери та цифрові зображення [9]. Зростаюча кількість смарт-пристроїв у нашому повсякденному житті створює різноманітність даних з різними категоріями та характеристиками. Під час розслідування цифрової криміналістики дані збираються та аналізуються, щоб допомогти слідчим виявити і запобігти несанкціонованому доступу до інформації [10]. У багатьох випадках дані та докази можуть бути видалені з пристрою після скоєння злочину. Цей процес є надзвичайно важливим для слідчих, оскільки він може допомогти встановити точний характер злочину та ідентифікувати жертв [11]. Однак, нестача людських ресурсів для проведення детальних розслідувань може значно збільшити тривалість цього процесу.

Ця наукова галузь фокусується на аналізі та збереженні даних, зібраних на різних носіях. Хоча її витоки можна простежити до 1980-х років, прискорений розвиток галузі почався в 1990-х роках з появою багатокористувацьких, багатозадачних і глобальних мереж [12]. Зі зростанням кількості кіберзагроз та атак цифрова криміналістика стала однією з найважливіших сфер кібербезпеки.

Цей розділ криміналістики також зосереджується на правових процедурах, пов'язаних з аналізом і захистом цифрової інформації. Цифрова криміналістика включає ідентифікацію та вилучення інформації з різних джерел, яку можна використовувати для оцінки даних у цивільному чи кримінальному процесі [13]. Цей процес передбачає застосування наукових і технологічних методів для аналізу даних, створених різними цифровими пристроями [14]. Основна мета цифрової криміналістики – збір доказів для встановлення фактів, пов'язаних з інцидентом. Під час розслідувань зазвичай ставлять питання за схемою «5WH» (англ. Who?

What? When? Why? Where? How?): хто був причетний, який стався інцидент і коли, чому, де і коли це сталося. Відповіді на ці питання допомагають слідчим підтвердити деталі інциденту [15].

1.2 Процес розслідування у цифровій криміналістиці

За даними Національного інституту стандартів і технологій, процес розслідування цифрової криміналістики складається з чотирьох основних етапів, які допомагають організаціям зрозуміти важливість їхніх розслідувань. Ці процедури можуть варіюватися в залежності від складності дослідження [16]. З розвитком цифрових технологій зростає кількість джерел даних, які можна зібрати. Ось короткий опис кожного етапу:

1) Збір даних: першим етапом є визначення потенційних джерел даних. Дані зазвичай збираються з персональних комп'ютерів, ноутбуків, мобільних пристроїв і серверів. Аналітики також повинні враховувати інші джерела даних, наприклад, журнали Інтернет-провайдера організації [17].

2) Експертиза: на другому етапі перевіряються зібрані дані. Використовуючи методи та інструменти цифрової криміналістики, з даних витягується необхідна інформація. Це включає виявлення файлів, які містять важливу інформацію, навіть якщо вона прихована за допомогою стиснення, контролю доступу чи шифрування [16][17].

3) Аналіз: цей етап включає виконання наукових процедур для ідентифікації людей, місць і подій та визначення їхніх зв'язків. Аналізу підлягають дані, зібрані з різних джерел. Наприклад, журнали IDS можуть містити інформацію про конкретного користувача, а журнали аудиту – про конкретний хост. Інструменти для керування подіями безпеки допомагають зіставити та зібрати ці дані [16].

4) Звітування: останнім етапом є звітування, яке включає аналіз зібраних даних і представлення результатів в офіційній документації. Хоча визначення точної причини події може бути складним, зібрана інформація допомагає аналітикам зрозуміти подію та запобігти її повторенню в майбутньому [13].

Більшість сучасної літератури та рекомендацій з цифрової криміналістики зосереджено на пошуку даних та інформації, наявної в існуючих системах. Враховуючи різноманітність цільової аудиторії та різні рівні знань у галузі обчислювальної техніки, однією з основних цілей є навчання практиків базовим процедурам пошуку та аналізу доказів у сучасних комп'ютерних системах. Більшість поточних робіт пояснює методи відновлення даних із систем у різних формах. Що стосується метаданих, деякі дослідження також розглядають криміналістичну цінність та якість знайденої інформації [8].

Під криміналістичною цінністю розуміється здатність робити висновки про події в системі на основі даних. Наприклад, мітки часу мають високу цінність для реконструкції подій, оскільки вони дозволяють упорядкувати операції з файлами на часовій шкалі, за умови, що мітки часу не були підроблені та що системний годинник правильний. Водночас інформація про контроль доступу має меншу цінність для реконструкції подій, оскільки відображає переважно статичні системні політики, а не окремі події. Інформація, яку можна отримати з даних про контроль доступу, включає перелік користувачів або груп, які могли мати доступ до об'єкта в системі. Для обґрунтованих висновків потрібні додаткові докази, такі як мітки часу або дані про входи в систему [8].

Якість інформації визначається її достовірністю. Наприклад, у деяких операційних системах мітки часу доступу та модифікації файлу можуть бути довільно змінені його власником, що знижує достовірність таких даних [8].

Хоча існує багато методів для керування великим обсягом даних, зібраних цифровим криміналістом, проте вони не є такими ефективними, як людський

мозок. Натомість дослідники використовують машинне навчання для більш ефективного аналізу та збору даних [18].

1.3 Історія метаданих цифрових файлів

Визначення терміну «метадані» часто трактують як «дані про дані». Більш інформативне визначення [19] описує їх як «інформацію про об'єкт, фізичний чи цифровий». Отже, метадані є описовими даними і основою для розробок у сфері цифрових збірок. Основні елементи метаданих включають назву, автора, рік публікації та подібні прості бібліографічні дані. Розширені структури метаданих також охоплюють технічні характеристики, властивості авторського права, анотації тощо.

Мета метаданих — «полегшити пошук, оцінку, отримання та використання» ресурсів. У випадку освітніх ресурсів метадані також мають сприяти «спільному використанню та обміну навчальними об'єктами, дозволяючи розробляти каталоги та переліки, враховуючи різноманітність культурних і мовних контекстів, у яких навчальні об'єкти та їх метадані будуть використані» [19].

Більшість обчислювальних систем мають певну форму довготривалого зберігання даних, які можна досліджувати для отримання доказів. Хоча це не обов'язково для кожної системи, зазвичай таке зберігання організоване у вигляді файлів, каталогів і метаданих. Метадані включають всі дані у файловій системі, які описують структуру та атрибути файлів і каталогів, такі як мітки часу, інформацію про контроль доступу, розмір файлу, а також дані про місце розташування та збір файлів або каталогів у файловій системі [8].

Спочатку метадані файлової системи не були розроблені для реконструкції подій, які відбулися в системі. Першим, хто використовував такі дані для моніторингу загроз, був Андерсон в 1980 році. Він запропонував використовувати

записи System Management Facilities (SMF), які зберігалися на мейнфреймах, таких як IBM OS/360.

У 1960-х роках більшість обчислювальних завдань виконувалися на мейнфреймах, і OS/360 була однією з домінуючих операційних систем. Інформація, що зберігалася на дисках серверів, описувала всю пакетну роботу користувача. Дані для завдань надходили з перфокарт або стрічкових носіїв. Записи, що зберігали інформацію про пакетні завдання, описували різні аспекти завдання, включаючи дані користувача, які можна було б розглядати як файли. Це включало тип файлу, мінімальний і максимальний розмір, час створення, доступу та модифікації, а також інформацію про виконання завдання, таку як час виконання, тривалість і використані ресурси. Порівняно з метаданими сучасних систем, тоді було більше інформації, корисної для судово-медичних цілей. Ця інформація досі присутня в сучасних системах, але зазвичай не записується або лише тимчасово, наприклад у файловій системі rgs [8].

Multics була першою операційною системою, яка забезпечувала ієрархічну файлову систему, що стала основою для більшості сучасних файлових систем. У документі про дизайн файлової системи Multics описують потребу користувачів зберігати свої дані в обчислювальному середовищі, а не на зовнішніх носіях, таких як перфокарти чи стрічки. Користувач мав повний контроль і право власності на свої дані, включаючи метадані. Було сформовано цілі дизайну, які передбачали зберігання інформації, якою рідко користуються, на повільніших пристроях, легкий доступ до інформації за потреби, а також доступність інформації для інших користувачів на контрольованій основі. Щоб визначити частоту використання інформації, було запропоновано таке рішення як мітки часу доступу. Час модифікації та створення був потрібний для системи резервного копіювання, яка закріплювала б новостворені та змінені файли для резервного копіювання на стрічку. Ця мотивація лягла в основу використання часу MAC (англ. «Modified, Accessed, Changed»). Для надання доступу іншим користувачам автори документу запропонували включити списки контролю доступу та дозволи

для кожного файлу. Усі інші метадані стосувалися фактичного розміщення файлу на диску [8].

Пізніше для UNIX було створено файлову систему, на яку вплинула операційна система Multics. Метадані для файлів зберігалися в індексному дескрипторі і включали розташування та розмір файлу, його тип (каталог чи файл), три мітки часу та інформацію про контроль доступу, що складалася з ідентифікатора користувача та групи, а також бітів захисту, які використовуються всіма сучасними варіантами UNIX.

MS-DOS, що з'явилася на початку 1980-х років, використовувала файлову систему FAT, яка відстежувала тип файлу, розмір, розташування та мітки часу. Простір, зарезервований для міток часу, варіювався від 2 до 4 байтів, що призводило до різного ступеня деталізації. Наприклад, час доступу вимірювався лише днями. Оскільки DOS не мала поняття про користувача, жодна інформація про користувача чи дозволи не зберігалася у FAT. Windows спочатку успадкувала файлову систему FAT, але з часом її обмеження стали очевидними, і була представлена NTFS, яка містить детальну інформацію про користувача та дозволи, а також мітки часу змінених, доступних, створених і змінених файлів [8].

1.4 Проблема класифікації метаданих

Метадані відіграють критичну роль у цифровій криміналістиці. Вони надають цінну інформацію про дані, не вимагаючи доступу до самих файлів, що є особливо важливим у контексті правових і етичних обмежень.

Для критично важливих систем може бути доцільним регулярно створювати моментальні знімки всього стану системи або її значних частин. Однак у більшості випадків обмежена ємність пам'яті та вимоги до продуктивності роблять це неможливим. Тому необхідно розглянути зменшення обсягу інформації та частоти записів, використовуючи вже існуючі механізми системи. Метадані

надають ключову інформацію про операції з файлами в системі. Зазвичай операційна система, механізм завантаження, виконувани файли програм, дані конфігурації, інформація про користувача та дані додатків зберігаються у файлах. Аналіз файлів дозволяє зрозуміти, які дії відбувалися в системі: доступ до файлу може вказати на виконання програм, а модифікація файлу – на зміни в системі. Хоча операції з файлами є лише частиною стану системи, вони можуть надати важливу інформацію для експерта із цифрової криміналістики [8].

Метадані файлів є логічним місцем для запису важливої інформації про стан системи з мінімальними обчислювальними витратами. Вони можуть включати такі атрибути, як мітки часу, інформацію про контроль доступу, розмір файлу, а також дані про місцезнаходження файлів або каталогів у файловій системі. Метадані можна розглядати як характеристики цифрового об'єкта, що мають значення для цифрової криміналістики [20].

Зберігання такої важливої інформації як метадані файлової системи для цифрової криміналістики має кілька переваг:

- Автоматичний збір та збереження: вся інформація, доступна системі, автоматично записується.
- Відсутність додаткових витрат на налаштування: інформація збирається без потреби в додаткових механізмах реєстрації.
- Безпосередня доступність: інформація зберігається безпосередньо з об'єктом інтересу, що виключає потребу в співвіднесенні різних системних журналів.
- Захист від фальсифікації: інформацію важче підробити, ніж файл. Якщо доступ до необробленого диска заборонено операційною системою, записана інформація захищена від усіх користувачів. Навіть якщо доступ дозволено,

зловмиснику доведеться пройти через файловою систему, щоб змінити або видалити дані, не зруйнувавши при цьому критично важливі для системи дані.

Таким чином, метадані у файлових системах стають цінним інструментом для цифрової криміналістики, забезпечуючи надійний спосіб фіксації важливих змін у системі [8].

Прикладом підтвердження важливості метаданих у цифровій криміналістиці може слугувати метод виявлення матеріалів сексуального насильства над дітьми (CSAM), який базується на метаданих [21]. У цьому випадку метадані суттєво прискорюють процес виявлення та блокування такого нелегального контенту.

У розглянутому методі, використання метаданих дозволяє обійти етичні та юридичні обмеження, які виникають при зборі та обробці зображень CSAM для навчання моделей машинного навчання. Метадані, як «дані про дані», не є безпосереднім записом злочину, тому їх використання не потребує дотримання суворих правових норм. Це робить їх особливо корисними для швидкого реагування в розслідуваннях [21].

Застосування моделей машинного навчання, навчених на метаданих, може бути надзвичайно ефективним. Наукові дослідження демонструють, що структури виявлення, які використовують шляхи до файлів замість самих файлів, досягають високої точності. Наприклад, у випадку структури виявлення CSAM, точність досягає значення 0.97, що демонструє стабільну поведінку під час супротивних атак. Це значно підвищує ефективність роботи правоохоронних органів, забезпечуючи швидке і точне виявлення нелегального контенту з мінімальною кількістю хибно позитивних результатів (0.002) навіть при різноманітних розподілах даних [21].

Таким чином, метадані є цінним інструментом у цифровій криміналістиці, забезпечуючи ефективність та оперативність розслідувань, особливо у випадках,

що потребують негайних дій. Вони допомагають аналітикам швидко ідентифікувати та оцінювати дані, що дозволяє оперативно приймати обґрунтовані рішення, що робить їх невід'ємною складовою сучасних підходів до розслідування цифрових злочинів.

Висновки до розділу 1

Було розглянуто основи цифрової криміналістики, підкреслюючи важливість правильного збору, збереження та аналізу цифрових доказів. Успішне проведення розслідувань значною мірою залежить від того, наскільки ретельно і точно здійснюється процес обробки цифрових даних. Цей розділ наголошує на важливості використання метаданих, які можуть слугувати неочевидними, але важливими доказами, що допомагають у веденні розслідування, реконструкції подій і підтвердженні фактів у судових справах.

Цифрова криміналістика охоплює широкий спектр дисциплін, від інформатики до права, що вимагає від фахівців глибокого розуміння різних аспектів збору та аналізу цифрових доказів. Особливу увагу приділено метаданим, які можуть надати цінну інформацію щодо хронології подій, доступу до файлів і потенційних маніпуляцій з даними. Метадані є ключовим елементом у визначенні автентичності та достовірності цифрових доказів, що є важливим для правосуддя.

Розділ також підкреслює необхідність належної підготовки і навчання фахівців, що працюють з цифровими доказами, оскільки неправильне використання інструментів чи непрофесійні дії можуть призвести до помилкових висновків. Таким чином, для забезпечення справедливого судочинства важливо забезпечити високий рівень компетентності у сфері цифрової криміналістики.

2 КЛАСИФІКАЦІЯ ЦИФРОВИХ ДОКАЗІВ ІЗ ВИКОРИСТАННЯМ МАШИННОГО НАВЧАННЯ

2.1 Класифікація цифрових доказів

Цифрові докази у цифровій криміналістиці – це будь-яка інформація, що зберігається або передається в електронному вигляді, яка може бути використана в суді. Вони включають файли, електронні листи, журнали системи, історію браузера, метадані та інші форми даних, які можуть допомогти у розслідуванні злочинів.

Цифровий доказ – це, як правило, абстракція якогось цифрового об'єкта чи події. Коли людина дає вказівку комп'ютеру виконати таке завдання, як надсилання електронної пошти, результати діяльності створюють залишки даних, які дають лише часткове уявлення про те, що сталося. Лише деякі результати діяльності, такі як повідомлення електронної пошти та журнали сервера, залишаються, щоб дати нам часткове уявлення про те, що сталося. Крім того, використання криміналістичного інструменту для відновлення видаленого файлу з носія включає кілька рівнів абстракції від магнітних полів на диску до літер і цифр, які можна спостерігати на екрані. Таким чином, потрібна інформація отримується не з фактичних даних, а лише із представлення, і кожен рівень абстракції може вносити помилки. Ця ситуація подібна до традиційного огляду місця злочину. У справі про вбивство можуть бути підказки, які можна використати для реконструкції подій, як-от складання пазла. Проте всі частини пазла недоступні, тому створити повну реконструкцію злочину неможливо [22].

Цифрові докази зазвичай непрямі, що ускладнює приписування діяльності комп'ютера особі. Тому цифрові докази можуть бути лише одним із компонентів надійного розслідування. Якщо справа базується на одній формі чи джерелі цифрових доказів, таких як позначки дати й часу на комп'ютерних файлах, тоді справа є неприйнятно слабкою. Без додаткової інформації можна обґрунтовано

стверджувати, що в той час комп'ютером користувався хтось інший. Наприклад, механізми захисту паролем на деяких комп'ютерах можна обійти, і багато комп'ютерів не потребують пароля, що дозволяє будь-кому використовувати їх. Подібним чином, якщо відповідач стверджує, що деякі цифрові докази, які виправдовують, не були зібрані з однієї системи, це вплине лише на слабку справу, яка не має підтверджуючих доказів провини з інших джерел [22].

Той факт, що цифровими доказами можна так легко маніпулювати або знищити, створює складні виклики для цифрових дослідників. Цифрові докази можуть бути змінені або стерті або навмисно зловмисниками, або випадково під час збирання, не залишаючи жодних очевидних ознак спотворення.

Проте, цифрові докази мають кілька функцій, які пом'якшують цю проблему:

- Цифрові докази можна точно скопіювати, а копію можна перевірити так, ніби це оригінал. Під час роботи з цифровими доказами загальноприйнятою практикою є перевірка копії, що дозволяє уникнути ризику зміни або пошкодження оригінального доказу.
- За допомогою відповідних інструментів дуже легко визначити, чи цифрові докази були змінені чи підроблені, порівнявши їх з оригінальною копією.
- Цифрові докази важко знищити. Навіть коли файл «видалено» або жорсткий диск відформатовано, цифровий доказ можна відновити.
- Коли злочинці намагаються знищити цифрові докази, копії та пов'язані з ними залишки можуть залишатися в місцях, про які вони не знали [29].

Цифрові докази у цифровій криміналістиці є важливими інструментами для розслідування злочинів. Їх можна класифікувати різними методами залежно від їхніх характеристик, походження та ролі у кримінальному розслідуванні. Нижче наведено деякі з основних підходів до класифікації цифрових доказів:

1. Класифікація за типом даних: цей підхід розрізняє цифрові докази на основі їхньої природи та формату. Він включає наступні категорії:

- Файлові системи та метадані: дані, збережені у файлових системах, такі як метадані файлів, журнали доступу, а також структури каталогів.
- Файли документів: текстові документи, електронні таблиці, презентації та інші файли, створені користувачем.
- Мультимедійні файли: зображення, аудіо- та відеофайли.
- Мережеві журнали: дані мережевих з'єднань, журнали доступу до мережевих ресурсів, файли логів веб-серверів.
- Комунікаційні дані: електронні листи, повідомлення в соціальних мережах, чати, SMS-повідомлення [22].

2. Класифікація за походженням: цей метод класифікує докази залежно від джерела, з якого вони були отримані:

- Локальні системи: дані, зібрані з персональних комп'ютерів, мобільних пристроїв та інших локальних систем користувачів.
- Мережеві ресурси: дані, отримані з серверів, мережевих сховищ, хмарних сервісів та інших віддалених ресурсів.
- Периферійні пристрої: дані, зібрані з периферійних пристроїв, таких як принтери, сканери, зовнішні жорсткі диски та USB-накопичувачі.

3. Класифікація за функціональністю: цей підхід розрізняє докази залежно від їхнього функціонального призначення та ролі у розслідуванні:

- Прямі докази: дані, які безпосередньо підтверджують факт злочину, наприклад, електронні листи з погрозами або відеозаписи злочину.
- Непрямі докази: дані, які допомагають встановити контекст або підтверджують інші докази, наприклад, журнали доступу або метадані файлів.
- Підтверджуючі докази: дані, які підтримують або спростовують інші докази, такі як журнали системних подій або записи аудиту.

4. Класифікація за доступністю та складністю відновлення: цей метод класифікує докази на основі їхньої доступності та методів, необхідних для їхнього відновлення:

- Доступні дані: дані, які легко доступні та можуть бути зібрані без спеціальних інструментів.
- Видалені або приховані дані: дані, які були видалені або приховані та потребують спеціальних криміналістичних інструментів для їхнього відновлення, наприклад, відновлення видалених файлів або даних із відформатованих дисків [22].
- Зашифровані дані: дані, які були зашифровані та потребують спеціальних методів для їхнього розшифрування та аналізу.

Ці методи класифікації допомагають цифровим криміналістам структурувати та аналізувати великі обсяги даних, отриманих під час розслідувань, забезпечуючи більш ефективне виявлення та використання цифрових доказів.

З розвитком технологій з'явилися нові можливості для кримінальних і комерційних розслідувань. Використання цифрових доказів дозволяє слідчим відстежувати історію транзакцій, повідомлень та інших форм цифрових медіа, пов'язаних з місцезнаходженням, банківськими рахунками, паспортами та іншими ідентифікаторами. Правоохоронні органи можуть відслідковувати електронні сліди злочинців через журнали дій і надавати докази для їх засудження [23].

Однак, прийняття цифрових доказів у суді часто викликає труднощі. Для того щоб цифрові докази були прийнятними, вони повинні відповідати критеріям автентичності та надійності. Це вимагає дисциплінованого і науково обґрунтованого підходу до їх аналізу, оскільки дані можуть бути легко змінені.

У судочинстві критично важливо оцінювати якість і достовірність будь-яких доказів, щоб уникнути необґрунтованих рішень. Точність судової експертизи завжди викликала занепокоєння, і тому необхідно встановити чіткі наукові стандарти для перевірки валідності та надійності методів судової експертизи. Методи, яким бракує відповідної автентифікації та надійної статистичної основи, неприйнятні в суді.

Цифрові докази часто стикаються з труднощами у відповідності стандартам наукових критеріїв у судах. Недовіра до процесу цифрової криміналістики та відсутність встановлених правил для оцінки дають адвокатам можливість оскаржувати докази в суді. Вони можуть знайти лазівки в процесі збору та порівняння доказів, створюючи розумні сумніви щодо їх точності та достовірності. Критики також ставлять під сумнів валідність цифрових криміналістичних методів і засобів [23].

Таким чином, класифікація та оцінка цифрових доказів є складним, але необхідним процесом для забезпечення справедливості у кримінальних справах. Завдяки належним науковим методам і стандартам цифрові докази можуть стати потужним інструментом у розслідуванні та судочинстві.

2.2 Існуючі методи аналізу цифрових доказів на основі машинного навчання

Існуючі традиційні цифрові криміналістичні методології та інструменти для дослідження даних викликають перерву в критичних за часом експериментах. Вони непридатні для збільшення щільності зберігання та обробки величезних обсягів даних у цифрових криміналістичних лабораторіях, і вони потенційно можуть порушувати правові обмеження щодо обшуку та вилучення. Коли мова йде про кримінальне правопорушення чи напад, смерть, викрадення, зникнення безвісти тощо, швидке виявлення, розслідування та представлення доступного мультимедійного вмісту, особливо на місці, розглядається як критичне завдання для всього дослідження. Крім того, кількість мультимедійних даних, необхідних для аналізу, швидко збільшується через збільшення щільності зберігання. Сучасні підходи не мають ресурсів, щоб встигати за такою швидкістю. Вилучені пристрої, які потребують криміналістичного аналізу, не є чимось незвичайним для цифрових криміналістичних лабораторій, коли вони мають кілька місяців [24].

Традиційні підходи до аналізу даних все більше не справляються з обсягом та складністю інформації, що підлягає обробці. У зв'язку з цим, методи на основі машинного навчання пропонують потужні інструменти для автоматизації та підвищення точності класифікації цифрових доказів.

Машинне навчання дозволяє системам вчитися на прикладах та вдосконалювати свої алгоритми на основі нових даних, що значно підвищує ефективність та адаптивність процесу розслідування. Цей підрозділ розглядає існуючі методи класифікації цифрових доказів, що використовують алгоритми машинного навчання, оцінює їх ефективність, переваги та недоліки, а також їхню роль у сучасній цифровій криміналістиці.

1. Метод на основі нейронних мереж для ідентифікації файлів, які були змінені комп'ютерною програмою [25]

Нейронні мережі стали потужним інструментом в цифровій криміналістиці, особливо для класифікації файлів, які були змінені певною комп'ютерною програмою. Метод базується на аналізі цифрових доказів, які залишаються програмою при доступі, оновленні, модифікації або видаленні файлів файлової системи. Виявлення набору файлів, які були змінені в результаті інциденту, може значно полегшити процес реконструкції подій, хоча цей процес розглядається як можливий напрямок майбутніх досліджень.

Метод на основі нейронних мереж використовує машинне навчання для класифікації цифрових доказів. Цей підхід базується на аналізі цифрових слідів, залишених програмами при доступі, оновленні, модифікації або видаленні файлів у файлової системі. Основні джерела даних для аналізу включають метадані файлової системи, журнали аудиту подій та інформацію з реєстру.

Для тренування моделі збираються дані з різних джерел: метадані файлової системи, записи системного журналу подій, та інформація з реєстру. Комбінування характеристик з цих трьох джерел дозволяє компенсувати відсутність чи пошкодження даних в одному з джерел.

Цей метод застосовується для ідентифікації файлів, які були змінені певною комп'ютерною програмою, що допомагає у відтворенні подій під час цифрових криміналістичних розслідувань. Використання нейронних мереж дозволяє автоматизувати процес аналізу та зменшити часові витрати на класифікацію цифрових доказів.

2. Метод для автоматизованої категоризації цифрових медіа [26]

Цей метод представляє новаторський підхід до цифрової криміналістики. Він об'єднує принципи цифрової криміналістики та машинного навчання для автоматизації категоризації цифрових медіа, зосереджуючись на швидкому виявленні релевантних доказів для прискорення розслідувань.

Методологія включає кілька основних етапів. Спочатку відбувається процес створення точних копій цифрових пристроїв для подальшого аналізу. Далі експерти займаються виділенням специфічних ознак, які можуть бути пов'язані з розслідуваним злочином. Ознаки можуть включати історію браузера, журнали подій системи, записи дзвінків тощо.

Наступними важливим етапами буде визначення контексту та пріоритету – оцінка важливості різних ознак та їх зв'язку з розслідуваним злочином. І власне класифікація даних – застосування алгоритмів машинного навчання для класифікації пристроїв за ступенем їхньої важливості для розслідування.

У статті розглядаються два приклади застосування методології:

1) Порухення авторських прав – було використано набір даних, що містить 13 цифрових медіа та 45 ознак, пов'язаних із порушенням авторських прав. Результати показали, що понад 93% пристроїв були правильно класифіковані.

2) Обмін дитячою порнографією – використовувався набір даних з 23 мобільних телефонів та 23 ознак, пов'язаних із обміном дитячою порнографією. Усі телефони були правильно класифіковані, що демонструє високу ефективність методології.

Метод класифікації на основі машинного навчання значно покращує ефективність цифрової криміналістики, дозволяючи швидко і точно виявляти релевантні цифрові докази, а саме цифрових медіа. Це сприяє зменшенню навантаження на аналітиків і підвищенню ефективності розслідувань.

3. Метод ідентифікації типу файлів [27]

Метод ідентифікації типу файлів, описаний у науковій статті, використовує методи машинного навчання та штучний інтелект в обчисленнях для вибору та класифікації ознак файлів. Основна мета цього методу полягає в підвищенні точності ідентифікації типу файлів у цифровій криміналістиці, навіть якщо

злочинці намагаються змінити розширення або сигнатури файлів, щоб приховати їх справжній зміст.

Методологія включає три основні етапи:

1) Етап виділення ознак. Використовується розподіл частот байтів (англ. Byte Frequency Distribution). Кількість появ кожного значення байта у файлі підраховується та створюється масив з елементами від 0 до 255. Кожен елемент нормалізується шляхом ділення на максимальне значення. Таким чином, кожен файл представляється набором з 256 ознак. Метою цього етапу є отримання унікального набору ознак для кожного типу файлів, що дозволяє надійно визначати тип файлу незалежно від його розширення чи зміненої сигнатури.

2) Етап вибору ознак. Застосовується «генетичний» алгоритм для вибору релевантних ознак з 256 можливих. Кожне рішення (хромосома) представляється бінарним вектором, де 1 вказує на вибір ознаки, а 0 — на відхилення. Оцінка кожного рішення проводиться за допомогою функції, яка використовує алгоритм вибору ознак на основі кореляції. Цей етап проводиться для того, щоб вибрати найважливіші ознаки, які мають найбільший вплив на точність класифікації типу файлів.

3) Етап класифікації. Використовуються різні алгоритми машинного навчання для класифікації типів файлів, включаючи: Деревя рішень (DT), Метод опорних векторів (SVM), Штучні нейронні мережі (ANN), Логістичну регресію (LR) і Метод k-найближчих сусідів (KNN).

У дослідженні була проведена оцінка ефективності методології шляхом застосування її до набору даних з різними типами файлів. Результати показали, що метод нейронних мереж забезпечує найвищу точність класифікації у порівнянні з іншими алгоритмами. Навіть у випадках, коли розширення файлів були змінені або сигнатури були модифіковані, методологія змогла точно визначити типи файлів з високою точністю.

Метод ідентифікації типу файлів на основі методів машинного навчання та штучного інтелекту значно підвищує ефективність і точність цифрової

криміналістики. Використання генетичних алгоритмів для вибору ознак та різноманітних алгоритмів машинного навчання для класифікації дозволяє надійно ідентифікувати тип файлів, навіть якщо їхні розширення або сигнатури були змінені. Це робить методологію особливо цінною в умовах навмисного приховування або модифікації цифрових доказів.

2.3 Ефективність існуючих методів аналізу цифрових доказів

Порівняння методів наведено нижче у порівняльній Таблиці 2.1 - Переваги й ефективність методів та їх недоліки й обмеження.

Таблиця 2.1 - Переваги й ефективність методів та їх недоліки й обмеження

Метод	Переваги й ефективність	Недоліки й обмеження
<p>Метод на основі нейронних мереж для ідентифікації файлів, які були змінені комп'ютерною програмою</p>	<p>Висока автоматизація: дозволяє автоматизувати процес класифікації, що значно знижує навантаження на експертів.</p> <p>Компенсація відсутніх даних: використання характеристик з трьох різних джерел забезпечує високу надійність даних, навіть якщо деякі з них відсутні або пошкоджені.</p> <p>Гнучкість моделі: модель може бути налаштована під специфічні завдання шляхом коригування параметрів нейронної мережі.</p>	<p>Висока похибка: найкраща отримана похибка у дослідженні склала 10.07%, що вважається досить високим показником для практичного застосування.</p> <p>Накладання цифрових доказів: на одній файлової системі можуть працювати кілька програм, що ускладнює розрізнення їхніх цифрових доказів.</p> <p>Вимоги до обчислювальних ресурсів: тренування нейронної мережі вимагає значних обчислювальних ресурсів і великих обсягів даних.</p>

Продовження таблиці 2.1

Метод	Переваги й ефективність	Недоліки й обмеження
<p>Метод для автоматизованої категоризації цифрових медіа</p>	<p>Автоматизація процесів: методологія знижує обсяг рутинної роботи, яку виконують слідчі, дозволяючи їм зосередитися на більш складних та критичних аспектах розслідування.</p> <p>Робота в реальному часі: швидка обробка та категоризація даних в реальному часі, що є критичним у термінових розслідуваннях.</p> <p>Висока точність класифікації: дозволяє мінімізувати кількість хибно-позитивних та хибно-негативних результатів, досягає точності понад 93% у випадках порушення авторських прав та 100% у випадках обміну дитячою порнографією, що підтверджує її високу ефективність.</p> <p>Контекстуальна оцінка: Методологія враховує контекст та пріоритети, що дозволяє більш точно визначати релевантність та важливість цифрових доказів.</p>	<p>Вимоги до обчислювальних ресурсів: використання складних алгоритмів машинного навчання може вимагати значних ресурсів.</p> <p>Необхідність якісних навчальних даних: недостатня або неякісна навчальна вибірка може знизити точність класифікації.</p> <p>Залежність від початкових налаштувань: початкові налаштування та параметри моделей можуть сильно впливати на результати класифікації.</p>

Кінець таблиці 2.1

Метод	Переваги й ефективність	Недоліки й обмеження
Метод ідентифікації типу файлів	<p>Висока точність ідентифікації: методологія показала високу точність навіть у випадках, коли розширення файлів були змінені або сигнатури були модифіковані, що підвищує надійність методу в реальних умовах.</p> <p>Адаптивність та гнучкість: метод здатний адаптуватися до різних типів файлів та сценаріїв використання. Це досягається за рахунок використання багатьох ознак та потужних алгоритмів класифікації. Застосування генетичних алгоритмів дозволяє автоматично вибирати найбільш релевантні ознаки для конкретного набору даних, що підвищує ефективність класифікації.</p> <p>Надійність при приховуванні або модифікації файлів: методологія показала свою ефективність у випадках навмисної зміни розширень або сигнатур файлів, що робить її особливо корисною в умовах навмисного приховування цифрових доказів.</p>	<p>Висока обчислювальна складність: застосування генетичних алгоритмів та алгоритмів машинного навчання може вимагати значних обчислювальних ресурсів, що може бути проблемою при обробці великих обсягів даних або в умовах обмежених ресурсів.</p> <p>Необхідність попереднього навчання: потрібні великі та репрезентативні набори даних для навчання моделей, що може бути складно забезпечити в деяких випадках.</p> <p>Складність налаштування та підтримки: методологія вимагає спеціалізованих знань для налаштування і підтримки, що може обмежувати її використання.</p>

Порівняльна таблиця дозволяє зробити висновок, що кожен з розглянутих методів класифікації цифрових доказів на основі машинного навчання має свої унікальні переваги та недоліки.

2.4 Визначення пріоритетних підозрілих файлів

Як відомо, традиційно найкращим автоматичним рішенням під час цифрового криміналістичного дослідження для пошуку файлів, які є відомими і містять зловмисний характер, це методи білого чи чорного списків на основі хешування. Відповідно якщо не відбувається потрібне зіставлення, то важливого для розслідування цифрового доказу таким шляхом не буде виявлено, і тому дослідникам доведеться вручну виконати пошук за ключовими словами або фільтрувати файли за певними критеріями. Тому використання хеш-зіставлення для виявлення незаконних файлів може виявити лише ті самі атрибути файлів або зі незначними змінами, тобто за допомогою приблизного зіставлення.

Це дослідження представляє підхід, який дозволяє здійснити швидкий і спрощений пошук підозрілих для криміналістичного розслідування файлів шляхом навчання моделей машинного навчання. Пошук відбувається безпосередньо експертами-слідчими на основі проведеного методу класифікації всіх файлів, де визначені як підозрілі файли класифікатором, отримують мітку «1» у базі файлів досліджуваної системи й окремо зберігаються від наявної бази файлів.

Процес класифікації відбувається безпосередньо на основі метаданих файлів. Наявні так звані незаконні метадані файлів, позначаються як підозрілі для моделі класифікації і заносяться до чорного списку. Ці характеристики файлів в чорному списку можна використовувати для навчання класифікаторів для виявлення попередніх невідомих підозрілих файлів у нових випадках.

Основна ідея методу полягає в тому, що особливо в чутливих до часу випадках цифрової криміналістики потрібно якнайшвидше віднайти необхідні цифрові докази. В таких злочинах як дитяча порнографія цифрові докази часто споріднені за неочевидними характеристиками, такими як час створення цифрового доказу, як наприклад фото чи відео знімки, які зазвичай створюються із дуже малою часовою різницею. Також глибина директорії, в якій знаходяться ці

цифрові докази, адже зазвичай фотознімки, які створені в один проміжок часу, зберігаються в одній директорії. Такі споріднені характеристики можна відслідкувати саме за допомогою метаданих цифрових файлів.

Тому описаний метод визначення підозрілих файлів базується на метаданих файлів, де при схожій спорідненості метаданих відомих незаконних файлів, визначається підозрілість досліджуваних цифрових доказів, де значення пріоритетності позначається мітками.

Модель на основі методів машинного навчання не може приймати остаточне рішення чи файл є цифровим доказом для певного криміналістичного розслідування, проте інтеграція такої моделі у поширені програмні засоби цифрової криміналістики може пришвидшити його хід, що завжди має цінне значення і буває критично необхідним.

2.5 Приклади застосування машинного навчання у цифровій криміналістиці

Зі зростанням кількості кіберзлочинів в останні роки, цифрова криміналістика стала ключовою сферою дослідження для отримання якісних доказів. Експерти цифрової криміналістики стикаються з викликами щодо збору та аналізу даних для реконструкції подій. Завдяки значній щоденній взаємодії між людьми, машинне навчання дає змогу дослідникам проводити розслідування ефективніше, використовуючи різноманітні алгоритми. Машинне навчання, як підгалузь штучного інтелекту, є науковою дисципліною, яка зосереджується на створенні комп'ютерних моделей і алгоритмів, здатних виконувати завдання без необхідності програмування, таких як навчання та тестування на наборах даних, що може бути корисним у розслідуваннях. Кожен алгоритм машинного навчання застосовується в певній сфері цифрової криміналістики залежно від його характеристик, таких як складність, обсяг даних, часові рамки, кореляція та послідовність.

Цифрові криміналісти використовують алгоритми машинного навчання для аналізу великих обсягів даних, збережених у різних середовищах і хмарних мережах. Ці дані можуть бути використані для прогнозування поведінки користувачів. Крім того, алгоритми машинного навчання можуть виконувати розпізнавання образів. Завдяки методам машинного навчання, слідчі застосовують набір правил і методів для виявлення шаблонів даних, що допомагають ідентифікувати потенційну злочинну діяльність. В цьому підрозділі розглядаються кілька алгоритмів, запропонованих для виявлення цифрових доказів і покращення процесу розслідування.

1) Дерево рішень (DT)

Було розроблено архітектурну структуру в дослідженні [28], яка об'єднує MapReduce, розподілену файлову систему Hadoop та алгоритм DT. Ця структура обробляє великі обсяги даних, які можуть бути зібрані та збережені. Вона складається з чотирьох етапів: захоплення мережевого трафіку, перетворення його у зрозумілий формат, фільтрація пакетів, аналіз даних на предмет шкідливих дій, а також представлення аналізу та візуалізація загроз. DT класифікує трафік як зловмисний або нешкідливий, що підвищує точність і ефективність на кожному етапі. Дослідження показали, що ця модель здатна виявити 99% зловмисного та нешкідливого трафіку.

У 2021 році було створено гібридний підхід для вирішення проблем із системою репутації IP. Вони поєднали різні методи аналізу даних, такі як машинне навчання, динамічний аналіз зловмисного програмного забезпечення та аналіз кіберзагроз. Використовуючи криміналістику великих даних, система може передбачити ймовірність конкретної атаки до її здійснення та класифікувати її за поведінковими характеристиками. Запропонована система була оцінена у порівнянні з існуючими системами репутації, використовуючи різні методи машинного навчання, такі як DT, SVM і Наївний Баєсів класифікатор (NB). DT

показав високі показники метрики повноти (англ. recall), F-міри та влучності (англ. precision) [29].

2) Метод k-найближчих сусідів (KNN)

Було представлено структуру для аналізу та виявлення DDoS-атак — розподілених атак на відмову в обслуговуванні, використовуючи алгоритми KNN і NB. Метод застосовує статистичні техніки для покращення продуктивності виявлення аномального мережевого трафіку. Алгоритм KNN базується на статистичних характеристиках наборів даних. У порівнянні з NB, алгоритм KNN демонструє кращі результати за точністю (англ. accuracy), повнотою (англ. recall) та влучністю (англ. precision) класифікації [30].

Також було запропоновано метод для визначення статі, який складається з трьох етапів: видобування даних, створення ознак та вибір найкращих класифікаторів.

Перший етап включає видобування ключових точок тіла з відео, далі — створення ознак тіла за допомогою OpenPose. OpenPose — це система оцінки пози, розроблена дослідниками з Університету Карнегі-Меллона (CMU), яка може виявляти та відстежувати людське тіло в режимі реального часу та точно визначати його позу в 3D-просторі. Вона добре відома тим, що є першою системою оцінки пози кількох людей у режимі реального часу, яка точно визначає ключові точки людського тіла, рук, обличчя та ніг (загалом 135 ключових точок) на окремих зображеннях [31].

Другий етап передбачає навчання чотирьох різних класифікаторів: KNN, Алгоритм адаптивного підсилювання (AB), Випадковий ліс (RF) та SVM.

Найбільш точним методом визначення статі, навіть у темряві, є RF, за яким слідує метод KNN [32].

2) Метод опорних векторів (SVM)

Сара Феррейра та інші запропонували метод для розрізнення підроблених і справжніх цифрових фотографій та відео [33]. Цей метод базується на алгоритмі SVM для вилучення ознак з даних, зібраних за допомогою дискретного перетворення Фур'є (DFT). Використовуючи бібліотеку Scikit-Learn у Python 3.9, обробка SVM дозволила створити модель класифікації для отриманих даних. Ця модель потім прогнозує автентичність фотографій у тестовому наборі даних. Було розроблено набір програм на Python для обробки ознак фотографій і вилучення кадрів з відео. Вони також використовувалися для створення моделі SVM, яка може класифікувати зображення. Алгоритм DFT-SVM використовувався для аналізу фотографій та відео, і результати показали, що він виконує аналіз швидше, ніж Згортова нейронна мережа (CNN). Низький час обробки та висока продуктивність методу DFT-SVM роблять його ідеальним інструментом для виявлення підробленого мультимедійного контенту на етапі цифрового судового аналізу. Результати дослідження були багатообіцяючими і відповідали процедурі попередніх досліджень.

3) Наївний Баєсів класифікатор (NB)

Антон Юдхана та інші проаналізували дані, зібрані з журнального мережевого трафіку, щоб встановити точність виявлення за типом DDoS-атак. За допомогою ПЗ Wireshark вони зібрали набори даних мережевого трафіку та виділили мережеві характеристики для виявлення шаблонів у даних. Після цього пакету вони виконали процедуру класифікації мережевих мереж, використовуючи NB, і навчили його за допомогою ANN, використовуючи кілька нейронів. Аналіз і тестування показали, що ANN досягла точності 95,2381%, тоді як NB мав точність

99,999%. Дослідники вважають, що використання ANN та NB в мережевій криміналістиці може сприяти підвищенню точності результатів під час розслідувань [34].

4) Метод k-середніх (K-Means)

Було розроблено структуру, яка використовує алгоритм K-Means для аналізу даних, зібраних з різних джерел. Дані, зібрані з різних видів аналізу кіберзлочинів, можна легко сортувати для полегшення виділення характеристик. На етапі кластеризації алгоритм виявляє взаємодії між характеристиками. Запропонований метод може надати дієві кроки для запобігання повторенню таких злочинів у майбутньому [35].

Науковець Руріаван та інші розробили систему, яка може ідентифікувати цифрові докази та класифікувати вміст носіїв інформації за допомогою алгоритму кластеризації K-Means. Дублюючий пристрій копіює та зберігає вказані користувачем дані. Запропонована система може використовуватися для відновлення файлів, збережених на носіях. Ця система допомагає судовим слідчим ефективніше підготувати відповідні докази [36].

5) Логістична регресія (LR)

Було визначено типи шкідливого програмного забезпечення, які найчастіше зустрічаються в ОС Windows і націлюються на реєстр. Шкідливе ПЗ може спричинити втрату дорогоцінного часу під час процесу розслідування. Вони надали цінну інформацію про те, як ці типи шкідливого ПЗ взаємодіють з реєстром. Дослідники протестували різні класифікатори, такі як ANN, DT і LR. Результати їх дослідження показали, що можна проводити аналіз цифрової криміналістики, використовуючи змінені часові позначки та методи машинного

навчання. Автори визначили 47 місць у реєстрі, які часто стають мішенню шкідливого ПЗ. Дослідники встановили, що алгоритм Дерево підсилення (англ. Boosted Tree) правильно класифікував понад 72% шкідливого ПЗ у їхньому дослідженні. Цей метод дозволяє слідчим швидко визначати, який тип шкідливого ПЗ присутній, а який – ні [37].

Також було запропоновано підхід міток, який можна використовувати для організації та аналізу електронної пошти. Цей метод може допомогти провести розслідування, пов'язані з незаконним використанням електронної пошти. Підхід реалізується в кілька етапів: обробка даних, яка включає видалення найчастіше повторюваних слів у реченні, таких як "ми", "є" тощо. Потім різні техніки машинного навчання витягують звичайні характеристики електронної пошти з шкідливих листів. За результатами експерименту алгоритм LR показав кращі результати у порівнянні з іншими техніками машинного навчання щодо точності та аналізу класифікації електронної пошти [38].

2.6 Обмеження машинного навчання у цифровій криміналістиці

Недостатня прозорість моделей і методології тестування є значною проблемою, яка впливає на розробку та впровадження моделей машинного навчання. Дослідницькі лабораторії часто створюють нові моделі, які можна швидко реалізувати в реальних умовах, але вони можуть виявитися невдалими в цих ситуаціях. Наявність інструментів та ресурсів для відтворення моделей може допомогти різним галузям та фахівцям швидше вирішувати свої проблеми, а також запобігати таким проблемам, як упередженість. Однак багато моделей машинного навчання повинні бути розроблені з необхідною прозорістю та методами тестування для забезпечення ефективної роботи судових експертів. Відсутність цих аспектів може завадити їм ефективно пояснювати різні результати своїх систем..

Однією з найбільших проблем глибоких алгоритмів навчання є їх інтерпретованість. Машинні моделі можуть бути дуже потужними, але також безсилими, якщо їх неможливо адекватно інтерпретувати. Тому важливо, щоб їх можна було застосовувати в реальних сценаріях.

Різні техніки, такі як DT та SVM, використовуються для аналізу та прогнозування анонімної поведінки користувачів у великих даних. Через складність ANN їм потрібні необхідні тренувальні дані для правильного виконання своїх функцій. Зі зростанням їх архітектури зростає і потреба в даних, що означає, що повторне використання даних не дасть бажаних результатів. Існуючі системи репутації можуть бути проблематичними через їх обмежену здатність виявляти нульові аномалії, залежність від внутрішніх джерел і високі витрати на управління. Відсутність джерел даних і якісних даних є суттєвою проблемою, яку потрібно вирішити. Хоча наявність достатньої кількості інформації іноді може бути рівнозначною її відсутності, надання неякісних даних може вплинути на точність моделі.

Переваги довіри до комп'ютерних алгоритмів численні. Люди значно виграли від їх здатності автоматизувати процеси та аналізувати великі обсяги даних. На жаль, алгоритми також можуть бути упередженими, оскільки вони створюються та навчаються людьми, що ускладнює усунення упередженості.

Висновки до розділу 2

У цьому розділі було детально розглянуто існуючі методи класифікації цифрових доказів, зокрема на основі їх типу, джерела, доступності та складності відновлення. Важливо зазначити, що методи класифікації за типом даних, згадані раніше, відіграють важливу роль у структуризації цифрових доказів.

Аналіз існуючих методів продемонстрував потенціал для автоматизації процесу ідентифікації та класифікації цифрових доказів за допомогою різних методів машинного навчання, що дозволяє швидше обробляти великі обсяги даних. Однак, важливо враховувати, що такі методи вимагають значних обчислювальних ресурсів та якісних навчальних даних для забезпечення високої точності класифікації.

Узагальнюючи, можна зазначити, що кожен з розглянутих методів класифікації цифрових доказів є важливим інструментом у цифровій криміналістиці, і їх комбінація може надати найкращі результати для розслідування злочинів. Подальші дослідження та розвиток технологій, таких як машинне навчання, матимуть значний вплив на покращення процесів збору, аналізу та класифікації цифрових доказів.

Подальший розвиток і використання нових методів класифікації цифрових доказів забезпечить більш глибоке розуміння та ефективне використання цифрових доказів у розслідуваннях, що в кінцевому підсумку сприятиме ефективнішому розв'язанню кримінальних справ.

3 РОЗРОБКА МЕТОДУ КЛАСИФІКАЦІЇ ЦИФРОВИХ ДОКАЗІВ НА ОСНОВІ МЕТАДАНИХ

3.1 Збір даних

Для створення методу класифікації цифрових доказів на основі метаданих цифрових файлів методами машинного навчання необхідно створити такий набір даних, що буде містити метадані файлів, які є притаманними звичайному цифровому носію, а також метадані файлів, що є підозрілими, і саме їх цей метод класифікації повинен віднайти.

У ході криміналістичного розслідування вилучають різноманітні електронні пристрої та цифрові носії, які можуть містити цифрові докази, необхідні для розслідування. Серед них можуть бути: комп'ютери та ноутбуки, мобільні телефони та планшети, зовнішні жорсткі диски та USB-накопичувачі, сервери, мережеве обладнання, хмарні сховища та інші цифрові носії.

Було прийнято рішення зосередитися на дослідженні комп'ютера з операційною системою Windows. Важливість цього вибору для створення набору даних для класифікації цифрових доказів було обґрунтовано за такими факторами:

1. Поширеність ОС Windows: Windows є однією з найпоширеніших операційних систем у світі, як у приватному, так і в корпоративному секторі. Це означає, що більшість комп'ютерних систем, які можуть бути об'єктами криміналістичних розслідувань, працюють під управлінням Windows.

2. Різноманітність типів даних: Windows підтримує широкий спектр форматів файлів і типів даних, від текстових документів до відео, аудіо, баз даних та спеціалізованих програмних файлів. Це дозволяє створити багатогранний набір даних, який охоплює різні аспекти можливих цифрових доказів.

3. Наявність необхідних метаданих: Windows зберігає детальні метадані для файлів, включаючи інформацію про створення, зміну, доступ і атрибути файлів. Це критично важливо для класифікації цифрових доказів, оскільки метадані можуть допомогти визначити час і спосіб використання файлів.

4. Інструменти аналізу: існує безліч інструментів для комп'ютерної криміналістики, спеціально розроблених для роботи з Windows, таких як FTK, EnCase, і Autopsy. Ці інструменти надають потужні можливості для збору, аналізу та відновлення даних з Windows-систем. Саме Autopsy буде використовуватися у подальшій розробці.

5. Легкість віртуалізації: Windows легко встановлюється і налаштовується у середовищах для створення віртуальних машин, таких як VirtualBox, що дозволяє створювати точні копії середовища для аналізу без ризику пошкодження оригінальних даних.

Для використання цифрових файлів ОС Windows як основу, було створено віртуальну машину за допомогою програмного забезпечення VirtualBox, що розроблене компанією Oracle Corporation.

Отже, за цифровий носій для дослідження взято віртуальний жорсткий диск, що є найоптимальнішим рішенням для створення набору даних. Віртуальний жорсткий диск є файлом, який імітує фізичний жорсткий диск і може містити операційні системи, програми, файли та інші дані, подібно до фізичних жорстких дисків.

Після встановлення віртуальної машини з ОС Windows було необхідно встановити на неї підозрілі файли для методу класифікації. Щоб отримати широкий спектр різних типів файлів та їх метаданих було обрано скористатися відкритим джерелом даних під назвою Digital Corpora [39] для дослідження. Цей

ресурс у вільному доступі розміщує різні дані саме для використання в освітніх дослідженнях комп'ютерної криміналістики. На віртуальну машину було завантажено близько 24000 файлів зі збірки Govdocs1.

The screenshot shows the Digital Corpora website. The main heading is "Govdocs1 – (nearly) 1 million freely-redistributable files". The text explains that in recent years, forensic research has involved the analysis of files or file fragments, and that this corpus was created to be freely available for research. It lists three ways the corpus is available: a set of 1000 directories, a set of 1000 ZIP files, and a tar file with 109,282 JPEG pictures. The right sidebar contains search bars, "Recent Posts" (including "CIRCL Forensics Exercises" and "Compiled bulk_extractor 2.0 ready for download"), and "Recent Comments" (including "Desmond on 2012 National Gallery DC Attack").

Рисунок 3.1 - Збір файлів

Для збору метаданих всіх цифрових файлів віртуального жорсткого диску було обрано скористатися провідною наскрізною платформою цифрової криміналістики з відкритим кодом — Autopsy. Цей інструмент надає швидке, ретельне й ефективне рішення для дослідження жорстких дисків. Для того, щоб здійснити аналіз віртуальної машини було необхідно спочатку конвертувати диск у відповідний формат для обробки ПЗ Autopsy. Було обрано зручний для цього завдання командний інтерфейс VBoxManage, якого надає ПЗ Oracle VM VirtualBox VirtualBox при встановленні за замовчуванням. Так виглядає виконана команда:

```
C:\Program Files\Oracle\VirtualBox>.\VBoxManage.exe" clonehd "D:\windows\windows\windows.vdi" D:\windows\windows\newdisk.img --format raw
0%..10%..20%..30%..40%..50%..60%..70%..80%..90%..100%
Clone medium created in format 'raw'. UUID: e8eed028-892b-4601-a77c-b982752af01b
C:\Program Files\Oracle\VirtualBox>
```

Рисунок 3.2 - Конвертація файлового образу у RAW формат

Потім виконується сканування файлового образу та збереження результатів, генеруючи звіт як показано на Рисунку 3.3 і його збереження у форматі CSV: Необхідний для дослідження збір даних виконано. Також додано стовпець Label для позначення файлів міткою «1», які вважаються підозрілими і надалі будуть зазначатися у дослідженні як нелегальні.

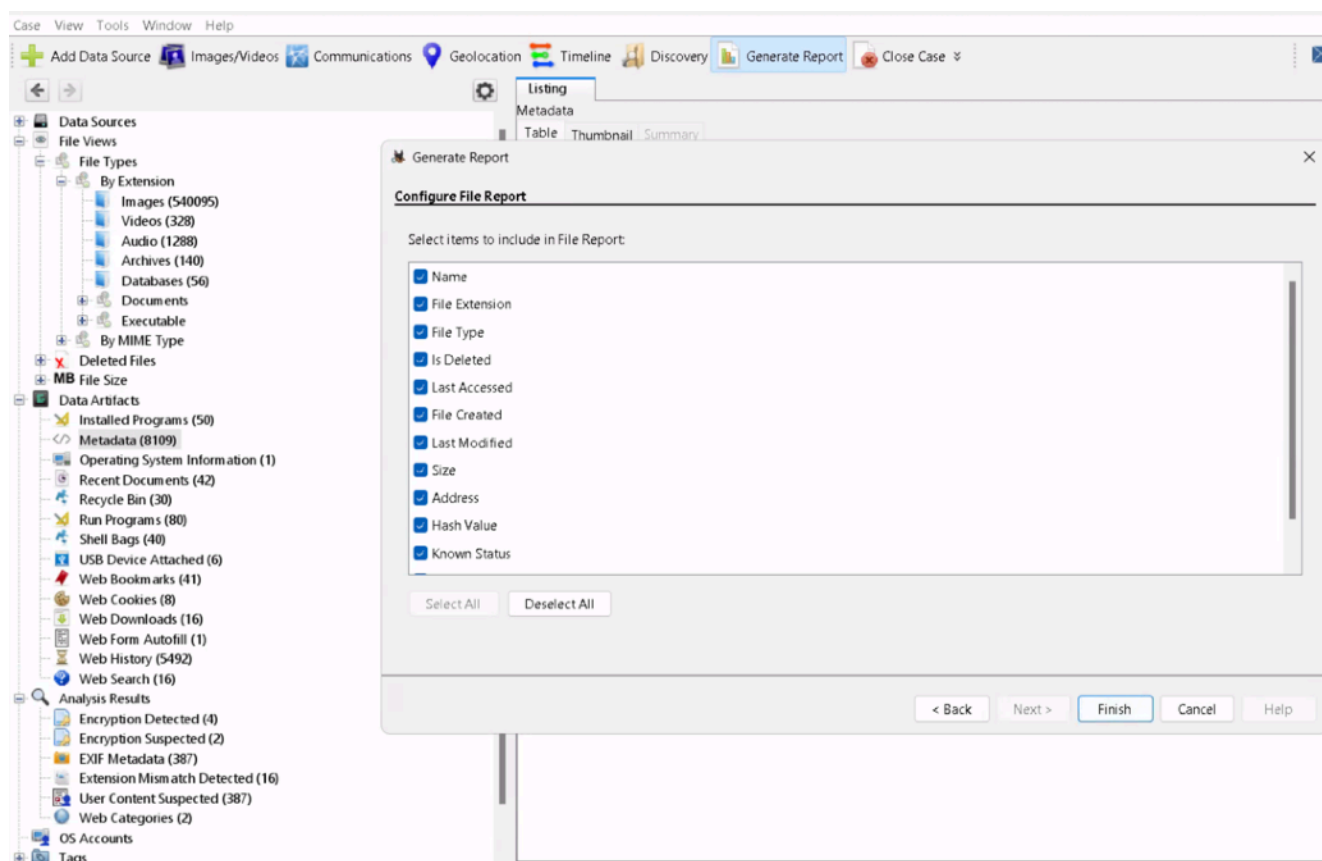


Рисунок 3.3 - Генерація звіту

3.2 Обробка даних

Перед створенням моделі для тренування і тестування методу класифікації, потрібно провести аналіз зібраних даних і відповідно підготувати їх для роботи. Розглянемо покрокові дії, що виконуються:

1) У програмному коді `preprocess.ipynb` мовою Python 3.10 завантажено бібліотеки `pandas`, `datetime` і з бібліотеки `scikit-learn` модуль `LabelEncoder` та файл метаданих файлів, що містить перелік таких стовпців зібраних даних: `Name`, `File`

Extension, File Type, Is Deleted, Last Accessed, File Created, Last Modified, Size, Address, Hash Value, Known Status, Permissions, Full Path, Label. Важливо зазначити, що Address містить числове значення глибини директорії, в якому міститься файл.

2) Проаналізовано кількість пустих значень стовпців датасету. Стовпець Hash Value містить понад 30% пустих значень, тому його було видалено. Інші стовпці мали не більш ніж 1% пустих значень, то було усунуто тільки їхні записи.

3) Видалено зайві для дослідження стовпці, такі як Known Status, Permissions, Is Deleted, Full Path, Name.

4) Змінено формат запису дати і часу із розширеного формату із зазначенням часового поясу — Центральноєвропейського літнього часу (CEST) до числового значення за допомогою функції `parse_datetime_with_timezone`. На Рисунку 3.4 продемонстровано приклад застосування даної функції:

```
def parse_datetime_with_timezone(timestamp_str):
    # Відокремлюємо дату та час від рядка та видаляємо часовий пояс
    timestamp_str_without_timezone = timestamp_str.split()[0] + " " + timestamp_str.split()[1]
    timestamp = pd.to_datetime(timestamp_str_without_timezone, format="%Y-%m-%d %H:%M:%S")
    unix_timestamp = timestamp.timestamp()
    return unix_timestamp

timestamp_str = "2024-05-17 05:32:22 CEST"
parsed_timestamp = parse_datetime_with_timezone(timestamp_str)
print(parsed_timestamp)

1715923942.0
```

Рисунок 3.4 - Функція обробки дати і часу даних

Функцію `parse_datetime_with_timezone` застосовано до стовпців Last Accessed, File Created, Last Modified.

5) Стовпці File Type та File Extension містять категоріальні змінні, тому їх потрібно перевести у числовий формат, а оскільки вони не мають природного

порядку, то можна використати LabelEncoder, де на Рисунку 3.5 зображено кодування цих стовпців і як в результаті обробки даних виглядають значення всіх стовпців:

```

1 label_encoder_extension = LabelEncoder()
2 label_encoder_type = LabelEncoder()
3
4 part_illegal['File Extension'] = label_encoder_extension.fit_transform(part_illegal['File Extension'])
5 part_illegal['File Type'] = label_encoder_type.fit_transform(part_illegal['File Type'])
6
7 print(part_illegal.head())

```

	Name	File Extension	File Type	Last Accessed	File Created	\
15	\$TxfLog.blf	1521	1	1.715924e+09	1.715924e+09	
26	BCD.LOG	1420	1	1.715924e+09	1.715924e+09	
27	BCD.LOG-slack	1421	1	1.715924e+09	1.715924e+09	
28	BCD.LOG1	1422	1	1.715924e+09	1.715924e+09	
29	BCD.LOG2	1424	1	1.715924e+09	1.715924e+09	

	Last Modified	Size	Address	Label
15	1.715900e+09	65536	32	0
26	1.715898e+09	21504	93	0
27	1.715898e+09	3072	93	0
28	1.715924e+09	0	94	0
29	1.715924e+09	0	95	0

Рисунок 3.5 - Обробка категоріальних змінних

б) Після обробки даних всі стовпці датасету містять числові значення. Змінений датасет зберігається в окремий файл формату CSV.

Створений набір даних загалом містить 300 000 записів із таким розподілом даних як зображено нижче на Рисунку 3.6:

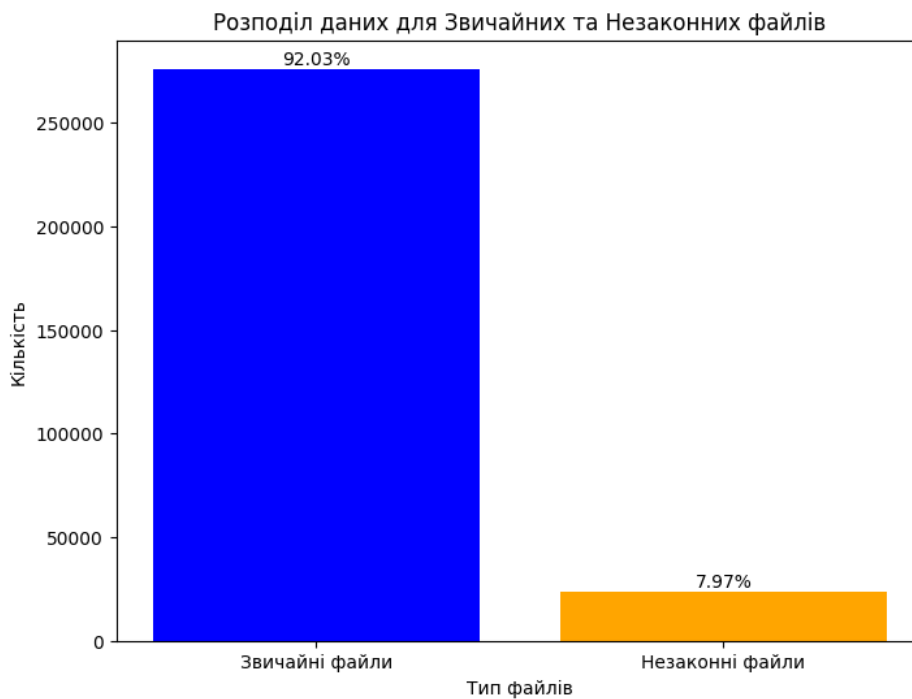


Рисунок 3.6 - Розподіл даних

Даний розподіл є доцільним, адже зазвичай кількість незаконних цифрових доказів підозрюваного є значно меншою за всі інші цифрові файли його комп'ютера чи іншого цифрового носія.

3.3 Створення моделі класифікації та аналіз результатів

Для створення моделі для задачі класифікації цифрових доказів використано підхід керованого машинного навчання, де правильними вихідним значення є мітки стовпчика Label.

Набір даних випадковим чином розділяється на тренувальний і тестовий набори, причому тестовий набір складає 20% від усіх даних.

Загалом, було випробувано чотири різні методи машинного навчання для задачі класифікації для порівняння їх ефективності:

- Логістична регресія (LR);
- Метод k-найближчих сусідів, де обране значення k=5 (KNN);
- Гаусів Наївний Баєсів класифікатор (GNB);
- Випадковий ліс (RF).

Для оцінки результатів окрім метрики точність (англ. accuracy) використовуються також влучність (англ. precision), повнота (англ. recall) та F1-міра (англ. F1-score). Оскільки кількість незаконних файлів значно менша за звичайні, то метрики точності може бути недостатньо, адже модель буде прогнозувати значення більшості класу для всіх прикладів і досягати високої точності класифікації, при цьому не враховуючи реальний баланс класів. Тому використання також інших метрик дає більш об'єктивну картину про ефективність моделей.

Таблиця 3.1 - Порівняння результатів класифікаторів

Метод / Метрика	Точність	Влучність	Повнота	F1-міра
Логістична регресія	0.98	0.80	0.99	0.89
Метод k-найближчих сусідів	0.97	0.80	0.95	0.87
Гаусів Наївний Баєсів класифікатор	0.98	0.80	0.99	0.89
Випадковий ліс	0.99	1.00	0.99	0.99

Судячи із результатів кожного класифікатора, можна зробити висновок, що Випадковий ліс найкраще класифікує цифрові докази на підозрілість і визначення їх пріоритету важливості для подальшого криміналістичного розслідування.

Додатково було побудовано для кожного методу ROC-криву, яка демонструє здатність класифікатора розрізняти між класами, незалежно від порогу

класифікації, та дозволяє зробити вибір між чутливістю та специфічністю залежно від конкретних потреб задачі для кожного із застосованих методів. Чим ближче до верхнього лівого кута ROC-кривої, тим кращою є модель.

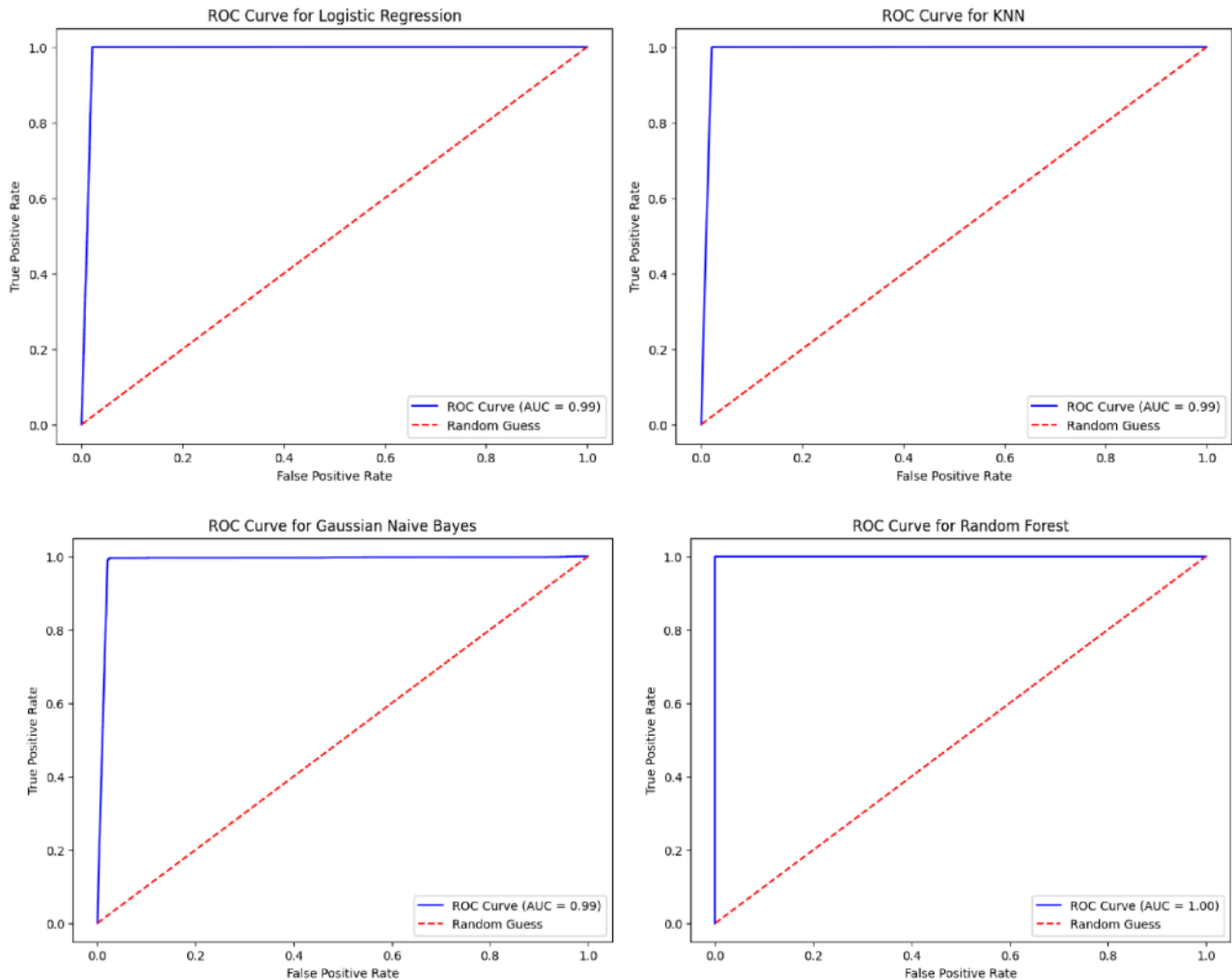


Рисунок 3.7 - ROC-криві класифікаторів

Після аналізу створених ROC-кривих видно, що кожний із методів машинного навчання добре класифікує цифрові докази.

3.4 Майбутні дослідження

Подальша робота для можливого розвитку та вдосконалення даного дослідження може зосереджуватися на таких напрямках:

1. Розробка додаткових сценаріїв. Під час даного дослідження використовувався сценарій, спрямований на виявлення та класифікацію файлів, які можуть використовуватися у цифровій криміналістиці у справах, що пов'язані, наприклад, з жорстоким поводженням з дітьми, де важливий акцент на таких метаданих як розмір файлу, глибина директорії, час створення тощо. Майбутні дослідження можуть передбачати розробку і оцінку більшої кількості сценаріїв розслідування. Наприклад, моделі для ідентифікації підозрілих файлів з різних джерел, таких як вкладення електронної пошти, хмарні облікові записи, USB-пристрої тощо, а також для виявлення сторонніх власників файлів.

2. Аналіз впливу розподілу даних. У майбутніх роботах рекомендовано створити більш різноманітні набори даних для тестування, що охоплюють широкий спектр співвідношень звичайних та незаконних файлів. Це дозволить краще зрозуміти вплив розподілу даних для дослідження і розробити стратегії для її подолання.

3. Адаптація моделей для реальних сценаріїв розслідування. Необхідно налаштувати моделі так, щоб вони краще відповідали потребам реальних криміналістичних розслідувань. Наприклад, важливо зосередитися на показниках запам'ятовування, а не лише на точності, оскільки пропуск будь-якого релевантного доказу є недопустимим. Це потребує налаштування моделей для забезпечення високої чутливості до підозрілих файлів.

Ці напрямки досліджень допоможуть удосконалити методи класифікації цифрових доказів, підвищити точність та ефективність розслідувань, а також сприятимуть розвитку нових підходів у галузі комп'ютерної криміналістики.

Висновки до розділу 3

У даній роботі було розроблено метод класифікації цифрових доказів на основі метаданих цифрових файлів за допомогою методів машинного навчання. На основі виконаного дослідження можна зробити наступні висновки:

- Обрання ОС Windows для створення набору даних є обґрунтованим вибором через її поширеність, різноманітність типів даних та наявність необхідних метаданих. Це дозволило створити релевантний набір даних, що охоплює різні аспекти можливих цифрових доказів. Встановлення віртуальної машини і використання джерела даних Digital Corpora для завантаження підозрілих файлів дозволили отримати великий та різноманітний набір даних.

- Було випробувано чотири різні методи машинного навчання: Логістична регресія, Метод k-найближчих сусідів, Гаусів Наївний Баєсів класифікатор та Випадковий ліс. Випадковий ліс показав найкращі результати за метриками точності, влучності, повноти та F1-міри. Порівняння результатів класифікаторів показало, що Випадковий ліс є найбільш ефективним для класифікації цифрових доказів.

- На основі отриманих результатів, майбутні дослідження можуть бути спрямовані на вдосконалення методу класифікації за рахунок використання реальних і різноманітних сценаріїв у цифровій криміналістиці для збору даних, а також навчання та адаптації до різних версій ОС і типів цифрових носіїв.

Загалом, розроблений метод класифікації цифрових доказів на основі метаданих продемонстрував свою ефективність та перспективність для використання у криміналістичних розслідуваннях. Впровадження подальших досліджень та удосконалень дозволить підвищити точність та надійність класифікації, що, у свою чергу, сприятиме покращенню процесу розслідування цифрових злочинів.

ВИСНОВКИ

У даному дослідженні було створено метод класифікації цифрових доказів на основі їх метаданих методами машинного навчання, що і було поставлено за мету роботи.

Створений метод класифікації покликаний пришвидшити процесу відбору релевантних файлів під час розслідувань у цифровій криміналістиці. Визначення підозрілості файлів базується на метаданих файлів, де при схожій спорідненості метаданих відомих незаконних файлів, визначається підозрілість досліджуваних цифрових доказів, де значення пріоритетності позначається мітками. Розроблена модель не приймає остаточного рішення чи файл є цифровим доказом, але класифікує його на підозрілість. Використання методів керованого машинного навчання дозволяє значно зменшити витрати часу та зусиль для фахівців, які займаються аналізом цифрових доказів.

Такий метод класифікації повинен бути застосований у чутливих до часу випадках цифрових розслідувань, де цифрові докази часто споріднені за неочевидними характеристиками — метаданими файлів, чим він і відрізняється від існуючих методів аналізу цифрових доказів. Прикладом може слугувати випадки злочинів, що пов'язані з дитячою порнографією, де створення фото чи відео знімків здійснено із короткою різницею в часі. Також глибини директорій, в яких розміщуються такі цифрові докази є наступною спорідненою характеристикою, яку потрібно відслідковувати в таких сценаріях.

Серед використаних методів машинного навчання в результаті Випадковий ліс показав найкращі результати за всіма метриками: Точність (англ. accuracy) 0.99, Влучність 1.00 (англ. precision), Повнота 0.99 (англ. recall) та F1-міра (англ. F1-score) 0.99.

Інтеграція розробленого методу класифікації у поширені програмні засоби цифрової криміналістики може значно пришвидшити хід цифрового криміналістичного розслідування, що внесе цінний соціально-економічний вклад у правоохоронні органи, що може значно підвищити ефективність розслідувань та сприяти боротьбі з кіберзлочинністю.

Майбутні дослідження можуть передбачати розробку і оцінку більшої кількості сценаріїв, щоб вони краще відповідали потребам реальних різноманітних криміналістичних розслідувань. Наприклад, моделі для ідентифікації підозрілих файлів з різних джерел, таких як вкладення електронної пошти, хмарні облікові записи тощо, а також для виявлення сторонніх власників файлів. Також важливо зосередитися на показниках запам'ятовування, а не лише на точності, оскільки пропуск будь-якого релевантного доказу є недопустимим. Це потребує налаштування моделей для забезпечення високої чутливості до підозрілих файлів.

ПЕРЕЛІК ДЖЕРЕЛ ПОСИЛАНЬ

- [1] David Lillis. Current Challenges and Future Research Areas for Digital Forensic Investigation [Текст] / David Lillis, Brett Becker, Tadhg O’Sullivan, Mark Scanlon // The 11th ADFSL Conference on Digital Forensics, Security and Law. - 2016 - С. 9-20.
- [2] Використання алгоритмів машинного навчання у комп’ютерному криміналістичному сортуванні для класифікації цифрових файлів на основі їх відповідності розслідуванню / Я. В. Лета, О. М. Барановський // Теоретичні і прикладні проблеми фізики, математики та інформатики: матеріали XXII Всеукраїнської науково-практичної конференції студентів, аспірантів та молодих вчених / Я. В. Лета, О. М. Барановський – Київ, 2024. – (Видавництво «Політехніка»). – С. 151–153.
- [3] Fahad Alanazi. The Value of Metadata in Digital Forensics [Текст] / Fahad Alanazi, Andrew Jones // European Intelligence and Security Informatics Conference, Manchester, UK - 2015 - С. 182-182.
- [4] Dansiger, A. L. Embedded Metadata: Friend or Foe to Our Digital Collections? [Текст] / Dansiger, A. L. // Library Student Journal, 6. - 2011.
- [5] Garfinkel S. Digital forensics XML and the DFXML toolset [Текст] / Garfinkel S. // Digital Investigation, 8(3), - 2012.
- [6] Garrie, D .B. Digital Forensic Evidence in the Courtroom: Understanding Content and Quality [Текст] / Garrie, D .B. // Nw.J.Tech. & Intell. Prop.12i. - 2014.
- [7] Eoghan Casey. Digital evidence and computer crime [Текст] / Eoghan Casey // Academic Press (2nd ed.), San Diego, CA, - 2004.
- [8] Florian Buchholz. On the role of file system metadata in digital forensics [Текст] / Florian Buchholz, Eugene Spafford // Digital Investigation, Volume 1, Issue 4, - 2004, - С. 298-309
- [9] Joakim Kavrestad. Fundamentals of Digital Forensics [Текст] / Joakim Kavrestad // Springer, - 2020.
- [10] Konstantinos Karampidis. A review of image steganalysis techniques for digital forensics [Текст] / Konstantinos Karampidis, Ergina Kavallieratou, Giorgos Papadourakis // Journal of information security and applications, - 2018, - С. 217–235.
- [11] Graeme Horsman. Tool testing and reliability issues in the field of digital forensics [Текст] / Graeme Horsman // Digital Investigation, - 2019, - С. 162-163.

- [12] Athanasios Dimitriadis. Digital forensics framework for reviewing and investigating cyber attack [Текст] / Athanasios Dimitriadis, Nenad Ivezic, Boonserm Kulvatunyou, Ioannis Mavridis // Array, - 2020.
- [13] Stefania Costantini. Digital forensics and investigations meet artificial intelligence [Текст] / Stefania Costantini, Giovanni De Gasperis, Raffaele Olivieri // Annals of Mathematics and Artificial Intelligence, - 2019, - C. 220.
- [14] Eoghan Casey. Handbook of digital forensics and investigation [Текст] / Eoghan Casey // Academic Press, - 2009.
- [15] Owen Defries Brady. Exploiting digital evidence artefacts: finding and joining digital dots [Текст] / Owen Defries Brady, - King's College London, - 2018.
- [16] Karen Kent. Guide to integrating forensic techniques into incident [Текст] / Karen Kent, Suzanne Chevalier, and Tim Grance // Tech. Rep. 800-86, - 2006.
- [17] Flora Amato. A semantic-based methodology for digital forensics analysis [Текст] / Flora Amato, Aniello Castiglione, Giovanni Cozzolino, and Fabio Narducci // Journal of Parallel and Distributed Computing, - 2020.
- [18] Godson Kalipe. Predicting malarial outbreak using machine learning and deep learning approach: a review and analysis [Текст] / Godson Kalipe, Vikas Gautham, and Rajat Kumar Behera // In 2018 International Conference on Information Technology (ICIT) IEEE, - 2018, - C. 33–38.
- [19] IEEE, 2001. Draft Standard for Learning Object Metadata [Текст]. - draft 6.1, - 2001.
- [20] Brian D. Carrier. Defining event reconstruction of digital crime scenes [Текст] / Brian D. Carrier, Eugene H. Spafford // Journal of Forensic Sciences, - 2004.
- [21] M. Pereira, Metadata-Based Detection of Child Sexual Abuse Material [Текст] / M. Pereira, R. Dodhia, H. Anderson, R. Brown // in IEEE Transactions on Dependable and Secure Computing.
- [22] Eoghan Casey. Digital Evidence and Computer Crime Forensic Science [Текст] / Eoghan Casey - Third Edition by Eoghan Casey - C. 25-26
- [23] Humaira Arshad. Digital Forensics: Review of Issues in Scientific Validation of Digital Evidence [Текст] / Humaira Arshad, Aman Bin Jantan, Oludare Isaac Abiodun // Journal of Information Processing Systems Vol.14, №2, - 2018.

- [24] Mohmed Afridhi. Digital Forensics Triage Classification Model using Hybrid Learning Approaches Article [Текст] / Mohmed Afridhi, Palanivel K. // International Journal of Innovative Research in Computer Science & Technology, - 2022.
- [25] Rami M. Mohammad. A Neural Network based Digital Forensics Classification [Текст] / Rami M. Mohammad // 2018 IEEE/ACS 15th International Conference on Computer Systems and Applications, - 2018.
- [26] Fabio Marturana. A Machine Learning-based Triage methodology for automated categorization of digital media [Текст] / Fabio Marturana, Simone Tacconi // Digital Investigation journal.
- [27] R. M. Mohammad. A Neural Network based Digital Forensics Classification [Текст] / R. M. Mohammad // IEEE 15th International Conference on Computer Systems and Applications (AICCSA), Aqaba, Jordan, - 2018.
- [28] Gurpal Singh Chhabra. Hadoop-based analytic framework for cyber forensics [Текст] / Gurpal Singh Chhabra, Varinderpal Singh, Maninder Singh // International Journal of Communication Systems, - 2018.
- [29] Nighat Usman. Intelligent dynamic malware detection using machine learning in ip reputation for forensics data analytics [Текст] / Nighat Usman, Saeeda Usman, Fazlullah Khan, Mian Ahmad Jan, Ahthasham Sajid, Mamoun Alazab, Paul Watters // Future Generation Computer Systems, - 2021.
- [30] Amit V Kachavimath. Distributed denial of service attack detection using naive bayes and k-nearest neighbor for network forensics [Текст] / Amit V Kachavimath, Shubhangeni Vijay Nazare, Sheetal S Akki // In 2020 2nd International conference on innovative mechanisms for industry applications (ICIMIA), IEEE, - 2020.
- [31] What is OpenPose? A Guide for Beginners [Электронный ресурс] - Режим доступа: <https://blog.roboflow.com/what-is-openpose/> - 2024.
- [32] Paola Barra. Gait analysis for gender classification in forensics [Текст] / Paola Barra, Carmen Bisogni, Michele Nappi, David Freire-Obregon, Modesto Castrillon-Santana // In International Conference on Dependability in Sensor, Cloud, and Big Data Systems and Applications, Springer, - 2019.
- [33] Sara Ferreira. Exposing manipulated photos and videos in digital forensics analysis [Текст] / Sara Ferreira, Mario Antunes, Manuel E Correia // Journal of Imaging, - 2021.
- [34] Anton Yudhana. DDoS classification using neural network and naive bayes methods for network forensics [Текст] / Anton Yudhana, Imam Riadi, Faizin Ridho // International Journal of Advanced Computer Science and Applications, 9(11), - 2018.

[35] Satya Sudha. Analysis and evaluation of integrated cyber crime offences [Текст] / Satya Sudha, Ch Rupa // In 2019 Innovations in Power and Advanced Computing Technologies, volume 1, IEEE, - 2019.

[36] Muhammad F. R. Development of digital evidence collector and file classification system with k-means algorithm [Текст] / Muhammad Faris Ruriawan, Bintaran Anggono, Isaac Anugerah Siahaan, Yudha Purwanto // In 2019 IEEE Asia Pacific Conference on Wireless and Mobile, IEEE, - 2019.

[37] Muhammad Ali. A proactive malicious software identification approach for digital forensic examiners [Текст] / Muhammad Ali, Stavros Shiaeles, Nathan Clarke, Dimitrios Kontogeorgi // Journal of Information Security and Applications, - 2019.

[38] Maryam Hina. Email classification and forensics analysis using machine learning [Текст] / Maryam Hina, Mohsan Ali, Abdul Rehman Javed, Gautam Srivastava, Thippa Reddy Gadekallu, Zunera Jalil // 2021 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computing, Scalable Computing & Communications, Internet of People and Smart City Innovation, IEEE, - 2021.

[39] Digital Corpora [Электронный ресурс] - Режим доступа: <https://digitalcorpora.org/> - 2024.

ДОДАТОК А

Код для обробки даних. Файл preprocess.ipynb:

```
import pandas as pd
from datetime import datetime
from sklearn.preprocessing import LabelEncoder
part_illegal = pd.read_csv('collected_metadata.csv')
print("Partly illegal: ", part_illegal.columns)

# Видалення непотрібних стовпців
part_illegal = part_illegal.drop('Known Status', axis=1)
part_illegal = part_illegal.drop('Is Deleted', axis=1)
part_illegal = part_illegal.drop('Full Path', axis=1)

def parse_datetime_with_timezone(timestamp_str):
    # Відокремлюємо дату та час від рядка та видаляємо часовий пояс
    timestamp_str_without_timezone = timestamp_str.split()[0] + " " + timestamp_str.split()[1]
    timestamp = pd.to_datetime(timestamp_str_without_timezone, format="%Y-%m-%d %H:%M:%S")
    unix_timestamp = timestamp.timestamp()
    return unix_timestamp

timestamp_str = "2024-05-17 05:32:22 CEST"
parsed_timestamp = parse_datetime_with_timezone(timestamp_str)
print(parsed_timestamp)
# Застосування функції для парсингу дат
part_illegal['Last Accessed'] = part_illegal['Last Accessed'].apply(parse_datetime_with_timezone)
part_illegal['File Created'] = part_illegal['File Created'].apply(parse_datetime_with_timezone)
part_illegal['Last Modified'] = part_illegal['Last Modified'].apply(parse_datetime_with_timezone)
total_rows = part_illegal.shape[0]
missing_values = part_illegal.isna().sum()
missing_percentage = (missing_values / total_rows) * 100
print(total_rows)
print(missing_percentage)

# 32.3% значень відсутні в цьому стовпчику, тому:
part_illegal = part_illegal.drop('Hash Value', axis=1)
part_illegal = part_illegal.drop('Permissions', axis=1)
part_illegal.dropna(subset=['File Extension'], inplace=True)
```

```
# Кодування категоріальних змінних
label_encoder_extension = LabelEncoder()
label_encoder_type = LabelEncoder()
part_illegal['File Extension'] = label_encoder_extension.fit_transform(part_illegal['File Extension'])
part_illegal['File Type'] = label_encoder_type.fit_transform(part_illegal['File Type'])
print(part_illegal.head())
part_illegal.to_csv('preprocessed_metadata.csv', index=False)
```

ДОДАТОК Б

Код для класифікації з найкращими результатами методом — Випадковий ліс.
Файл `random_forest_classifier.ipynb`:

```
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score, roc_curve, roc_auc_score
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
part_illegal = pd.read_csv('preprocessed_metadata.csv')

X = part_illegal[['File Type', 'File Extension', 'Last Accessed', 'File Created', 'Last Modified', 'Size', 'Address']]
y = part_illegal['Label']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
clf = RandomForestClassifier()
clf.fit(X_train, y_train)
y_pred_rf = clf.predict(X_test)
precision_rf = precision_score(y_test, y_pred_rf)
recall_rf = recall_score(y_test, y_pred_rf)
f1_score_rf = f1_score(y_test, y_pred_rf)
accuracy_rf = accuracy_score(y_test, y_pred_rf)
print("\nRandom Forest Results:")
print(f'Precision: {precision_rf}')
print(f'Recall: {recall_rf}')
print(f'F1-score: {f1_score_rf}')
print(f'Accuracy: {accuracy_rf}')

# Отримання ймовірностей приналежності до незаконних файлів
y_prob = clf.predict_proba(X_test)[:, 1]
# Для ROC-кривої
fpr, tpr, thresholds = roc_curve(y_test, y_prob)
# Розрахунок площі під ROC-кривою (AUC)
auc = roc_auc_score(y_test, y_prob)
# Побудова ROC-кривої
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='blue', label=f'ROC Curve (AUC = {auc:.2f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--', label='Random Guess')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve for Random Forest')
plt.legend()
```

```
plt.show()
X_part_illegal = part_illegal[['File Type', 'File Extension', 'Last Accessed', 'File Created', 'Last Modified', 'Size', 'Address']]
part_illegal_predictions = clf.predict(X_part_illegal)
part_illegal['Label_predicted'] = part_illegal_predictions
correct_predict_rf = sum(part_illegal['Label'] == part_illegal['Label_predicted'])
wrong_predict_rf = sum(part_illegal['Label'] != part_illegal['Label_predicted'])
print(correct_predict_rf, wrong_predict_rf)
```