



ВИЗУАЛИЗАЦИЯ И ПРЕДОБРАБОТКА СОБЫТИЙ ИНФОРМАЦИОННОЙ БЕЗОПАСНОСТИ НА ОСНОВЕ ДАННЫХ CICIDS17

И. Ю. Зеличенко^{1*}

¹Санкт-Петербургский федеральный исследовательский центр Российской академии наук,
Санкт-Петербург, 199178, Российская Федерация

*Адрес для переписки: zelichenok.igor@gmail.com

Аннотация—В настоящее время атаки на компьютерные сети продолжают развиваться со скоростью, опережающей способность специалистов по информационной безопасности создавать новые сигнатуры атак. Эта статья иллюстрирует подход к предобработке сырых данных и визуализации событий информационной безопасности в актуальном наборе. Показано как предобработка и первичное извлечение знаний для дальнейшего использования обработанного набора данных в моделях машинного обучения могут быть использованы при проектировании моделей машинного обучения для систем обнаружения вторжений. Отличительной чертой работы является то, что в качестве исследуемого набора данных, был взят наиболее актуальный набор CICIDS17. Хотя традиционно считаются популярными такие наборы как DARPA2000 и KDD-99, которым уже больше 20 лет. В статье также описываются критерии и характеристики, которые имеет набор.

Ключевые слова—информационная безопасность, большие данные, визуализация, анализ сетевого трафика, системы обнаружения вторжений.

Информация о статье

УДК 004.02.

Язык статьи – русский.

Поступила в редакцию 08.12.2021, принята к печати 20.12.2021.

Ссылка для цитирования: Зеличенко И. Ю. Визуализация и предобработка событий информационной безопасности на основе данных CICIDS17 // Информационные технологии и телекоммуникации. 2021. Том 9. № 4. С. 49–55. DOI 10.31854/2307-1303-2021-9-4-49-55.



VISUALIZATION AND PROCESSING OF INFORMATION SECURITY EVENTS BASED ON CICIDS DATA 17

I. Zelichenok^{1*}

¹St. Petersburg Federal Research Center of the Russian Academy of Sciences,
St. Petersburg, 199178, Russian Federation

*Corresponding author: zelichenok.igor@gmail.com

Abstract—At present, attacks on computer networks continue to develop at a speed that outstrips the ability of information security specialists to create new attack signatures. This article illustrates an approach to preprocessing raw data and visualizing information security events in a live dataset. It is shown how preprocessing and primary knowledge extraction for further use of the processed dataset in machine learning models can be used in the design of machine learning models for intrusion detection systems. A distinctive feature of the work is that the most relevant set CICIDS17 was taken as the studied dataset. Although traditionally considered popular such kits as DARPA2000 and KDD-99, which are more than 20 years old. The article also describes the criteria and characteristics that the set has.

Keywords—information security, big data, visualization, network traffic analysis, intrusion detection systems.

Article info

Article in Russian.

Received 08.12.2021, accepted 20.12.2021.

For citation: Zelichenok I.: Visualization and Processing of Information Security Events Based on CICIDS Data 17 // Telecom IT. 2021. Vol. 9. Iss. 4. pp. 49–55 (in Russian). DOI 10.31854/2307-1303-2021-9-4-49-55.



Введение

Системы обнаружения вторжений (IDS) и системы предотвращения вторжений (IPS) являются наиболее важными инструментами защиты от сложных и постоянно растущих сетевых атак. Из-за отсутствия надежных наборов данных для тестирования и валидации подходы к обнаружению вторжений на основе аномалий страдают от роста ресурсоёмкости.

Также, многие проблемы, связанные с внедрением элементов искусственного интеллекта в существующие системы информационной безопасности, требуют комплексного подхода к их решению. Например, в работе [1] представлена математическая основа для адаптивного моделирования и симуляции противостояния на основе совместно развивающихся агентов для кибервойны. Такая разработка потенциально может включать в себя методы машинного обучения для автоматической симуляции действий участников. Работа [2] посвящена разработке архитектуры анализа защищенности компьютерных сетей на основе имитаций действий злоумышленников. В статье представлены продемонстрированы обобщенные модели атак, модели анализируемой системы и оценки уровня защищенности, которые были полезны при анализе текущего набора данных. В работе [3] были представлены рекомендации по проектированию безопасных сетей и IDS.

В дополнение к существующим исследованиям, в данной работе представлен пример предобработки и визуализации для нечасто используемых в научной среде данных CICIDS17¹, подготовленных Канадским институтом кибербезопасности [4]. В русскоязычной научной среде для исследований зачастую используются данные DARPA2000² и их производный набор KDD-99³, которые, в силу возраста, во многом теряют свою актуальность. Данная работа призвана частично закрыть пробел в области реализации разработок моделей машинного обучения для обнаружения сетевых атак при помощи использования актуальных наборов данных.

Статья состоит из четырех разделов. Раздел «Обзор датасета» посвящен разбору набора данных и выводу основных правил, которыми руководствовались сотрудники института при создании используемого набора данных. В разделе «Статистический анализ» проведен количественный и качественный разбор датасета, визуально продемонстрированы типы атак и распределение атак по кластерам. В исследовании использовались популярные библиотеки машинного обучения, код был написан на языке Python версии 3.8. При разработке моделей использовался инструментарий, включающий в себя среду разработки (IDE), Pandas и NumPy для работы с данными, библиотеку для машинного обучения Ski-kit learn, и вспомогательных инструментов для обработки данных.

¹ CICIDS17 dataset. Available at: <https://www.unb.ca/cic/datasets/ids-2017.html>

² DARPA2000 dataset. Available at: <https://www.ll.mit.edu/r-d/datasets/2000-darpa-intrusion-detection-scenario-specific-datasets>

³ KDD-CUP 99 dataset. Available at: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>



Обзор набора данных CICIDS17

Разработчиком набора данных является Канадский институт кибербезопасности, всего опубликовано несколько наборов (2017 и 2019 годов). Набор данных CICIDS2017 собран специалистами института в результате мониторинга сетевого трафика для прототипа в изолированной среде. В сети моделировались действия для 25 легальных пользователей, а также вредоносные действия нарушителей. Набор объединяет более 50 Гб «сырых» данных в формате PCAP и включает 8 предобработанных файлов в формате CSV, содержащих размеченные сессии с выделенными признаками в разные дни наблюдения. Набор имеет следующие характеристики:

- Полная конфигурация сети: включает шлюзы, брандмауэры, коммутаторы, маршрутизаторы и наличие множества операционных систем, таких как Windows, Ubuntu и Mac OS X.
- Полный трафик: наличие агента профилирования пользователей, 12 разных машин в сети – объекте атаки и реальные атаки из сети – источника атаки.
- Помеченный набор данных: показаны ярлыки безобидных и атакующих для каждого дня. Кроме того, подробности о времени атаки опубликованы в документе набора данных.
- Полное взаимодействие: показаны атаки как внутри, так и между внутренними LAN. Имелись две разные сети и связь через Интернет.
- Полный захват: использовался зеркальный порт и весь трафик был захвачен и записан на сервере хранения.
- Доступные протоколы: наличие всех распространенных доступных протоколов, таких как HTTP, HTTPS, FTP, SSH и протоколы электронной почты.
- Разнообразие атак: в наборе представлены разнообразные атаки.
- Неоднородность: захват сетевого трафика с главного коммутатора, дампа памяти и системных вызовов со всех машин-жертв во время выполнения атак.
- Набор функций: извлечено более 80 функций сетевого потока из сгенерированного сетевого трафика с помощью CICFlowMeter и предоставлен набор данных сетевого потока в виде файла CSV.
- Метаданные: полностью объяснен набор данных, который включает время, атаки, потоки и метки.

Таким образом, можно сделать вывод о пригодности набора данных для проведения экспериментальных исследований. Однако для данного набора присущ дисбаланс классов.

Статистический анализ

Процесс проектирования включал в себя предобработку данных для повышения производительности и точности получаемой модели. Было проанализировано распределение атак по записям. Подавляющее большинство атак в наборе данных составляют DoS атаки (рис. 1, см. ниже).

Затем шел процесс проверки набора на пустые и неструктурированные значения, после чего набор данных был сохранен в отдельный файл формата CSV.

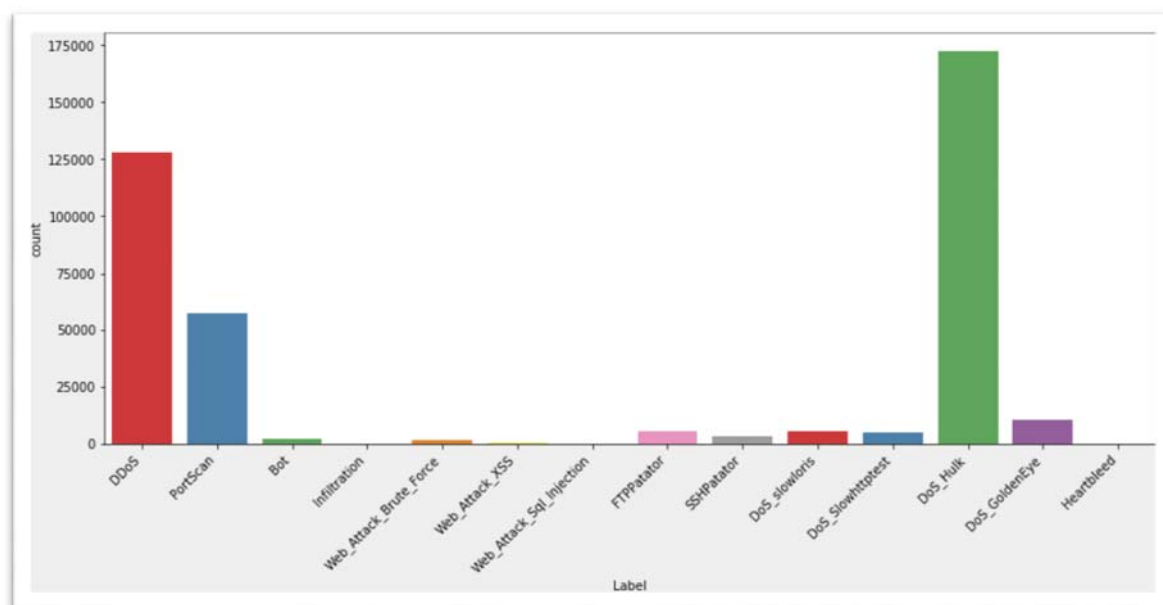


Рис. 1. Распределение атак по типам в наборе данных

Набор данных содержит 78 функций и разделен на 15 категорий (14 атак и 1 «нормальное» состояние). Следовательно, следующим шагом, предпринятым в рамках экспериментального исследования, была визуализация набора данных в пространстве признаков [5]. Для этого был использован анализ главных компонентов (PCA) для уменьшения размерности, а затем сокращенный набор данных был передан в t-SNE для визуального представления в 2D-пространстве.

Благодаря встроенному в библиотеку машинного обучения Ski-kit learn классу `sklearn.ensemble.RandomForestClassifier` [6], выявляющему зависимость разных признаков друг от друга, был проведен анализ влияния признаков на конечный результат. После изъятия из датасета свойств с минимальной важностью, признаковое пространство сократилось до следующих свойств (табл.).

Таблица.

Анализируемые поля в наборе данных

Название поля	Краткое описание
Average Packet Size	Средний размер передаваемого сообщения.
Flow Bytes/s	Скорость передачи данных.
Max Packet Length	Максимальная длина пакета.
Fwd Packet Length Mean	Средняя длина исходящего пакета.
Fwd IAT Min	Минимальная задержка между сообщениями.
Total Length of Fwd Packets	Суммарная длина исходящих сообщений.
Fwd IAT Std	Среднеквадратичное отклонение от суммарной длины пакетов.
Flow IAT Mean	Среднее значение задержки между сообщениями.
Fwd Packet Length Max	Максимальная длина исходящего сообщения.
Fwd Header Length	Суммарная длина заголовков сообщений.



Были выбраны 10000 случайных строк из набора данных для визуализации. Задан случайный Random Seed для большей достоверности [7]. Далее был выполнен анализ главных компонент. С девятнадцатью компонентами коэффициент дисперсии остался равен 99 %.

На рис. 2 показано распределение данных в 2D-пространстве. Очевидно, что атаки в пространстве плохо отделены от нормального состояния. Кластеры атак трудно увидеть, вместо этого они находятся в том же месте, что и точки данных «нормального состояния».

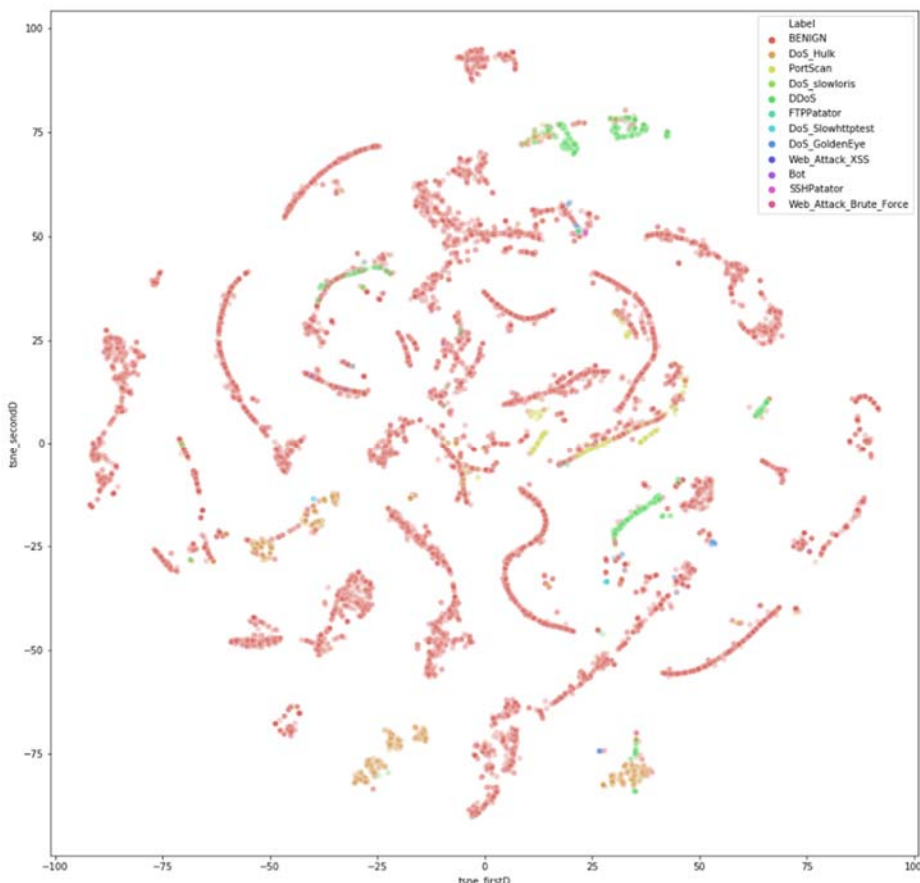


Рис. 2. Распределение данных в двумерном пространстве

Выводы

В работе представлен метод предобработки и первичного извлечения знаний из наборов данных для дальнейшего использования, обработанного датасета в моделях машинного обучения. Новизна работы заключается в выборе исследуемого набора данных, который можно назвать более актуальным, нежели другие часто используемые наборы, такие как DARPA2000 и KDD-99. Согласно ранее описанным критериям, CICIDS17 имеет все то, что требуется для разработки надежной системы обнаружения вторжений. Алгоритмы анализа и предобработки наборов данных могут уточняться, дополняться и дорабатываться.

Работа выполнена при частичной финансовой поддержке бюджетной темы 0073-2019-0002.



Литература

1. Kotenko I. Agent-Based Modeling and Simulation of Cyber-Warfare between Malefactors and Security Agents in Internet // 19th European Simulation Multiconference "Simulation in wider Europe". ECMS 2005. Riga, Latvia, 1–4 June. 2005. С. 533–543.
2. Котенко И. В., Степашкин М. В., Богданов В. С. Архитектуры и модели компонентов активного анализа защищенности на основе имитации действий злоумышленников // Проблемы информационной безопасности. Компьютерные системы. 2006. № 2. С. 7–24.
3. Котенко И. В., Десницкий В. А., Чечулин А. А. Исследование технологии проектирования безопасных встроенных систем в проекте Европейского сообщества SecFutur // Защита информации. Инсайд. 2011. № 3. С. 68–75.
4. Iman Sharafaldin, Arash Habibi Lashkari, and Ali A. Ghorbani Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization // 4th International Conference on Information Systems Security and Privacy (ICISSP). Portugal. 2018.
5. Виткова Л. А., Зеличенко И. Ю. Методика мониторинга и диагностики локальных инцидентов с потенциалом протестной мобилизации // Информатизация и связь. 2021. № 5. С. 90–96.
6. Горюнов М. Н., Мацкевич А. Г., Рыболовлев Д. А. Синтез модели машинного обучения для обнаружения компьютерных атак на основе набора данных CICIDS2017 // Труды института системного программирования РАН. 2020. Т. 32 (5). С. 81–94. DOI: 10.15514/ISPRAS-2020-32(5)-6.
7. Зеличенко И. Ю., Пирмагомедов Р. Я. Применение методов машинного обучения для анализа физической активности пользователя смартфона // Информационные технологии и телекоммуникации. 2020. Т. 8. № 2. С. 92–108. DOI: 10.31854/2307-1303-2020-8-2-92-108.

References

1. Kotenko I. Agent-Based Modeling and Simulation of Cyber-Warfare between Malefactors and Security Agents in Internet // 19th European Simulation Multiconference "Simulation in wider Europe". ECMS 2005. Riga, Latvia, 1–4 June. 2005. С. 533–543.
2. Kotenko I.V., Stepashkin M.V., Bogdanov V.S. Architectures and Models of Active Vulnerabilities Analysis Based on Simulation of Malefactors' Actions // Information Security Problems. Computer Systems. 2006. No. 2. pp. 7–24.
3. Kotenko I. V., Desnickij V. A., Chechulin A. A. Issledovanie tekhnologii proektirovaniya bezopasnyh vstroennyh sistem v proekte Evropejskogo soobshchestva SecFutur // Zashchita informacii. Insajd. 2011. № 3. S. 68–75.
4. Iman Sharafaldin, Arash Habibi Lashkari, and Ali A. Ghorbani Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization // 4th International Conference on Information Systems Security and Privacy (ICISSP). Portugal. 2018.
5. Vitkova L. A., Zelichenok I. U. Ethod for Monitoring and Diagnosing Local Incidents with the Potential for Protest Mobiliza // Informatization and communication. 2021. No. 5. pp. 90–96.
6. Goryunov M. N., Matskevich A. G., Rybolovlev D. A. Synthesis of a Machine Learning Model for Detecting Computer Attacks Based on the CICIDS2017 Dataset. Proceedings of the Institute for System Programming of the RAS (Proceedings of ISP RAS). 2020;32(5):81–94. (In Russ.) [https://doi.org/10.15514/ISPRAS-2020-32\(5\)-6](https://doi.org/10.15514/ISPRAS-2020-32(5)-6).
7. Zelichenok I., Pirmagomedov R.: Tutorial on Using Machine Learning for Activity Recognition Via a Smartphone // Telecom IT. 2020. Vol. 8. Iss. 2. pp. 92–108 (in Russian). DOI 10.31854/2307-1303-2020-8-2-92-108.

Зеличенко Игорь Юрьевич

аспирант Санкт-Петербургского Федерального исследовательского центра Российской академии наук, zelichenok.igor@gmail.com

Zelichenok Igor Yu.

postgraduate student, Saint-Petersburg Federal Research Center, Russian Academy of Sciences, zelichenok.igor@gmail.com