

Моделирование системы информационной безопасности на основе анализа системных журналов

Ю.В. Ефимова, А.Г. Гаврилов

Чистопольский филиал «Восток» Казанского национального исследовательского технического университета им. А.Н. Туполева – КАИ, Чистополь

Аннотация: В статье дано описание метода идентификации аномальных действий пользователей корпоративных компьютерных систем на основе анализа лог-файлов. Предложенный метод основывается на кластеризации событий системного журнала алгоритмом IPLoM и построении матрицы счета событий для ее дальнейшего анализа с использованием методов машинного обучения.

Ключевые слова: машинное обучение, информационная безопасность, анализ данных, мониторинг системного журнала, лог-файл, IPLoM.

В настоящее время в крупномасштабных корпоративных системах актуальным является решение задач обеспечения информационной безопасности. Согласно исследованиям существуют пять поколений кибератак и киберзащит. Кибератаки первого поколения сводятся к прямому нацеливанию на конкретные хосты, злоумышленники эксплуатируют незакрытые уязвимости и небезопасные протоколы. Кибератаки второго поколения основаны на использовании автоматических утилит для сканирования уязвимостей и уязвимых протоколов одновременно на нескольких сотнях и тысячах хостах корпорации. Методы защиты от кибератак первого и второго поколения основаны на создании сильного сетевого периметра. Третье поколение кибератак акцентируется на проникновении внутрь корпоративной сети и захват контроля над системами кибербезопасности и привилегированными аккаунтами. Соответственно, информационная защита проектируется по принципу «изнутри наружу», т.е. в первую очередь защищаются наиболее критичные секторы корпоративной компьютерной структуры. Кроме того, внутри предприятия также выстраивается дополнительный периметр безопасности, а все аккаунты защищаются многофакторной аутентификацией. При этом следует учитывать, что каждое поколение атак изначально проводилось только ограниченной группой продвинутых хакеров, а по истечению времени становились доступны всем. Так, кибератаки четвертого поколения нацелены, во-первых, на то, чтобы заставить защитные системы от атак третьего поколения генерировать большое число предупреждающих сообщений, с целью отвлечения внимания системных администраторов. Во-вторых, на то, чтобы пока системные администраторы решают поступившие проблемы, запустить

специальное вредоносное ПО, предназначенное для автоматического поиска и использования слабых мест в киберзащите корпорации третьего поколения. Кибератаки пятого поколения подразумевают не только автоматическое использование вредоносного ПО, но и постоянную смену методов атаки, которые при этом способны самостоятельно адаптироваться к противодействующим им средствам защиты. Стоит учитывать, что кибератаки пятого поколения проводятся только правительственными хакерами, так как их реализация требует больших затрат: финансовых, ресурсных и временных [1,2]. Таким образом, актуальным является создание автоматизированных систем для идентификации угроз кибератак третьего и четвертого поколения, в том числе и от внутрикорпоративных злоумышленников.

При этом разграничение поведения легальных пользователей крупных корпоративных систем от аномальных действий злоумышленников, приводящих к реализации угроз информационной безопасности различного рода на основе имеющихся уязвимостей становится все более актуальной задачей [3]. Это связано с тем, что все большее количество стран переходит на электронный документооборот и доступ к коммерческой информации может повлечь за собой финансовые, трудовые и временные потери[4].

Одним из подходов к диагностике критических состояний системы с точки зрения кибербезопасности является постоянный анализ системных журналов в режиме реального времени, поскольку информация, имеющаяся в данных журналах, отражает состояние системы, ее ресурсы, а также действия пользователей.

При этом традиционная ручная обработка событий системных журналов становится трудно осуществимой по следующим причинам: ежедневно информационные системы производят достаточно большой объем данных[5], что делает невозможным анализ этих данных с помощью простых подходов к поиску критических выражений, например, с использованием поиска «по словарю»; велика вероятность ложно-повторных срабатываний фильтрующих модулей из-за избыточности системных журналов в крупных корпоративных сетях, в результате многократного дублирования действий пользователей при реализации отказоустойчивых механизмов. Кроме того, при разработке систем анализа аномального поведения пользователей для крупных корпораций, имеющих территориально-распределенную архитектуру внутренней сети и обладающих параллельным характером электронного документооборота, а также использующее облачные сервисы хранения, разработчикам достаточно сложно в полной мере оценить

поведение системы для выявления критических ситуаций из имеющихся данных системных журналов.

Целью разработки является создание программного комплекса мониторинга системных журналов, производящего интеллектуальный анализ данных, основанный на технологиях математической статистики и машинного обучения для распознавания аномальных состояний системы, что особенно актуально в службах корпоративной информационной безопасности и технической поддержки для диагностики действий пользователей.

В зависимости от области применения программного комплекса будут различаться ключевые признаки аномальных состояний, что требует достаточно длительного обучения, предшествующего его использованию на практике. Алгоритм обнаружения подозрительных действий пользователей включают в себя несколько базовых блоков: получение лог-файлов; анализ и фильтрация информации; выделение ключевых признаков и разбиение событий системного журнала на два кластера – события, не связанные с реальными сбоями системы и критические аномалии.

Получение лог-файлов. Для качественного анализа записей системных журналов при формировании обучающей выборки необходимо учитывать, что исходный лог-файл должен содержать информацию о действиях пользователей не только служебного типа в рамках обычных бизнес-процессов, но и действия пользователей, оцениваемые как неприемлемые с точки зрения информационной безопасности и нарушающие, например, существующую политику разграничения доступа. Системный журнал содержит записи о действиях, производимых пользователями и метки времени, что позволяет использовать их для критической оценки и проведения интеллектуального анализа с целью выявления неправомерных действий. Пример фрагмента лог-файла показан на рис.1.

```
94 Aug 5 2019 09:24:34 ZTE %%01SFM/4/SFUINCHANNELOPEN(I)
[94]: Command SERDES interface input chip 0 channel5 is opened!
95 Aug 5 2019 09:24:34 ZTE %%01SFM/4/BOOTMODE(I) [95]: Slot 14
has startup with Normal mode.
96 Aug 5 2019 09:24:32 ZTE %%01SFM/4/SFUREG(I) [96]: Command
registered successfully.
97 Aug 5 2019 09:23:37 ZTE %%01SFM/3/PWRONFINISH(I) [97]:
SlotD12, board power-on finish!
```

Рисунок 1 – Пример лог-файла

Анализ и фильтрация информации. Вторым этапом обнаружения аномалий является анализ и фильтрация системного журнала. Записи о действиях пользователя не имеют фиксированной длины и могут значительно различаться по своей структуре. При обработке записей системных журналов необходимо выделить структурные части каждой строки. Каждая запись содержит в себе информацию, которая достаточно часто значительно отличается структурно. Для выделения структурных частей в автоматическом режиме возможно использовать готовые программные инструменты, например, Splunk, Logstash и др. Данные программные продукты при этом обладают достаточно сложными процедурными решениями для настройки и использования правил кластеризации. Наиболее гибким решением для анализа системных журналов будет использование методов кластеризации итерационного типа, когда исходные данные последовательно разделяются в отдельные кластеры в зависимости от их формата. Среди существующих методов кластеризации больших объемов данных высокой точностью характеризуется метод Iterative Partitioning Log Mining (IPLoM) [6,7].

В основе метода лежит последовательная итерационная обработка каждой записи системного журнала и выделения фреймов. Строка записи рассматривается как набор данных, обладающих различными атрибутами. В качестве атрибутов в случае обработки системных журналов выступают тестовые словосочетания и отдельные слова, составляющие эту запись. Каждый фрейм характеризуется постоянными величинами и переменными. К переменным величинам относятся шаблонные словосочетания временных характеристик действий пользователей, например, метки времени, идентификаторы устройств и пользователей, получающих доступ к ресурсам компьютерной системы. К постоянным величинам можно отнести некоторые стандартные выражения, используемые для описания событий, метки действий пользователей, например, получен доступ, завершена загрузка и т.д.

В результате обработки системного журнала переменные величины отделяются от постоянных, а постоянные величины группируются до получения фреймов. При построении множеств постоянных атрибутов сложность автоматизации заключается в том, что порядок следования отдельных атрибутов достаточно часто имеет критическое значение.

При выполнении нескольких итераций алгоритма кластеризации формируются фреймы, в которых переменные величины заменены знаком «*», а постоянные величины оставлены в исходном текстовом формате.

Пример фрейма представлен на рисунке 2.

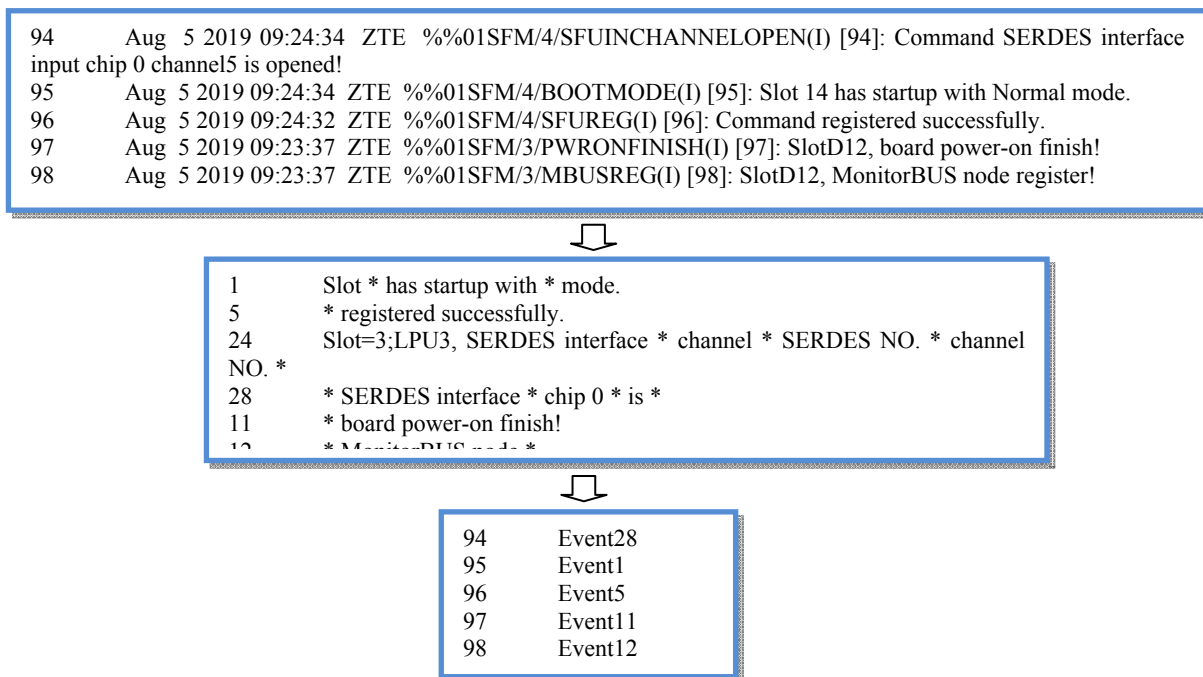


Рисунок 2 – Получение фреймов лог-файла

Анализируя набор данных, полученный после выполнения нескольких итераций алгоритма кластеризации видно, что каждый сформированный фрейм содержит информацию о некотором конкретном событии, отраженном в системном журнале, и имеет одинаковый размер. В результате анализа, реализовано разделение записей системного журнала по размеру события. После чего, различные кластеры сформированы в соответствии с длиной записи. Можно сделать предположение, что дальнейшая группировка по размеру полученных фреймов получит множество всех записей, относящихся к одному и тому же типу событий. При этом стоит учитывать, что при анализе записей реальных журналах возможно ошибочное присвоение записи заданного номера фрейма, основываясь только на длину этой записи.

На втором этапе разбиения на фреймы M проанализированных записей системного журнала можно представить в виде матрицы размерностью $K \times M$, где K – максимальная длина записей в данном кластере, n – количество записей длиной k , M – количество

записей в исходном системном журнале. Таким образом, общее количество проанализированных записей в системном журнале будет равно

$$m = \sum_{k=1}^K n_k,$$

где n_k – количество записей в n кластере, k – длина записи, $k = \overline{1, K}$.

Далее матрица упорядочивается по возрастанию количества постоянных величин в строке, показано на рис. 3. Тогда элементы нескольких первых строк будут являться уникальными точками для построения кластеров событий, т.е. строки содержат только константы, определяющие статус события либо как стандартное действие пользователя в рамках трудовых обязанностей, либо как подозрительное. Все остальные записи матрицы разделяются на кластеры по признаку наличия в определенной позиции фрейма уникальной точки разбиения, т.е. каждый сгенерированный кластер имеет одну и то же уникальную точку в определенном столбце.

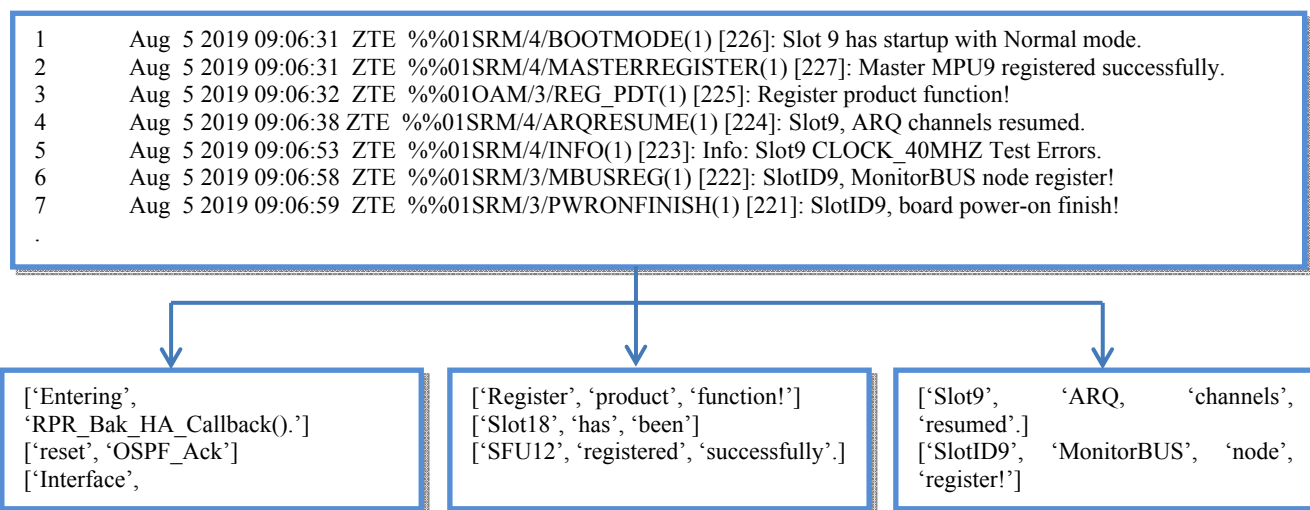


Рисунок 3 – Разбиение фрейма для получения уникальных точек

Для определения наличия скрытых закономерностей между уникальными точками в столбцах выбираются первые столбцы для кластеров с количеством уникальных точек более одной и разделяются на основе отношений между ними, показано на рис. 4.

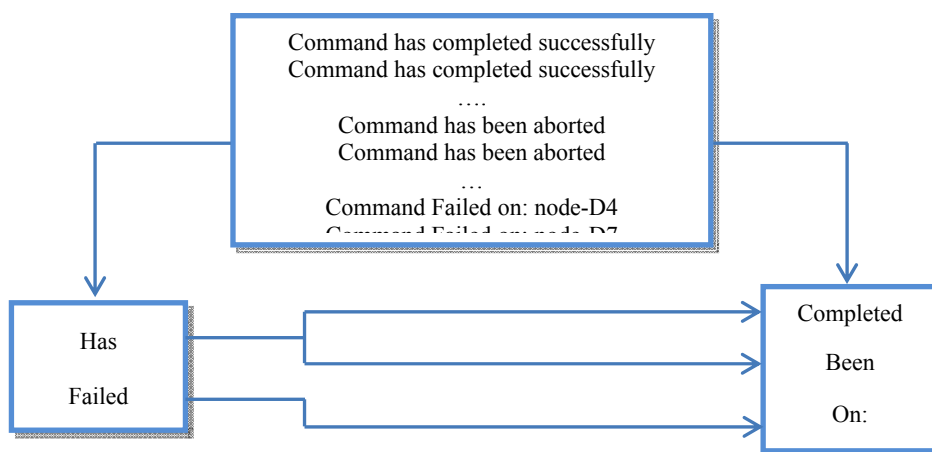


Рисунок 4 – Этап поиска отношений между уникальными точками

Отношения отображения могут отличаться в зависимости от того, какие именно уникальные точки используются, возможно, получение несколько типов отношений 1:1, 1:M, M:1, M:M. Например, в случае отношений 1:M или M:1 позиция M отношение «ко многим» может представлять собой либо переменные значения (предполагается, что имеется один тип событий) или постоянные значения (предполагается, что каждое значение фактически представляет собой другой тип события). Для отношений 1:M и M:1 необходимо сначала решить, содержит ли столбец на стороне M константы или переменные. Установление нижней и верхней границы для разделения отношений 1:M или M:1 реализуется на этапе проектирования для каждой системы рассматривается в частном порядке. В случае отношений M:M шаг разделения повторяется до получения отношения 1:M или пропускается.

Этап анализа и фильтрации информации, а также разделения всего системного журнала на отдельные кластеры по уникальным точкам завершается проведением экспертной оценки качества полученных кластеров.

Экспертная оценка завершает обработку всех кластеров, созданных на предыдущих этапах на основе системного журнала, по результатам которой получены отдельные фреймы для любого ключевого события. Для каждого столбца в кластере подсчитывается количество уникальных точек. Если в столбце есть только одна уникальная точка, то она считается постоянной. В противном случае уникальные точки в столбце являются переменными и будут заменены символами «*» в итоговом журнале фреймов.

Выделение ключевых признаков. По завершению этапа фильтрации лог-файлов имеем фреймы, содержащие необработанные данные системных журналов. Необходимо

сформировать векторное представление признаков аномальных событий, используемое в дальнейшем для обнаружения искомых критических ситуаций. Алгоритм получения вектора признаков следующий: фрейм разделяется на семантические структуры, далее применяются методы машинного обучения и интеллектуального анализа, позволяющие разделить на кластеры все события системного журнала с указанием частоты появления того или иного события.

Для определения наличия критических событий исходный системный журнал соотносит с каждым событием, полученным на этапе кластеризации. Все записи журнала разделяют по идентификатору пользователя, реализующего некоторые действия, так называемое разделение с использованием «окна сеанса». Далее строится матрица событий для каждого уникального идентификатора с указанием тех фреймов событий, которые сгенерированы пользователем. Каждый фрейм, полученный на этапе анализа и фильтрации, нумеруется для более простого представления в матрице событий. Например, в результате выполнения служебных трудовых обязанностей некоторым пользователем в системный журнал был занесен ряд событий. Тогда формируется последовательность в которой подсчитывается количество появлений каждого события в течении «окна сеанса» пользователя с данным идентификатором. Например, если данная последовательность содержит следующие значения [1,0,4,0,0,0,0], это означает, что в этой последовательности первое событие произошло один раз, третье четыре раза, а остальные события не были зафиксированы.

Аналогичным образом каждое событие системного журнала преобразуется в вектор счета событий. Наконец, множество векторов составляют матрицу событий X , где X_{ij} отражает сколько раз событие j произошло в i -том векторе счета событий, последним столбцом указывается признак наличия аномалии в данном векторе. В случае отсутствия аномалии записывается «0», а при отнесении события к аномальному – «1».

Полученная в результате предшествующих преобразований матрица делает возможным на основе вектора семантических признаков классифицировать любое событие, входящие в лог-файл, по степени аномальности с точки зрения информационной безопасности. Для этого используются методы машинного обучения, например, инструмент Scikit-learn, представляющий собой пакет Python для интеллектуального анализа данных, в том числе классификацию с использованием алгоритмов машинного обучения[8,9].

Результатом применения разработанного программного комплекса мониторинга системных журналов для распознавания подозрительных действий пользователей компьютерных систем является верное выявление критических ситуаций в 90% рассмотренных случаев, что подтверждает его достаточную эффективность.

Сферой применения разработанной системы можно считать крупные корпоративные холдинги, обладающие обширным серверным системным журналом благодаря большому объему сетевого трафика; системы связи между различными отделами на крупных промышленных предприятиях, желающих защитить информацию о новейших технических разработках с целью сохранения коммерческой тайны; системы дистанционного и традиционного обучения для предотвращения утечки диагностирующих материалов и защиты от подмены одного обучающегося другим при прохождении итоговой аттестации; для повышения эффективности управления предприятием [10]. Кроме того, реализованный алгоритм обнаружения пригоден для классификации аномальных состояний системы в режиме реального времени в рамках модулей помощи принятия решения для использования системными администраторами и в службах технической поддержки.

Литература

1. Караев А. Эволюция кибератак и киберзащиты. Часть 1: первые три поколения // Системный администратор, 2019. №5. С. 32-37.
2. Караев А. Эволюция кибератак и киберзащиты. Часть 2: последние два поколения // Системный администратор, 2019. №6. С. 16-21.
3. Грошева Е. К., Невмержицкий П. И. Информационная безопасность: современные реалии // Бизнес-образование в экономике знаний. 2017. №3. С. 35-38.
4. Аникин И. В. Управление внутренними рисками информационной безопасности корпоративных информационных сетей // Научно-технические ведомости СПбГПУ. Информатика. Телекоммуникации. Управление. 2009. №3. С. 35-40.
5. Закирова Г.Ф., Ефимова Ю.В. Автоматизированная система учета клиентов в компании связи // В сборнике: Теория и практика системной динамики Материалы конференции VIII Всероссийской конференции (с международным участием). Ответственный редактор: Олейник А.Г. 2019. С. 84-88.

6. A. Makanju, A. Zincir-Heywood, and E. Milios. Clustering event logs using iterative partitioning. In KDD'09: Proc. Of International Conference on Knowledge Discovery and Data Mining, 2009. DOI: 10.1145/1557019.1557154 pp. 1255-1263.

7. R. Vaarandi. Mining event logs with slct andloghound. In Proceedings of the 2008 IEEE/IFIP Network Operations and Management Symposium. 2008. pp 1071–1074.

8. Кавчук Д.А., Тумоян Е. П., Евстафьев Г. А. Интеллектуальный подход к анализу рисков и уязвимостей информационных систем // Известия ЮФУ. Технические науки. 2013. №12. С. 79-86.

9. Братченко А. И., Бутусов И. В., Кобелян А. М., Романов А. А. Применение методов теории нечетких множеств к оценке рисков нарушения критически важных свойств защищаемых ресурсов автоматизированных систем управления // Вопросы кибербезопасности. 2019. №1. С. 18-24.

10. Свирина А.А. Повышение эффективности управления при изменении принципов управления предприятием // Российское предпринимательство. 2007. №6. С. 63-66.

References

1. Karaev A. Evolyuciya kiberatak i kiberzashchity. CHast' 1: pervye tri pokoleniya. Sistemnyj administrator, 2019, №5. pp. 32-37.

2. Karaev A. Evolyuciya kiberatak i kiberzashchity. CHast' 2: poslednie dva pokoleniya. Sistemnyj administrator, 2019, №6. pp. 16-21.

3. Grosheva E. K., Nevmerzchickij P. I. Biznes-obrazovanie v ekonomike znaniy. 2017. №3. pp. 35-38.

4. Anikin I. V. Nauchno-tekhnicheskie vedomosti SPbGPU. Informatika. Telekommunikacii. Upravlenie. 2009. №3. pp. 35-40.

5. Zakirova G.F., Efimova YU.V. Avtomatizirovannaya sistema ucheta klientov v kompanii svyazi. V sbornike: Teoriya i praktika sistemnoj dinamiki Materialy konferencii VIII Vserossijskoj konferencii (s mezhdunarodnym uchastiem). Otvetstvennyj redaktor: A.G. Olejnik. 2019. pp. 84-88.

6. A. Makanju, A. Zincir-Heywood, and E. Milios. Clustering event logs using iterative partitioning. In KDD'09: Proc. Of International Conference on Knowledge Discovery and Data Mining, 2009. DOI: 10.1145/1557019.1557154 pp. 1255-1263.

7. R. Vaarandi. Mining event logs with slct andloghound. In Proceedings of the 2008 IEEE/IFIP Network Operations and Management Symposium. 2008. pp 1071–1074.



8. Kavchuk D.A., Tumoyan E. P., Evstaf'ev G. A. Izvestiya YUFU. Tekhnicheskie nauki. 2013. №12. pp. 79-86.
9. Bratchenko A. I., Butusov I. V., Kobelyan A. M., Romanov A. A. Voprosy kiberbezopasnosti. 2019. №1. pp. 18-24.
10. Svirina A.A. Rossijskoe predprinimatel'stvo. 2007. №6. pp. 63-66.