

Д. В. КОМАШИНСКИЙ, И. В. КОТЕНКО

## МЕТОД ИЗВЛЕЧЕНИЯ СТРУКТУРНЫХ ПРИЗНАКОВ ВРЕДОНОСНОГО ПРОГРАММНОГО ОБЕСПЕЧЕНИЯ В ЗАДАЧЕ ЕГО ОБНАРУЖЕНИЯ

Представлен метод извлечения из набора данных статических структурных признаков, развивающий ранее предложенный авторами подход к обнаружению вредоносных документов на основе применения методов машинного обучения. Метод базируется на использовании структурных особенностей элементов вредоносных документов и связей между ними.

*Ключевые слова:* извлечение признаков, вредоносное программное обеспечение, классификация.

**Введение.** Одна из важнейших задач компьютерной безопасности — обнаружение вредоносных программ. За последнее десятилетие темп появления нового, ранее неизвестного вредоносного программного обеспечения (ВПО) существенно изменился. Так, в 2001 г. появлялось не более 8—10 новых экземпляров подобных программ в сутки [1], а осенью 2011 г. этот показатель составил порядка 40 000 экземпляров [2]. Необходимость анализа постоянно увеличивающегося потока данных привела к разработке дополнительных, эвристических подходов к обнаружению ВПО. Эта задача решается, в том числе, с помощью методов машинного обучения (Machine Learning — ML) и интеллектуального анализа данных (Data Mining — DM) [1, 3]. В настоящее время существуют промышленные программные системы, используемые поставщиками услуг по информационной безопасности для предварительной классификации потока входных данных и разделения их на категории [3].

Применение методов ML и DM имеет характерные особенности. Они определяются рядом вопросов, на которые исследователь должен ответить при планировании экспериментов. К таковым, в первую очередь, относятся следующие проблемы:

- определение исходного набора данных, используемого для подготовки обучающей и тестовой выборок;
- выбор методов обучения, формирующих конечную модель обнаружения;
- выбор процедур извлечения из набора данных признаков, определяющих пространство используемых атрибутов.

Вопрос выбора исходного набора данных решается за счет применения как открытых наборов, используемых в академических исследованиях [4], так и наборов, формируемых компаниями, которые занимаются проблемами информационной безопасности [3]. Вопрос применения того или иного метода обучения свойственен не только задаче обнаружения ВПО. Существует ряд работ (см., например, [5]), содержащих рекомендации по его решению, определяемые, как правило, типом решаемой задачи, характером используемых в процессе обучения данных и требованиями к ресурсопотреблению процедуры обучения модели. Таким образом, во многом успех задачи обнаружения ВПО определяется процедурой извлечения признаков (Feature Extraction — FS).

Одним из открытых вопросов в области обнаружения ВПО является своевременное выявление Web-ресурсов, на которых функционируют так называемые „пакеты эксплуатации уязвимостей“ (Exploit Kits). Подобные вредоносные программы при обращении пользовательского приложения по зараженной ссылке проверяют наличие на стороне клиента уязвимых приложений и перенаправляют его на новый, специально подготовленный полиморфный

файловый контейнер, обработка которого на атакуемой стороне приводит к „срабатыванию уязвимости“ и эксплуатации ее злоумышленником.

Авторами настоящей статьи ранее был предложен подход [6] к обнаружению подобных вредоносных файлов формата PDF (Portable Document Format) [7]. Суть данного метода заключается в использовании ряда статических признаков, свойственных, согласно отчетам исследовательского сообщества [8, 9], вредоносным документам, для формирования многомерного пространства атрибутов, характеризующих любой документ этого формата. Предложенный подход позволяет оценить применимость методов DM для задачи обнаружения вредоносных документов, эффективность отдельных методов обучения, а также эффективность применения групп признаков. Было показано, что совместное использование некоторых групп статических признаков и методов обучения позволяет создать модели обнаружения с показателями точности выше 90 %. Вместе с тем была отмечена необходимость дальнейшего исследования проблемы в контексте поиска дополнительных статических свойств, присущих вредоносным документам [6].

В настоящей статье рассматривается метод, развивающий предложенный подход за счет извлечения дополнительных статических признаков из файлов формата PDF. Метод основан на использовании структурных особенностей элементов документа и связей между ними.

В методологии обнаружения ВПО, которой посвящено множество публикаций, различают две основные группы используемых признаков — статическую и динамическую [10]. Статическая группа признаков включает в свой состав данные, извлекаемые из документов на основе анализа их содержимого и структуры. Динамические (также известные как поведенческие) признаки извлекаются из анализируемых документов с помощью процедуры их обработки интерпретирующей средой (исполнения) или моделью этой среды (эмуляции). В работах [10, 11] представлен анализ применимости для обнаружения ВПО байтовых цепочек (*n*-грамм). Структурные особенности различных форматов вредоносных файлов и их значимость обсуждаются в работах [10, 12]. Использование цепочек минимальных логических элементов файлов (например, опкодов машинных инструкций и их обобщений) рассматривается в работах [13,14], а пример анализа более крупных логических элементов опасных файлов (линейных блоков трансляции, внутренних процедур и их взаимосвязей с библиотечными процедурами) приведен в работах [14,15]. В частности, задачи сравнения графов взаимосвязей отдельных элементов вредоносных файлов обсуждаются в работе [15].

Предлагаемый в данной статье подход к извлечению признаков основан на типовых особенностях современных форматов документов, представляющих сложную иерархическую совокупность связанных элементов. Основным отличием рассматриваемого метода от существующих подходов к анализу структуры взаимосвязей документов является обнаружение отдельных пар элементов документа с последовательным формированием отдельного признака, характеризующего совокупностью свойств 1) главного элемента пары, 2) атрибута его связи с подчиненным элементом, 3) подчиненного элемента пары и длиной минимального маршрута от главного элемента пары до элемента графа, представляющего начальную структуру (так называемую „точку входа“) анализируемого документа.

**Формирование модели обнаружения ВПО.** Извлечение признаков из исходного набора данных является первым шагом общего процесса использования методов ML и DM для формирования модели обнаружения ВПО [16]. IDEF0-модель данного процесса представлена на рис. 1.

Процесс формирования модели содержит следующие шаги.

1. *Извлечение признаков*, формирующих базовый набор атрибутов, характеризующий общее пространство признаков, которые могут быть использованы при обучении модели. Как было показано ранее, именно этот шаг отвечает за извлечение из набора данных необходимой

для исследователя семантики атрибутов, определяющей идею общего подхода к обнаружению ВПО.

2. *Выделение значимых признаков*, обеспечивающее выбор из базового набора наиболее значимых атрибутов. Данный шаг позволяет уменьшить количество признаков, используемых для обучения модели, и обеспечивает необходимые для обучения приемлемые расходы вычислительных ресурсов.

3. *Обучение модели*, позволяющие сформировать модель обнаружения ВПО. Традиционно на данном шаге применяется один из выбранных исследователем методов классификации или кластеризации потока данных.

4. *Проверка (оценивание) модели*, обеспечивающая выходной контроль качественных параметров модели — точности обнаружения, времени работы и ресурсопотребления.

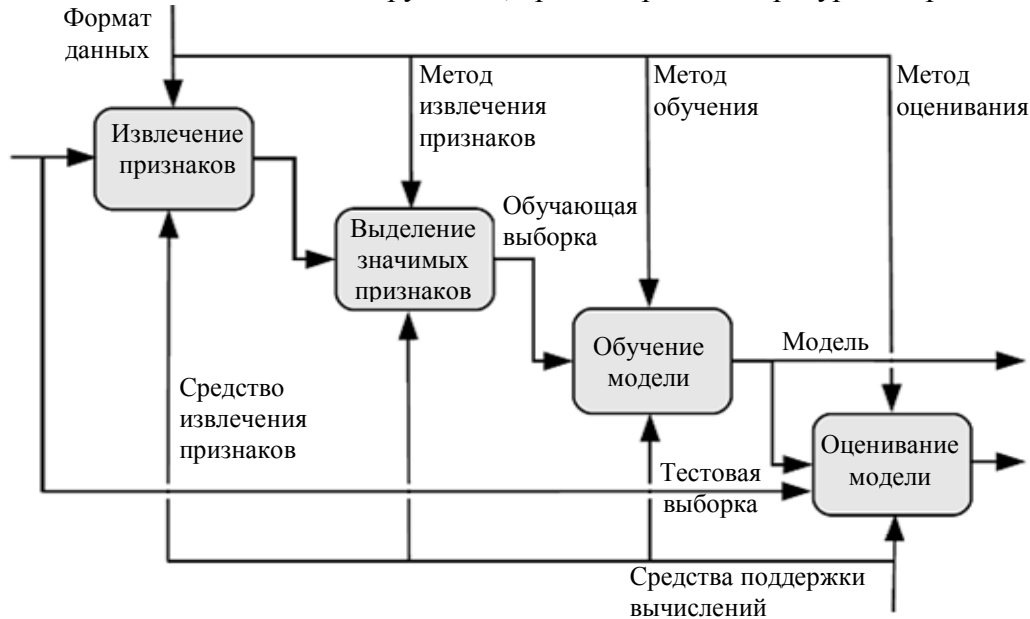


Рис. 1

**Метод извлечения признаков.** Общие особенности формата PDF представлены в работе [7]. В рамках рассматриваемого метода особый интерес представляют три внутренние группы сущностей, составляющих документ: это корневые элементы иерархии (стартовый объект и элементы таблиц взаимосвязей), косвенные объекты (*indirect objects*), хранящие основное содержание документа, и связи, характеризующие структурные отношения (роли) объектов. Пример структуры вредоносного документа представлен на рис. 2.

Точка входа в документ — стартовый объект (*Trailer*) — имеет ссылку на косвенный объект 1.0.obj с ролью *Root*. Данный объект имеет тип *Catalog* и ссылается на последующие в иерархии объекты, имеющие определенные роли (например, роль *Kids*, представляющая связь коллекции страниц документа в косвенном объекте 5.0.obj с косвенным объектом 8.0.obj представления отдельной страницы документа). В соответствии со спецификацией формата [7] косвенный объект также может иметь дополнительные внутренние свойства, определяющие его тип (общее свойство *Type*), структурные (например, свойство *Count* объекта 5.0.obj) и функциональные (свойство *S* объекта 10.0.obj) особенности.

Одной из структурных особенностей вредоносных документов является, как показано ниже, наличие объявленных ссылок на несуществующие косвенные объекты, представленные на рис. 2 в виде элементов иерархии, очерченных тонкими линиями.

Метод извлечения структурных признаков включает выполнение следующих шагов.

1. Построение структуры документа, представленной в виде ориентированного графа  $G=(V,E)$ , где множество  $V$  вершин определяется существующими ссылками на корневой и косвенные объекты, а множество  $E$  дуг определено множеством пар связанных вершин.

2. Каждый элемент множества  $E$  рассматривается в виде отдельного структурного признака, содержащего 1) свойства начальной вершины дуги, 2) тип ее связи с конечной вершиной и 3) свойства конечной вершины дуги.

3. При формировании конечного множества признаков каждый из элементов получает дополнительный атрибут, определяемый длиной минимального маршрута от начальной вершины соответствующего ему элемента множества  $V$  до вершины, характеризующей точку входа документа.

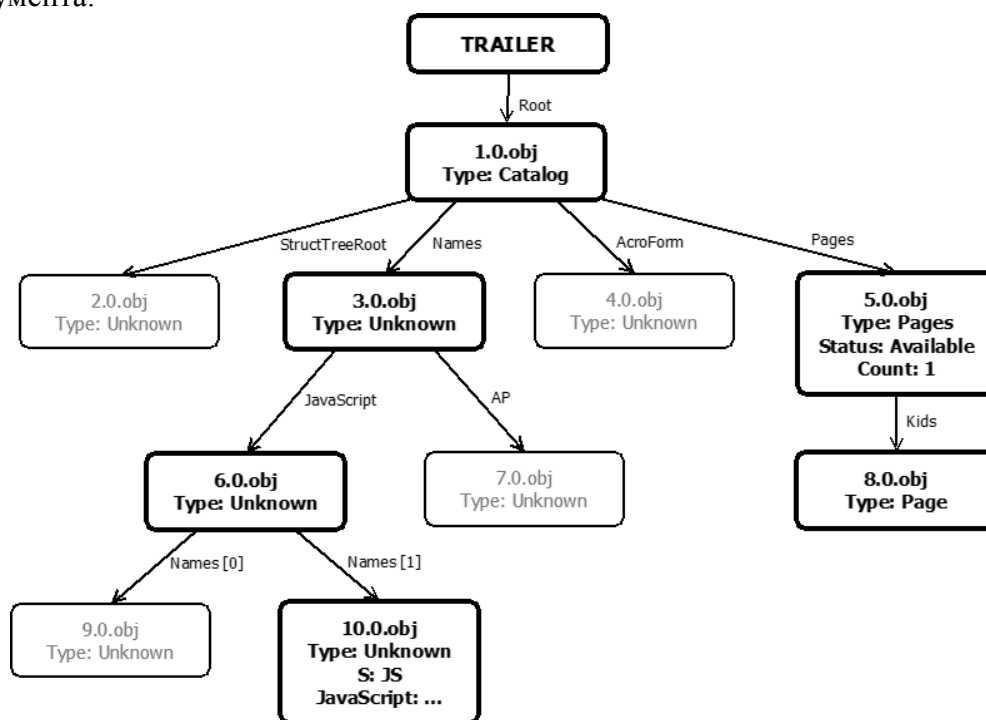


Рис. 2

Изложенный метод позволяет составить набор статических структурных признаков, наличие или отсутствие каждого из которых отображается булевой величиной в конечном характеризующем векторе каждого объекта обучающего набора. Первоначальная оценка эффективности метода была произведена на наборе данных, использованном авторами в работе [6]. В ходе экспериментов была разработана модель обнаружения вредоносных документов на основе метода классификации Decision Tree (дерево решений). Эксперименты проводились в соответствии с общей схемой выполнения работ [6] и поддерживались с помощью программных пакетов Open PDF Analysis Framework [17] и RapidMiner [18].

Точность классификации данных при использовании метода 10-кратной кросспроверки составила 94,9 %. Сопоставимый результат был получен в работе [6] при анализе применимости группы статических признаков, характеризующей используемые во вредоносных документах методы компрессии данных.

Результаты экспериментов позволили выявить ряд скрытых свойств вредоносных документов:

— значимость наличия ссылок на несуществующие в документе косвенные объекты и значимость изолированных цепочек взаимосвязанных объектов; наличие данных аномалий объясняется, по-видимому, спецификой программной реализации генераторов вредоносных документов: судя по полученным результатам, злоумышленники нацелены на создание вредоносных документов, имеющих существенные структурные отклонения, что препятствует их правильному разбору существующими средствами анализа;

— значительное количество устойчивых последовательностей признаков, имеющих отдельные переменные атрибуты косвенных взаимосвязанных объектов; обнаружение

современных вредоносных документов может быть осуществлено по признакам, не относящимся к структурным аномалиям или динамическому анализу вложений в документы; возможно, это объясняется спецификой реализации генераторов вредоносных документов, использующих при формировании новых поколений средств создания ВПО этого класса ограниченный набор базовых незначительно измененных программных компонентов.

**Заключение.** Представленный метод извлечения статических структурных признаков документов позволяет формировать системы обнаружения ВПО с показателями эффективности, сравнимыми с известными решениями [6, 8], а также выявить дополнительные, скрытые данные о структурных особенностях вредоносных документов.

Предложенный метод может быть применен также для документов альтернативных форматов (например, OLE2 и HTML). Использование данного подхода в комплексе с другими средствами обнаружения современных Web-угроз и противодействия им позволит эффективно обнаруживать потенциально опасные ресурсы сети Интернет.

Проведенные исследования предопределяют задачи для последующей серии экспериментов, направленных на более точное выявление значимых признаков отдельных структурных элементов документа и их взаимосвязей.

Статья подготовлена по результатам исследований, проводимых при финансовой поддержке Российского фонда фундаментальных исследований (проект №10-01-00826-а), программы фундаментальных исследований Отделения нанотехнологий и информационных технологий РАН (проект № 2.2) и государственного контракта 11.519.11.4008, а также при частичной финансовой поддержке, осуществляемой в рамках проектов Евросоюза “SecFutur” и MASSIF.

#### СПИСОК ЛИТЕРАТУРЫ

1. *Schultz M., Eskin E., Zadok E., Stolfo S.* Data mining methods for detection of new malicious executables // Proc. of the IEEE Symp. on Security and Privacy. Washington, DS, 2001. P. 38—49.
2. A Look at One Day of Malware Samples [Электронный ресурс]: <<http://blogs.mcafee.com/mcafee-labs/a-look-at-one-day-of-malware-samples>>.
3. *Ye Y., Li T.* Automatic malware categorization using cluster ensemble // Proc. of the 16th ACM SIGKDD Intern. Conf. on Knowledge Discovery and Data Mining. New York, NY, 2010. P. 95—104.
4. VXHeavens [Электронный ресурс]: <<http://www.vxheavens.com>>.
5. *Gibert K., Sanchez-Marré M., Codina V.* Choosing the right data mining technique: classification of methods and intelligent recommendation // Proc. of the IEMSS 5th Biennial Meeting: Intern. Congress on Environmental Modelling and Software. Ottawa, 2010. P. 1933—1940.
6. *Комашинский Д. В., Котенко И. В.* Обнаружение вредоносных документов формата PDF на основе интеллектуального анализа данных // Проблемы информационной безопасности. Компьютерные системы. 2012. № 1. С. 19—35.
7. International Organization for Standardization, Portable Document Format, ISO 32000-1:2008 [Электронный ресурс]: <[http://www.images.adobe.com/www.adobe.com/content/dam/Adobe/en/devnet/pdf/pdfs/PDF32000\\_2008.pdf](http://www.images.adobe.com/www.adobe.com/content/dam/Adobe/en/devnet/pdf/pdfs/PDF32000_2008.pdf)>.
8. *Kubec J., Sejtko J.* X is not enough! Grab the PDF by the tail! // Proc. of Virus Bulletin Annual Conf., Barselona, Oct. 2011. P. 128—135.
9. *Blonce A., Filiol E., Frayssignes L.* Portable document format (PDF) security analysis and malware threats // Presentations of Europe BlackHat Conf. Amsterdam, 2008.
10. *Kolter J. Z., Maloof M. A.* Learning to detect malicious executables in the wild // Proc. of the 10th ACM SIGKDD Intern. Conf. on Knowledge Discovery and Data Mining. New York, NY: ACM, 2004. P. 470—478.
11. *Masud M. M., Khan L., Thuraisingham B.* Feature-based techniques for auto-detection of novel email worms // Proc. of the 11th Pacific-Asia Conf. on Knowledge Discovery and Data Mining. Springer-Verlag Berlin, Heidelberg, 2007. P. 205—216.

12. *Shahzad F., Farooq M.* ELF-Miner: Using structural knowledge and data mining methods to detect new (linux) malicious executables // Knowledge and Information Systems. 2011. Vol. 30 (3). P. 589—612.
13. *Siddiqui M., Wang M., Lee J.* Detecting Internet worms using data mining techniques // J. of Systemics, Cybernetics and Informatics. 2008. Vol. 6, N 6. P. 48—53.
14. *Ye Y., Li T., Huang K., Jiang Q., Chen Y.* Hierarchical associative classifier (HAC) for malware detection from the large and imbalanced gray list // J. of Intelligent Information Systems. 2010. Vol. 35, Iss. 1. P. 1—20.
15. *Lanzi A., Balzarotti D., Kruegel C., Christodorescu M., Kirda E.* AccessMiner: Using System-Centric Models for Malware Protection. New York, NY: ACM, 2010. P. 399—412.
16. *Комашинский Д. В., Котенко И. В.* Концептуальные основы использования методов Data Mining для обнаружения вредоносного программного обеспечения // Защита информации. Инсайд. 2010. № 2. С. 74—82.
17. Open PDF Analysis Framework [Электронный ресурс]: <<http://code.google.com/p/opaf/>>.
18. Rapid – I Rapid Miner 5 [Электронный ресурс]: <<http://rapid-i.com/content/view/181/190/>>.

**Сведения об авторах****Дмитрий Владимирович Комашинский**

— аспирант; СПИИРАН, лаборатория проблем компьютерной безопасности; E-mail: komashinskiy@gmail.com

**Игорь Витальевич Котенко**

— д-р техн. наук, профессор; СПИИРАН, лаборатория проблем компьютерной безопасности; E-mail: ivkote@comsec.spb.ru

Рекомендована СПИИРАН

Поступила в редакцию  
10.06.12 г.