

Отчет о проверке на заимствования №1



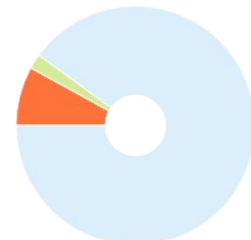
Автор: Соленая Оксана Ярославовна osolenaya@list.ru / ID: 12
Проверяющий: Соленая Оксана Ярославовна (osolenaya@list.ru / ID: 12)
Организация: Санкт-Петербургский государственный университет аэрокосмического приборостроения
 Отчет предоставлен сервисом «Антиплагиат» - <http://guap.antiplagiat.ru>

ИНФОРМАЦИЯ О ДОКУМЕНТЕ

№ документа: 1098
 Начало загрузки: 06.12.2019 19:04:23
 Длительность загрузки: 00:00:24
 Корректировка от 06.12.2019 19:29:15
 Имя исходного файла:
 Диссертация_Салахутдинова
 Размер текста: 7186 кБ
 Символов в тексте: 231646
 Слов в тексте: 27432
 Число предложений: 1666

ИНФОРМАЦИЯ ОБ ОТЧЕТЕ

Последний готовый отчет (ред.)
 Начало проверки: 06.12.2019 19:04:48
 Длительность проверки: 00:00:43
 Комментарии: [Автосохраненная версия]
 Модули поиска: Модуль поиска ИПС "Адилет", Модуль выделения библиографических записей, Сводная коллекция ЭБС, Коллекция РГБ, Цитирование, Модуль поиска переводных заимствований, Модуль поиска переводных заимствований по eLibrary (EnRu), Модуль поиска переводных заимствований по интернет (EnRu), Модуль поиска переводных заимствований по Wiley (RuEn), Коллекция eLIBRARY.RU, Коллекция ГАРАНТ, Модуль поиска "ГУАП", Модуль поиска Интернет, Коллекция Медицина, Модуль поиска перефразирований eLIBRARY.RU, Модуль поиска перефразирований Интернет, Коллекция Патенты, Модуль поиска общепотребительных выражений, Коллекция Wiley



ЗАИМСТВОВАНИЯ	ЦИТИРОВАНИЯ	ОРИГИНАЛЬНОСТЬ
7,81% 	2,04% 	90,15% 

Заимствования — доля всех найденных текстовых пересечений, за исключением тех, которые система отнесла к цитированиям, по отношению к общему объему документа.
 Цитирования — доля текстовых пересечений, которые не являются авторскими, но система посчитала их использование корректным, по отношению к общему объему документа. Сюда относятся оформленные по ГОСТу цитаты; общепотребительные выражения; фрагменты текста, найденные в источниках из коллекций нормативно-правовой документации.
 Текстовое пересечение — фрагмент текста проверяемого документа, совпадающий или почти совпадающий с фрагментом текста источника.
 Источник — документ, проиндексированный в системе и содержащийся в модуле поиска, по которому проводится проверка.
 Оригинальность — доля фрагментов текста проверяемого документа, не обнаруженных ни в одном источнике, по которым шла проверка, по отношению к общему объему документа.
 Заимствования, цитирования и оригинальность являются отдельными показателями и в сумме дают 100%, что соответствует всему тексту проверяемого документа.
 Обращаем Ваше внимание, что система находит текстовые пересечения проверяемого документа с проиндексированными в системе текстовыми источниками. При этом система является вспомогательным инструментом, определение корректности и правомерности заимствований или цитирований, а также авторства текстовых фрагментов проверяемого документа остается в компетенции проверяющего.

№	Доля в отчете	Доля в тексте	Источник	Ссылка	Актуален на	Модуль поиска	Блоков в отчете	Блоков в тексте
[01]	0%	4,53%	ИДЕНТИФИКАЦИИ ПРОГРАММНОГО О.	http://elibrary.ru	14 Окт 2019	Коллекция eLIBRARY.RU	0	36
[02]	0%	3,21%	ПОДХОД К ВЫБОРУ ИНФОРМАТИВНО...	https://yandex.ru	31 Окт 2018	Модуль поиска Интернет	0	58
[03]	0%	2,51%	Скачать PDF	https://openbooks.ifmo.ru	24 Окт 2019	Модуль поиска Интернет	0	58
[04]	0%	2,34%	Подход к выбору информативного пр...	http://elibrary.ru	16 Июл 2018	Коллекция eLIBRARY.RU	0	44
[05]	0%	1,55%	ИССЛЕДОВАНИЕ ВЛИЯНИЯ ВЫБОРА П...	http://elibrary.ru	27 Мая 2019	Коллекция eLIBRARY.RU	0	33
[06]	0%	1,09%	https://esu.citis.ru/dissertation/CTF0KYQ...	https://esu.citis.ru	10 Мая 2018	Модуль поиска Интернет	0	24
[07]	0,52%	0,8%	Воробьева, Алиса Андреевна Методик...	http://dlib.rsl.ru	15 Дек 2017	Коллекция РГБ	1214	20
[08]	0,4%	0,7%	Годовой отчет 2018 года	http://spiiras.nw.ru	21 Окт 2019	Модуль поиска Интернет	926	15
[09]	0,37%	0,69%	Модели и методы обнаружения наруш.	http://spiiras.nw.ru	06 Ноя 2018	Модуль поиска Интернет	867	12
[10]	0%	0,68%	Викснин, Илья Игоревич Модели и ме...	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	0	13
[11]	0,51%	0,64%	Математические модели и алгоритмы...	http://netess.ru	30 Янв 2017	Модуль поиска перефразирований Интернет	1191	5
[12]	0%	0,63%	RESEARCH at ITMO University	http://research.ifmo.ru	29 Ноя 2019	Модуль поиска Интернет	0	6
[13]	0,03%	0,62%	Нормативное обеспечение информац..	http://dlib.rsl.ru	04 Дек 2017	Коллекция РГБ	66	13
[14]	0%	0,6%	Комплексная защита информации в о...	https://book.ru	03 Июл 2017	Сводная коллекция ЭБС	0	12
[15]	0%	0,6%	Комплексная защита информации в о...	https://book.ru	03 Июл 2017	Сводная коллекция ЭБС	0	12
[16]	0,34%	0,59%	https://esu.citis.ru/dissertation/RNHXAO...	https://esu.citis.ru	20 Мар 2018	Модуль поиска Интернет	784	14

[17]	0,57%	0,57%	On challenges in evaluating malware clus..	https://cs.unc.edu	07 Янв 2018	Модуль поиска переводных заимствований	1330	2
[18]	0%	0,55%	Защита информационных технологий..	https://book.ru	03 Июл 2017	Сводная коллекция ЭБС	0	10
[19]	0%	0,52%	Ложников, Павел Сергеевич Методоло..	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	0	10
[20]	0,3%	0,5%	Objdump & readelf使用-静默梧桐-ChinaU..	http://blog.chinaunix.net	02 Апр 2018	Модуль поиска Интернет	690	15
[21]	0,46%	0,46%	Анализ методов распознавания вредо...	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	1077	3
[22]	0,18%	0,43%	Постановление администрации муниц..	http://ivo.garant.ru	14 Мая 2018	Коллекция ГАРАНТ	411	14
[23]	0%	0,43%	Постановление администрации Алекса..	http://municipal.garant.ru	15 Мая 2018	Коллекция ГАРАНТ	0	14
[24]	0,11%	0,42%	Сорокин, Иван Витальевич диссертаци..	http://dlib.rsl.ru	раньше 2011	Коллекция РГБ	254	10
[25]	0,12%	0,42%	Лившиц, Илья Иосифович Модели и м..	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	267	10
[26]	0%	0,4%	Методика защиты информации в орга..	https://book.ru	03 Июл 2017	Сводная коллекция ЭБС	10	9
[27]	0,07%	0,4%	Кацупеев, Андрей Александрович Сист..	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	163	5
[28]	0,1%	0,39%	Антонов, Алексей Евгеньевич диссериа..	http://dlib.rsl.ru	25 Дек 2015	Коллекция РГБ	228	9
[29]	0,14%	0,38%	Прикладная информатика	http://ibooks.ru	раньше 2011	Сводная коллекция ЭБС	334	8
[30]	0,03%	0,38%	Идентификация типа файла на основе...	https://yandex.ru	28 Дек 2018	Модуль поиска Интернет	80	8
[31]	0%	0,38%	не указано	https://narfu.ru	15 Дек 2018	Модуль поиска Интернет	0	9
[32]	0%	0,37%	Нурдинов Руслан Артурович. Модель ...	https://docplayer.ru	14 Мая 2019	Модуль поиска Интернет	0	7
[33]	0,13%	0,36%	Комашинский, Дмитрий Владимирови..	http://dlib.rsl.ru	25 Дек 2015	Коллекция РГБ	301	9
[34]	0%	0,3%	Юрьева, Радда Алексеевна Метод и мо..	http://dlib.rsl.ru	15 Апр 2018	Коллекция РГБ	0	9
[35]	0,17%	0,3%	Национальный стандарт РФ ГОСТ Р ИС..	http://ivo.garant.ru	13 Янв 2017	Коллекция ГАРАНТ	391	7
[36]	0,11%	0,29%	All publications Faculty of Technological ...	http://ftmi.ifmo.ru	09 Фев 2019	Модуль поиска Интернет	265	6
[37]	0,29%	0,29%	Идентификация типа файла на основе...	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	670	2
[38]	0,15%	0,28%	Киричек, Руслан Валентинович Разраб..	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	359	4
[39]	0,01%	0,28%	Модели и алгоритмы анализа информ...	https://ugatu.su	29 Окт 2019	Модуль поиска Интернет	29	5
[40]	0,12%	0,27%	Приказ Управления ветеринарии г. Се...	http://ivo.garant.ru	14 Авг 2018	Коллекция ГАРАНТ	272	8
[41]	0,03%	0,27%	Эдель, Дмитрий Александрович диссер..	http://dlib.rsl.ru	раньше 2011	Коллекция РГБ	76	5
[42]	0,09%	0,25%	Конфиденциальная информация	https://knowledge.allbest.ru	28 Фев 2019	Модуль поиска Интернет	202	5
[43]	0,04%	0,25%	Безопасность информации в деятельн..	http://elibrary.ru	23 Сен 2015	Коллекция eLIBRARY.RU	104	5
[44]	0%	0,25%	Исаев, Александр Сергеевич Метод и м..	http://dlib.rsl.ru	22 Авг 2019	Коллекция РГБ	0	5
[45]	0,02%	0,24%	Агеев, Сергей Александрович Модели,...	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	42	5
[46]	0,23%	0,23%	2.1.2. Starting A Process	http://d3s.mff.cuni.cz	07 Янв 2018	Модуль поиска переводных заимствований	530	3
[47]	0%	0,23%	Динамический анализ программ с цел..	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	0	1
[48]	0,23%	0,23%	Полный текст статьи в формате pdf (1/...	http://ispras.ru	01 Янв 2017	Модуль поиска перефразирований Интернет	522	1
[49]	0%	0,21%	Хакимов, Дмитрий Валерьевич Автом...	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	10	6
[50]	0%	0,21%	https://esu.citis.ru/dissertation/2rWl000...	https://esu.citis.ru	21 Мар 2018	Модуль поиска Интернет	0	3
[51]	0,2%	0,2%	Malware detection using bilayer behavio...	https://doi.org	31 Авг 2018	Модуль поиска Интернет	472	5
[52]	0,06%	0,2%	264399	http://e.lanbook.com	10 Мар 2016	Сводная коллекция ЭБС	133	5
[53]	0%	0,2%	Новохрестов, Алексей Константинович.	http://dlib.rsl.ru	15 Окт 2019	Коллекция РГБ	0	3
[54]	0,1%	0,19%	Высокопроизводительные параллельн.	https://docplayer.ru	19 Июн 2019	Модуль поиска Интернет	242	4

[55]	0,05%	0,19%	Динамический анализ программ с цел..	http://elibrary.ru	25 Дек 2016	Коллекция eLIBRARY.RU	117	4
[56]	0,12%	0,18%	1117	http://e.lanbook.com	09 Мар 2016	Сводная коллекция ЭБС	288	6
[57]	0,17%	0,17%	62575	http://e.lanbook.com	09 Мар 2016	Сводная коллекция ЭБС	391	7
[58]	0%	0,15%	Волокитина, Евгения Сергеевна диссер.	http://dlib.rsl.ru	28 Фев 2014	Коллекция РГБ	0	3
[59]	0%	0,15%	ОСНОВНЫЕ НАПРАВЛЕНИЯ РАЗВИТИЯ..	http://elibrary.ru	раньше 2011	Коллекция eLIBRARY.RU	0	4
[60]	0%	0,14%	МЕТОД ИДЕНТИФИКАЦИИ ИСПОЛНЯЕ..	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	0	2
[61]	0,04%	0,13%	http://books.ifmo.ru/file/pdf/1570.pdf	http://books.ifmo.ru	14 Мая 2019	Модуль поиска Интернет	99	3
[62]	0%	0,13%	Актуальные проблемы инфотелекомм..	http://elibrary.ru	05 Дек 2015	Коллекция eLIBRARY.RU	0	4
[63]	0,07%	0,12%	Постановление администрации Ребри...	http://municipal.garant.ru	16 Июл 2019	Коллекция ГАРАНТ	173	5
[64]	0,04%	0,12%	66085	http://e.lanbook.com	09 Мар 2016	Сводная коллекция ЭБС	96	4
[65]	0,12%	0,12%	XGBoost使用说明 - CSDN博客	http://blog.csdn.net	08 Апр 2018	Модуль поиска Интернет	278	3
[66]	0%	0,12%	Информационные системы и техноло...	http://studentlibrary.ru	27 Ноя 2017	Сводная коллекция ЭБС	0	3
[67]	0,12%	0,12%	Методы и средства анализа информат...	http://swsys.ru	05 Янв 2017	Модуль поиска перефразирований Интернет	271	1
[68]	0,04%	0,11%	Путин, Евгений Олегович Глубокие ген.	http://dlib.rsl.ru	22 Фев 2019	Коллекция РГБ	104	3
[69]	0,08%	0,11%	Бытко, Сергей Юрьевич Эффективнос...	http://dlib.rsl.ru	14 Июн 2019	Коллекция РГБ	184	3
[70]	0,11%	0,11%	Подвербный, Вячеслав Анатольевич д...	http://dlib.rsl.ru	раньше 2011	Коллекция РГБ	258	2
[71]	0%	0,11%	АНАЛИЗ ИНФОРМАЦИОННОЙ БЕЗОПАС...	http://elibrary.ru	05 Авг 2016	Коллекция eLIBRARY.RU	0	3
[72]	0%	0,11%	О подписании Соглашения об обеспеч.	http://adilet.zan.kz	21 Янв 2016	Модуль поиска ИПС "Адилет"	0	3
[73]	0%	0,1%	Идентификация типа файла на основе...	http://elibrary.ru	17 Дек 2016	Коллекция eLIBRARY.RU	0	1
[74]	0%	0,1%	Predictive bitmaps	http://halobates.de	05 Янв 2018	Модуль поиска переводных заимствований	0	1
[75]	0%	0,1%	[Gzip'd Text 322 KB]	http://os.inf.tu-dresden.de	07 Янв 2018	Модуль поиска переводных заимствований	0	1
[76]	0%	0,1%	Формирование словаря терминов тео...	http://elibrary.ru	раньше 2011	Коллекция eLIBRARY.RU	0	2
[77]	0,01%	0,1%	"Функциональная сигнатура компьют...	http://tusur.ru	30 Янв 2017	Модуль поиска перефразирований Интернет	25	1
[78]	0,09%	0,09%	ПРИМЕНЕНИЕ ИСКУССТВЕННЫХ НЕЙР..	http://elibrary.ru	11 Мая 2018	Коллекция eLIBRARY.RU	215	2
[79]	0%	0,09%	ОЦЕНКА РЕФОРМ В СФЕРЕ ПОЛОВЫХ П.	http://elibrary.ru	раньше 2011	Коллекция eLIBRARY.RU	0	3
[80]	0,03%	0,09%	АНАЛИЗ И КЛАССИФИКАЦИЯ МЕТОДО..	http://elibrary.ru	05 Авг 2016	Коллекция eLIBRARY.RU	64	2
[81]	0,09%	0,09%	https://sourceware.org/bugzilla/show_b...	https://sourceware.org	04 Июн 2018	Модуль поиска Интернет	209	1
[82]	0,09%	0,09%	Анализ методов распознавания вредо...	http://elibrary.ru	04 Фев 2015	Коллекция eLIBRARY.RU	208	2
[83]	0%	0,09%	РЕЦЕНЗИЯ НА КНИГУ «СЕМЬ БЕЗОПАС...	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	0	1
[84]	0%	0,08%	РЕЦЕНЗИЯ НА КНИГУ «СЕМЬ БЕЗОПАС...	http://elibrary.ru	03 Мая 2017	Коллекция eLIBRARY.RU	0	1
[85]	0%	0,08%	Методические рекомендации для орга...	http://ivo.garant.ru	14 Янв 2017	Коллекция ГАРАНТ	0	2
[86]	0,03%	0,08%	Lecture 4 Notes - Computer Science 473 ...	https://studyblue.com	08 Янв 2018	Модуль поиска переводных заимствований	60	1
[87]	0,08%	0,08%	Методы и средства анализа информат...	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	185	1
[88]	0%	0,08%	CS590-SWS/527 Software Security - PDF	http://docplayer.net	08 Янв 2018	Модуль поиска переводных заимствований	0	1
[89]	0,07%	0,07%	Параллельная реализация алгоритма ...	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	173	1
[90]	0%	0,07%	aft	http://people.redhat.com	06 Янв 2018	Модуль поиска переводных заимствований	0	1

[91]	0,08%	0,07%	ЭБС «Издательство «Лань»	https://e.lanbook.com	05 Янв 2017	Модуль поиска перефразирований Интернет	177	1
[92]	0,07%	0,07%	Постановление Администрации муниц...	http://ivo.garant.ru	10 Апр 2019	Коллекция ГАРАНТ	162	1
[93]	0,07%	0,07%	РЕАЛИЗАЦИЯ ПАРАЛЛЕЛЬНОГО АЛГО...	http://cyberleninka.ru	30 Янв 2017	Модуль поиска перефразирований Интернет	165	1
[94]	0,03%	0,06%	Тютин Д.В. Налоговое право. Курс лекц.	http://ivo.garant.ru	12 Янв 2017	Коллекция ГАРАНТ	79	2
[95]	0,06%	0,06%	АЛГОРИТМ ВЫРАВНИВАНИЯ ПОСЛЕДО.	http://elibrary.ru	02 Янв 2018	Модуль поиска перефразирований eLIBRARY.RU	145	1
[96]	0,06%	0,06%	Рекомендации по межгосударственно...	http://ivo.garant.ru	14 Янв 2017	Коллекция ГАРАНТ	144	1
[97]	0%	0,06%	Apparatus and method accumulating ca...	http://freepatentsonline.com	04 Ноя 2016	Коллекция Патенты	0	1
[98]	0%	0,06%	Method and system for object detection...	http://freepatentsonline.com	05 Ноя 2016	Коллекция Патенты	0	1
[99]	0%	0,06%	Apparatus and method accumulating ca...	http://freepatentsonline.com	06 Ноя 2016	Коллекция Патенты	0	1
[100]	0%	0,06%	Method and system for object detection...	http://freepatentsonline.com	06 Ноя 2016	Коллекция Патенты	0	1
[101]	0%	0,06%	ELASTOMERIC OPTICAL TACTILE SENSOR...	http://freepatentsonline.com	05 Ноя 2016	Коллекция Патенты	0	1
[102]	0%	0,06%	From word models to executable model...	https://doi.org	14 Ноя 2019	Коллекция Wiley	0	1
[103]	0%	0,06%	автореферат	http://mpei.ru	05 Янв 2017	Модуль поиска перефразирований Интернет	0	1
[104]	0%	0,05%	Threatened species impact assessments...	https://doi.org	14 Ноя 2019	Коллекция Wiley	0	1
[105]	0,05%	0,05%	Predicting species distributions from mu...	https://doi.org	14 Ноя 2019	Коллекция Wiley	127	1
[106]	0%	0,05%	Effects of sample size on the performanc...	https://doi.org	14 Ноя 2019	Коллекция Wiley	0	1
[107]	0%	0,05%	Incorporating habitat use in models of fa...	https://doi.org	14 Ноя 2019	Коллекция Wiley	0	1
[108]	0%	0,05%	СТАТИСТИЧЕСКИЕ МЕТОДЫ В ПЕДАГОГ.	http://elibrary.ru	раньше 2011	Коллекция eLIBRARY.RU	0	1
[109]	0%	0,05%	Способ формирования испытательны...	http://findpatent.ru	25 Июн 2015	Коллекция Патенты	0	1
[110]	0%	0,05%	СПОСОБ ФОРМИРОВАНИЯ ИСПЫТАТЕ...	http://ru-patent.info	25 Июн 2015	Коллекция Патенты	0	1
[111]	0,05%	0,05%	METHODS AND SYSTEMS FOR INDOOR N...	http://freepatentsonline.com	06 Ноя 2016	Коллекция Патенты	125	1
[112]	0%	0,05%	SD-WLB: An SDN-aided mechanism for w...	https://doi.org	14 Ноя 2019	Коллекция Wiley	0	1
[113]	0%	0,05%	Способ защиты информационно-вычи.	http://findpatent.ru	24 Июн 2015	Коллекция Патенты	0	1
[114]	0%	0,05%	Способ защиты вычислительных сетей.	http://findpatent.ru	раньше 2011	Коллекция Патенты	0	1
[115]	0,05%	0,05%	Learnable non-darwinian evolution - Geo...	http://freepatentsonline.com	09 Ноя 2016	Коллекция Патенты	115	1
[116]	0%	0,05%	Способ передачи и приема дискретно...	http://findpatent.ru	24 Июн 2015	Коллекция Патенты	0	1
[117]	0,05%	0,05%	Comparing deep belief networks with su...	https://doi.org	07 Янв 2019	Модуль поиска переводных заимствований по Wiley (RuEn)	114	1
[118]	0%	0,05%	MIPRO'2015. 38th International Convent...	http://elibrary.ru	11 Июл 2015	Коллекция eLIBRARY.RU	0	1
[119]	0%	0,05%	SYSTEM AND METHOD FOR ROBUST EVA...	http://freepatentsonline.com	04 Ноя 2016	Коллекция Патенты	0	2
[120]	0,05%	0,05%	Распределение стран мира по ВВП на д.	http://elibrary.ru	29 Апр 2017	Коллекция eLIBRARY.RU	107	1
[121]	0%	0,04%	An Introduction to Boosting and Leverag...	http://elibrary.ru	раньше 2011	Коллекция eLIBRARY.RU	0	1
[122]	0%	0,04%	Способ моделирования сетей связи. П...	http://findpatent.ru	24 Июн 2015	Коллекция Патенты	0	1
[123]	0%	0,04%	SYSTEM AND METHOD FOR PROFILING J...	http://freepatentsonline.com	06 Ноя 2016	Коллекция Патенты	0	1
[124]	0,04%	0,04%	Identification and expression of the 11β...	https://doi.org	14 Ноя 2019	Коллекция Wiley	97	1
[125]	0,04%	0,04%	ПРИМЕНЕНИЕ ЛОГЛИНЕЙНОГО АНАЛ...	http://elibrary.ru	27 Дек 2016	Коллекция eLIBRARY.RU	93	1
[126]	0%	0,04%	Защита от взлома: сокеты, эксплойты,...	http://studentlibrary.ru	19 Дек 2016	Коллекция Медицина	0	1
[127]	0%	0,04%	Обоснование прогностической оценк...	http://emll.ru	20 Дек 2016	Коллекция Медицина	0	1
[128]	0%	0,04%	Алгоритмы и модели ограничения дос...	http://studentlibrary.ru	20 Дек 2016	Коллекция Медицина	0	1

[129]	0%	0,04%	Постановление Правления ОАО "Газпр..	http://ivo.garant.ru	01 Мар 2018	Коллекция ГАРАНТ	0	1
[130]	0,04%	0,04%	LoGos: Internet-Explorer-Based Maliciou...	https://doi.org	14 Ноя 2019	Коллекция Wiley	82	1
[131]	0%	0,03%	ФИПС - Федеральное государственное...	http://www1.fips.ru	25 Июн 2015	Коллекция Патенты	0	1
[132]	0,01%	0,03%	Automated Analysis of Text in Graduate ...	https://doi.org	14 Ноя 2019	Коллекция Wiley	19	1
[133]	0%	0,03%	Trait-based risk assessment for invasive ...	https://doi.org	14 Ноя 2019	Коллекция Wiley	0	1
[134]	0%	0,03%	Определен порядок формирования за...	http://ivo.garant.ru	21 Фев 2019	Коллекция ГАРАНТ	5	2
[135]	0,02%	0,02%	Comparative analysis of respiratory mot...	https://doi.org	14 Ноя 2019	Коллекция Wiley	56	1
[136]	0,02%	0,02%	Кашепов В.П., Голованова Н.А., Гравин..	http://ivo.garant.ru	21 Фев 2019	Коллекция ГАРАНТ	36	1
[137]	0%	0%	не указано	не указано	раньше 2011	Модуль выделения библиографических записей	0	3
[138]	0,46%	0%	не указано	не указано	раньше 2011	Цитирование	1057	13
[139]	0,86%	0%	не указано	не указано	раньше 2011	Модуль поиска общеупотребительных выражений	1991	115

Текст документа

Федеральное государственное бюджетное учреждение науки **139**

Санкт-Петербургский институт информатики и автоматизации

Российской академии наук

(СПИИРАН)

На правах рукописи **25**

Салахутдинова Ксения Иркиновна

Методика идентификации исполняемых файлов на основе статического

анализа характеристик дизассемблированного кода программ **24**

Специальность 05.13.19 – Методы и системы защиты информации,

информационная безопасность

Диссертация на соискание ученой степени

кандидата технических наук

Научный руководитель:

доктор технических наук, профессор **33**

Лебедев Илья Сергеевич

Санкт-Петербург – 2019 **9**

2

Оглавление

Введение 6

Глава 1. Постановка задачи обеспечения **9** безопасности автоматизированной

системы от потенциальных угроз со стороны несанкционированно

установленного программного обеспечения 15

1.1 Обзор проблематики современных подходов к идентификации

программного обеспечения 15

1.2 Модель угроз информационной безопасности при обработке информации

конфиденциального характера в информационной системе 18

1.3 Модель вероятного нарушителя информационной безопасности и оценка

его возможностей 22

Выводы по главе 1 25

Глава 2. Обзор существующих методов идентификации программного

обеспечения 27

2.1 Методы идентификации, основанные на статическом характере анализа

характеристик программного обеспечения 29

2.2 Методы идентификации, основанные на динамическом характере сбора

характеристик программного обеспечения 37

2.3 Методы идентификации, основанные на сборе информации путем

обращения к встроенным функциям операционной системы или задействования собственного программного агента	39
2.4 Постановка задачи исследования	48
Выводы по главе 2	52
Глава 3. Разработка методики идентификации установленного программного обеспечения на основе статических характеристик программного кода файла	54
3.1 Анализ структуры и характеристик исполняемого файла	55
3	
3.1.1 Представление исходного кода на языке ассемблера	59
3.1.2 Анализ особенностей ассемблерного кода программного обеспечения	61
3.2 Подход к представлению программного обеспечения.....	63
3.2.1 Модель представления исполняемого файла	65
3.2.2 Идентификационное признаковое пространство	66
3.2.3 Модель представления программного обеспечения	68
3.3 Метод формирования сигнатур исполняемых файлов, обладающий достаточной гибкостью по распознаванию различных версий программ, а также позволяющий наиболее точно производить идентификацию	69
3.3.1 Критерии и алгоритм отбора наиболее информативных признаков	70
3.3.2 Формирование унифицированных сигнатур программ (УСП), используемого в дальнейшем в статистическом методе сравнения сигнатур и с использованием искусственной нейронной сети	74
3.3.3 Формирование сигнатур программ (СП), используемого в дальнейшем в методе сравнения сигнатур на основе градиентного бустинга деревьев решений	81
3.4 Метод сравнения сигнатуры идентифицируемого исполняемого файла (СИИФ) с эталонной сигнатурой программы (ЭСП)	82
3.4.1 Статистические методы сравнения сигнатур на основе критерия однородности хи-квадрат и критерия Колмогорова.....	83
3.4.2 Метод сравнения сигнатур на основе машинного обучения.....	86
3.4.3 Оценка вычислительной сложности различных методов сравнения СИИФ с ЭСП	91
3.5 Методика идентификации программного обеспечения на основе статических характеристик программного кода файла	92
4	
3.5.1 Предварительный этап методики идентификации программного обеспечения. Сбор и формирование обучающей выборки и составление АЭСП, содержащий ЭСП или унифицированные ЭСП	94
3.5.2 Основные этапы методики идентификации исполняемого файла ...	96
3.6 Оценка точности идентификации программного обеспечения	98
3.7 Ограничения методики	100
Выводы по главе 3	100
Глава 4. Экспериментальная проверка точности идентификации программного обеспечения и анализ полученных результатов	102
4.1 Входные данные экспериментов. Обучающая и тестовая выборки	102
4.2 Определение точности идентификации при использовании разработанных методов сравнения СИИФ с ЭСП	103
4.2.1 Точность идентификации при использовании статистических методов сравнения сигнатур	103
4.2.2 Точность идентификации при использовании метода сравнения сигнатур на основе машинного обучения	111
4.2.3 Повышение точности идентификации путем использования аддитивного критерия Фишберна	117
4.3 Определение итоговой точности идентификации программного	

обеспечения на основе предложенной методики	120
4.4 Использование результатов исследования для повышения информационной безопасности автоматизированных систем	123 27
Выводы по главе 4	126
Заключение	128
Список сокращений и условных обозначений	131
Список 38 использованной 27 литературы	133

5	
Приложение 1. Значения информативности для 118 ассемблерных команд	142
Приложение 2. Числовые идентификаторы эталонных программ	143
Приложение 3. Коды сравнения сигнатур с помощью библиотек градиентного бустинга деревьев решений	145
Приложение 4. Свидетельство о регистрации программы для ЭВМ	150
Приложение 5. Копии актов внедрения.....	156
Приложение 6. Публикации соискателя по теме Диссертации	159

6

Введение

Актуальность работы. Свободно распространяемое программное обеспечение (ПО) является неотъемлемой частью современных информационных систем, эксплуатируемых в различных секторах экономики, оно позволяет расширить возможности по осуществлению процессов анализа, управления, принятия решений.

Применение открытого ПО обуславливает необходимость разработки дополнительных методов, систем и средств защиты информации. Возможные дефекты программного обеспечения, наличие не декларированных возможностей, нелегальное использование интеллектуальной собственности, применение специальных программ, направленных на преодоление установленной защиты, могут привести к росту числа уязвимостей и повлиять на информационную безопасность систем. 139

В связи с этим возникает необходимость 139 решения ряда задач идентификации, верификации и валидации программного обеспечения. Предлагаемые решения ориентированы, в основном, на отслеживание фиксированного состояния кода программ на носителях и в оперативной памяти, что не всегда позволяет оперативно определить санкционированные модификации, изменения версий. Данная работа направлена на решение задачи идентификации программного обеспечения, в условии изменения версий распространяемого ПО, на основе процедуры построения информативной модели в виде математического кортежа по выбранному признаковому пространству, характеристики которой позволяют найти однозначное соответствие между анализируемой последовательностью и хранящимся эталоном исполняемого файла.

Степень разработанности темы. Различным подходам идентификации программного обеспечения посвятили свои работы такие отечественные ученые как Казарин О. В., Копыльцов А.В., Сорокин И.В., Антонов А. Е., Федулов А. С., зарубежные – Kornblum J.D., Long C., Ebringer T., Sun L., Boztas S., Schultz M. G., Eskin E., Santos I. и другие.

7

Существующие подходы к идентификации исполняемых файлов, по способу анализа характеристик программы, можно разделить на статические и динамические. К динамическим относятся методы, рассматривающие программу как процесс выполнения некоторого алгоритма или как последовательную смену состояний на графе переходов программы 21. Статические, в свою очередь, делятся на 139 не сигнатурные (методы проверки целостности): сравнение контрольных сумм, с полной копией данных, вычисление CRC (Cyclic Redundancy Check), хэш-сумм, использование имитовставки, цифровой подписи; и сигнатурные (методы

сравнения признаковых последовательностей): выделение «магического числа», сопоставление характерных последовательностей байтов и др.

В качестве информативных признаков файла в отечественных и зарубежных научных работах рассматриваются такие признаки как операторы в тексте программ (Казарин О. В.), энтропия сегментов файлов, цифровые изображения (Копыльцов А.В., Сорокин И.В.), блоки в сегментированной программе (Антонов А.Е., Федулов А.С.), хеш-функции и кусочно-зависимое хеширование (Kornblum J.D., Long C., Guoyin W., Baier H., Frank B.), статистические профили (Ebringer T., Sun L., Boztas S.), бинарные профили, последовательность строк, шестнадцатеричные дампы (Schultz M. G., Eskin E., Zadok F., Stolfo S. J.), n-граммы (Santos I., Brezo F., Ugarte-Pedrero X., Bringas P. G.) и многие другие.

Таким образом, существующие подходы к идентификации установленного программного обеспечения обладают рядом ограничений и недостаточными возможностями по обеспечению комплексной реализации мер по информационной безопасности.

Рассматриваемая в работе методика идентификации программного обеспечения направлена на распознавание программы, не основываясь на ее целостности. В результате создается достаточная гибкость, позволяющая определять исполняемые файлы, ранее не задействованные в процессе формирования эталонных сигнатур. Процесс идентификации сравнивает две сигнатуры: эталонную – созданную на основе обучающей выборки и сигнатуру

8
идентифицируемого исполняемого файла – создаваемую непосредственно перед этапом сравнения.

Все возрастающее количество версий программного обеспечения, обуславливает сложность полноценной комплексной защиты информации. В следствии чего, возникает необходимость в увеличении показателей качества идентификации программ. Противоречие, возникающее между недостаточной точностью идентификации ПО существующими методами и необходимостью по обеспечению достаточного уровня защищенности информации в информационной системе, обуславливает актуальность данного исследования.

Целью работы является увеличение точности идентификации установленного программного обеспечения за счет разработки и обоснования научно-методического аппарата по идентификации исполняемых файлов, устанавливаемых на средства вычислительной техники, обеспечивающего увеличение точности идентификации в условиях наличия различных версий, большого числа различных наименований программ и ограниченности числа объектов обучающей выборки.

Для достижения поставленной цели были решены следующие задачи **139** :

1. **134** Исследование и анализ существующих подходов и методов идентификации программного обеспечения, используемых отечественными и зарубежными исследователями.
2. Разработка модели нарушителя информационной безопасности организации.
3. Разработка метода выделения признакового пространства исполняемых файлов с последующим созданием сигнатур на его основе.
4. Разработка метода сравнения сигнатуры идентифицируемого исполняемого файла с эталонными сигнатурами программ, обеспечивающего увеличение точности идентификации.
5. Разработка методики идентификации исполняемых файлов на основе созданного архива эталонных сигнатур программ с использованием статистического подхода и машинного обучения.

6. Проведение вычислительного эксперимента и обоснование

применимости разработанной методики идентификации программного обеспечения.

В соответствии с заявленными целями и задачами работы: объектом

исследования является 7 разнообразие версий программного обеспечения, а предметом исследования методы идентификации программного обеспечения на основе статического анализа характеристик дизассемблированного кода программ.

В результате проведенных исследований были получены следующие результаты, составляющие научную новизну диссертации:

1. Метод формирования сигнатур исполняемых файлов, основанный на построении частотного распределения каждой из градаций выделенной характеристики исполняемых файлов, отличающийся от существующих использованием ряда отобранных наиболее информативных и устойчивых в проявлении ассемблерных команд.
2. Метод сравнения сигнатур идентифицируемых исполняемых файлов с ранее сформированными эталонными сигнатурами программ, отличающийся от известных применением комбинированного подхода использования алгоритма машинного обучения и аддитивного критерия, способствующего снижению числа ошибочных результатов классификации и обеспечивающего увеличение точности от совокупного использования признакового пространства, а также учитывающий ряд изменений в коде исполняемых файлов и позволяющий идентифицировать не рассматриваемые на этапе обучения версии программ.
3. Разработанная методика идентификации ПО, основанная на комбинированном анализе характеристик дизассемблированного кода программ, отличающаяся от известных, применением уникального сформированного признакового пространства и теории полезности для принятия решения на основе аддитивного критерия, что позволяет распознавать версии программ, ранее не задействованных в создании эталонных сигнатур исполняемых файлов

Теоретическую и практическую значимость работы 139 составляют разработанные методы и методика по идентификации программного обеспечения,

10 позволяющие обнаруживать нарушения информационной безопасности при обработке конфиденциальной информации, вызванные несанкционированной установкой программ. Проведенные вычислительные эксперименты подтверждают результативность предложенных решений на реальных данных. Методология и методы исследования поставленных задач включали теорию информационной безопасности, методы математической статистики, теорию предпочтений, методы машинного обучения, экспериментальные методы исследования.

На защиту выносятся следующие положения:

1. Метод формирования эталонных сигнатур программ и сигнатур идентифицируемых исполняемых файлов, основанный на статическом подходе анализа характеристик дизассемблированных кодов программ, обеспечивающий создание уникальных по форме и амплитуде частотных распределений, построенных на использовании ряда отобранных наиболее информативных ассемблерных команд, устойчиво проявляющихся в различных программах.
2. Метод сравнения сигнатур идентифицируемых исполняемых файлов с ранее сформированными эталонными сигнатурами программ при помощи комбинированного подхода использования алгоритма машинного обучения – градиентного бустинга деревьев решений и аддитивного критерия Фишберна, позволяющий достигать наименьшего числа ошибок неверной классификации и максимизировать эмерджентность совокупности признакового пространства.
3. Методика идентификации исполняемых файлов, включающая разработанные методы по формированию и сравнению сигнатур идентифицируемых исполняемых файлов с эталонными сигнатурами программ из архива эталонных сигнатур программ, обеспечивающая увеличение точности

идентификации при комбинированном анализе характеристик из их дизассемблированного представления.

Обоснованность и достоверность полученных результатов **7** подтверждается

использованием апробированного математического аппарата **7**; проведением

сравнительного анализа с существующими методами; серией практических **7**

11

экспериментов по идентификации **7** исполняемых файлов; проверкой адекватности

положений и выводов; согласованностью результатов, полученных при

теоретическом исследовании с результатами проведенных экспериментов **7**;

практической апробацией результатов исследования в докладах и публикациях на

отечественных и зарубежных научных конференциях.

Результаты практической апробации работы подтверждают адекватность и

корректность разработанных методов. Основные результаты работы были

представлены на следующих конференциях: 18th, 20th Conference of Open

Innovations Association FRUCT and **9** ISPIT 2017 seminar, 2016, 2017гг.; IV, VI, VII,

VIII Всероссийский конгресс молодых ученых, 2015, 2017, 2018, 2019гг.; 11th IEEE

International Conference on Application of Information and Communication

Technologies, AICT **16** 2017, 2017гг.; IX, X Санкт-Петербургская межрегиональная

конференция «Информационная безопасность регионов России (ИБРР **7** -2017)»,

2015, 2017гг.; Региональная информатика "РИ-2016", 2016гг.; XLVII, XLVIII

Научная и учебно-методическая конференция Университета ИТМО, 2018, 2019гг.;

International Conference on Next Generation Wired/Wireless Networking **8** Conference on

Internet of Things and Smart Spaces NEW2AN **8** 2018, ruSMART, 2018гг.; 28-я научно-

техническая конференция. Методы и технические средства обеспечения

безопасности информации, МитСОБИ **8**, 2019г.

Результаты, полученные в диссертации, были реализованы в рамках

выполнения следующих НИР: Проект по программе Президиума РАН No 0073-

2018-0008 « Теория и распределенные алгоритмы самоорганизации группового

поведения агентов в автономной миссии **8** », 2018,2019гг.; Проект по программе

Президиума РАН No 0073-2018-0007 « Разработка масштабируемых устойчивых

алгоритмов построения семантических моделей больших данных и их

использование для решения прикладных задач кластеризации и машинного

обучения», **138** 2018,2019гг.; НИР-ФУНД «Разработка методов интеллектуального

управления киберфизическими системами с использованием квантовых

технологий» No617026 (2017-2018гг.); НИР-ФУНД «Разработка методов создания

и внедрения киберфизических систем» No 619296 (2018-2019гг.). Результаты

12

работы использовались при разработке системы мониторинга состояния

внутренних сетей компании АО «НПК «ТРИСТАН». Полученные результаты

используются в образовательном процессе факультета БИТ Университета ИТМО

по направлениям подготовки **139** бакалавриата 10.03.01 и магистратуры 10.04.01 по

дисциплинам «Организация и управление службой защиты информации», «Теория

вероятностей», «Методы цифровой обработки видеозображений», «Управление

информационной безопасностью».

Публикации. По результатам диссертационного исследования автором

опубликовано 32 работы, из них статей в журналах, рекомендованных ВАК РФ **7** – 8

входящих в базы цитирования Web of Science и Scopus – 8, свидетельств о

государственной регистрации программы для ЭВМ **139** – 6, в прочих изданиях – 10.

Личный вклад автора. Результаты диссертационной работы получены

автором самостоятельно. Автором проведен анализ существующих методов

идентификации **7** программного обеспечения. Проанализированы условия и

ограничения применения каждого из методов.

Структура и объем диссертации. Диссертационная работа содержит

введение, 4 раздела, заключение, список литературы и 6 приложений. Объем

работы составляет 163 страницы. Работа включает 26 рисунков, 17 таблиц. Список

литературы содержит 101 наименование.

Краткое содержание работы.

Во введении обоснован выбор темы настоящей работы и ее актуальность,

представлена степень разработанности темы, определены объект, предмет и цель

исследования; описаны частные задачи, обоснована теоретическая и практическая

значимость получаемых результатов; раскрыты принципы используемых подходов

и разработанной методики; сформулированы положения, выносимые на защиту и

проведена апробация результатов исследования.

В первой главе представлен **11** анализ существующей проблематики

современного подхода к идентификации программного обеспечения. На основе

выделения объекта защиты, цели реализации угроз, потенциальных уязвимостей и

видов ущерба разработаны модель угроз и нарушителя информационной

13

безопасности при обработке информации конфиденциального характера в

информационной системе.

Во второй главе произведен обзор современных решений в области

идентификации исполняемых файлов, представленных в отечественной и

зарубежной литературе **139**. Проанализированы различные подходы к сбору

характеристик файлов и методы сравнения наборов данных характеристик для

нескольких файлов. Выявлены достоинства и недостатки данных методов.

Сформулирована постановка задачи обеспечения безопасности АС от

потенциальных угроз со стороны нелегитимно установленного ПО.

В третьей главе произведен анализ структуры и характеристик ELF-файла,

его дизассемблирование и представление исходного кода на низкоуровневом языке

ассемблера. Описаны подход к представлению ПО и модель представления

исполняемого файла в настоящем исследовании. Произведено выделение

признакового пространства и описано его дальнейшее использование в различных

методах формирования сигнатур.

Представлено подробное описание разработанных методов формирования

сигнатур – эталонных и идентифицируемых. Сформулированы методы сравнения

сигнатур – основанные как на статистическом подходе, так и на машинном

обучении. Описана методика идентификации исполняемых файлов на основе

статического анализа характеристик дизассемблированного кода программ,

включающая в себя также и этап постобработки результатов сравнения сигнатур

для увеличения точности идентификации исполняемых файлов. Представлены

ограничения, накладываемые на методику и условия ее использования.

Сделан вывод о том, что **139** необходимо экспериментальное подтверждение

работоспособности представленной методики и входящих в нее методов, а также

выделение наиболее результативных из них, путем сравнения достигаемых итогов

точности идентификации. А также сравнение получаемых результатов с

результатами методов отечественных и зарубежных исследователей.

В четвертой главе с целью проверки точности идентификации исполняемых

файлов при помощи разработанной методики была проведена серия

14

экспериментов, направленных на каждый разработанный метод сравнения

сигнатур файлов в отдельности. Полученные результаты прошли сравнение с

результатами других исследователей и на данном основании **сделан вывод о том,**

что **139** в условиях **139** наличия множества классов и ограниченности объема обучающей

выборки, разработанная методика обеспечивает более высокую точность

идентификации.

На практике методика идентификации программного обеспечения может

быть применена для повышения безопасности информационной системы организации, путем внедрения ее в качестве технической меры по аудиту программного обеспечения на электронных носителях информации, она позволяет осуществлять автоматизированную идентификацию ELF-файлов в соответствии с формируемыми архивами легитимных или нелегитимных программ. Выделен ряд смежных задач, решаемых разработанной методикой.

В заключении приведены основные результаты и выводы по разработанной методике идентификации исполняемых файлов на основе статического анализа характеристик дизассемблированного кода программ.

15

Глава 1. Постановка задачи обеспечения безопасности автоматизированной системы от потенциальных угроз со стороны несанкционированно установленного программного обеспечения

Под идентификацией программного обеспечения понимается процедура построения некоторой информативной модели (модель в виде математического кортежа) программы (сигнатуры) по выбранному признаковому пространству (ассемблерным командам), характеристики которой позволяли бы, с заданной точностью, найти соответствие между рассматриваемой (идентифицируемой) программой и, предопределенной ранее на этапе формирования архива сигнатур, конкретной программой. Другими словами, идентифицировать исполняемый файл означает распознать его как ту или иную программу. Под программным обеспечением рассматриваются исполнимые и компоуемые файлы 32x и 64x разрядностей формата ELF [1] в Linux операционных системах для архитектур процессоров x86 и x86-64.

1.1 Обзор проблематики современных подходов к идентификации программного обеспечения

Стремительный прогресс информационных технологий, их повсеместное внедрение во все сферы человеческой деятельности [139] приводит к тому, что современная организация бизнес процессов становится полностью зависима от надлежащего функционирования информационных систем (ИС). Область, отвечающая за безопасность информации обрабатываемой в таких системах и ее ресурсов, является областью информационной безопасности (ИБ).

Понятие ИБ трактуется различными ГОСТами по-разному, может определяться как процесс: « защита конфиденциальности, целостности и доступности информации; кроме того, сюда могут быть отнесены и другие свойства, например аутентичность, подотчетность, неотказуемость и надежность» [2], [35] как свойство: « свойство информации сохранять [138

16 конфиденциальность, целостность и доступность» [138 3], как состояние: « состояние защищенности информации [данных], при котором обеспечены ее [их] конфиденциальность, доступность и целостность» [138 4]. [40] Однако, из всех определений вытекает главная цель информационной безопасности – обеспечение трех наиболее важных свойств защищаемой информации: конфиденциальности, целостности и доступности.

Основным документом, описывающим порядок и принципы обеспечения информационной безопасности в организации, является политика ИБ, определение которой вытекает из двух понятий ГОСТа Р ИСО/МЭК 27002-2012, а именно

информационной безопасности « [13] Защита конфиденциальности, целостности и доступности информации; кроме того, сюда могут быть отнесены и другие свойства, например, аутентичность, подотчетность, неотказуемость и надежность.» [2] и [35] политики «Общее намерение и направление, официально выраженное руководством». Таким образом, политика информационной безопасности является совокупностью документированных управленческих решений, направленных на защиту [63] необходимых свойств информации и

ассоциированных с ней ресурсов. 63

Дополнительно необходимо отметить, что представляет собой понятие конфиденциальной информации, ведь по законам Российской Федерации, защите [5] подлежит секретная [6] и конфиденциальная информация [7]. Ранее понятие конфиденциальной информации давалось в ныне утратившем силу ФЗ от 20 февраля 1995 г. № 24-ФЗ: « документированная информация, доступ к которой ограничивается в соответствии с законодательством Российской Федерации» [138 8], теперь же в соответствии с ФЗ от 139 27 июля 2006 г. № 149-ФЗ понятие конфиденциальной информации вытекает из двух определений:

конфиденциальность информации «обязательное для выполнения лицом, получившим доступ к определенной информации, требование не передавать такую информацию третьим лицам без согласия ее обладателя» и 40 информация «сведения (сообщения, данные) независимо от формы их представления» [5]. 22 Таким образом, конфиденциальная информация – это сведения, независимо от формы их 42

17

представления, которые не могут быть переданы лицом, получившим доступ к данным сведениям, третьим лицам без согласия их правообладателя 42 .

При взаимодействии пользователя с информационной системой организации, ее ресурсами и обрабатываемой внутри информацией наступает потребность в регулировании такого взаимодействия, реализуемое за счет политики ИБ, мер и средств по обеспечению безопасности [9]. В настоящей работе идентификация ПО рассматривается как техническая мера обеспечения безопасности, подкрепляющая собой организационную меру, часто выраженную в положении установленной политики безопасности организации о запрете на несанкционированное установление ПО.

С каждым годом все более доступным становится доступ к различным страницам, сайтам сети Интернет, ориентированным на информирование интернет-пользователей о существующих уязвимостях операционных систем, программного обеспечения и эксплуатации нетривиальных способов обхода установленных на рабочих местах систем защиты информации. Так же в свою очередь не малую роль играет то, что Linux ОС и большое количество программного обеспечения для них, представляют собой свободное программное обеспечение, пользователи которого, в частности, имеют права на его свободное использование и, что самое важное, модификацию. Все вышеописанное ведет к тому, что аудит электронных носителей информации на предмет выявления несанкционированно установленного ПО становится важной и все более актуальной задачей. Разработанная методика позволяет выявить нарушения установленной политики безопасности при обработке конфиденциальной информации.

Действия пользователей АС, направленные против установленной политики безопасности в организации, способны стать источником роста числа уязвимостей системы и повлиять на ее ИБ. Причиной этого являются возможные дефекты ПО, наличие не декларированных возможностей, нелегальное использование интеллектуальной собственности, а также использование специальных программ, направленных на преодоление установленной защиты либо противоправных

18

действий внутри сети Интранет или Интернет. Последнее особенно актуально в области преступлений, связанных с компьютерной информацией [10].

Существующие подходы к идентификации установленного программного обеспечения обладают рядом недостатков и ограничений, обеспечивающие низкую точность идентификации версий не вредоносного ПО. Возникает задача повышения точности идентификации ПО.

1.2 Модель угроз информационной безопасности при обработке информации конфиденциального характера в информационной системе

В соответствии с описанными выше определениями ИБ, в обобщенном виде формулирующие поддержание наиболее важных свойств защищаемой информации: конфиденциальности, целостности и доступности, приведем ниже описание объекта и целей реализации угроз относительно задачи данного исследования (рисунок 1).

По определению объектом защиты информации [13] является « 26 информация или носитель информации, или информационный процесс, которые необходимо защищать в соответствии с целью защиты информации» [138 4], 22 таким образом для данного исследования объектом защиты являются сведения конфиденциального характера, обрабатываемые в автоматизированной системе, для которых необходимо предотвратить утечку, несанкционированные и непреднамеренные воздействия. Так же в качестве объекта может выступать сама автоматизированная система, ее ресурсы, электронные носители информации, и программное обеспечение, представляющее собой потенциальный источник уязвимостей системы [11].

19

Рисунок 1 – Взаимосвязь факторов, угроз и появление новых уязвимостей при несанкционированной установке ПО

Целью реализации угроз является осуществление личных интересов пользователем АС, посредством преднамеренной установки несанкционированного ПО, способного привести к материальному или моральному ущербу организации.

В качестве потенциальных угроз, способных привести к нижеописанному ущербу, выступают следующие угрозы, возникающие при нарушении организационных мер установленной политики ИБ организации: негативное воздействие на защищаемую информацию, как конфиденциального характера, так и открытого типа; сбой в работе программного и аппаратного обеспечения, а следовательно, и всей ИС организации;

20

необоснованное повышение расходов ресурсов АС (загрузка процессора, захват оперативной памяти и памяти на внешних носителях 41); негативное влияние на работоспособность персонала; административное и уголовное преследование со стороны компетентных органов.

Преодоление организационных мер приводит к появлению уязвимостей в АС, а следовательно, повышению рисков ИБ организации. Несанкционированно установленное ПО может не только стать лазейкой для противоправных действий нарушителя ИБ, но также привести к критическим ошибкам системы, способным остановить отложенные бизнес-процессы, замедлить работу АС, повлиять на работоспособность персонала, путем переноса внимания работника от непосредственных обязанностей, деморализации коллективных отношений и т.д. Так использование нелегального ПО нарушает законодательство Российской Федерации 139 и влечет наложение административной (если стоимость контрафактных экземпляров программ ЭВМ составит менее 50'000 рублей) или уголовной ответственности (если стоимость контрафактных экземпляров программ ЭВМ составит более 50'000 рублей) за нарушение авторских и смежных прав 139 , которая предусматривается статьей 146 Уголовного кодекса Российской Федерации 139 (УК РФ). При этом ответственность за противоправное использование ПО, в основном, возлагается на руководителя организации, принявшего управленческое решение по его использованию или же физическое лицо, самовольно установившее такое ПО. Отдельно необходимо отметить, что преступления по статье 146 УК РФ относятся к преступлениям публичного обвинения и не требуют подачи заявления с потерпевшей стороны. Ответственность же варьируется от административного штрафа в размере 1'500-2'000 рублей для физического лица до лишения свободы

на срок до 6 лет **136** со штрафом в размере до **139** 500'000 рублей **61**. Так, согласно исследованию BSA | The Software Alliance [12], было отмечено, что использование нелегального ПО в 2017 году составило 62%. При этом, коммерческая стоимость нелегального ПО, составила 1,291 миллиарда долларов США. По

21

данным **отчета о видах наказания по наиболее тяжкому преступлению (без учета сложения)** за **69** 2017 год к уголовной ответственности по статье 146 УК РФ было привлечено 482 человека [13].

Причинами возникновения уязвимостей могут являться:

наличие дефектов/не декларированных возможностей установленного

ПО;

нарушение российского законодательства в сфере охраны прав

интеллектуальной собственности;

нарушение лицензионного соглашения/ отсутствие лицензии на ПО;

направленное воздействие нарушителем на преодоление установленной

защиты, либо свершение противоправных действий внутри сети Интранет или

Интернет.

В качестве потенциального ущерба от реализации угроз может быть

причинен:

моральный и материальный ущерб деловой репутации организации,

возникающий в результате выявления надзорными органами нарушений

использования чужой интеллектуальной собственности;

моральный, физический или материальный ущерб, связанный с

разглашением персональных данных отдельных лиц **43**, который может возникнуть в

результате использования ПО имеющего дефекты или направленного на

преодоление установленной защиты;

материальный ущерб от разглашения защищаемой информации, который

может возникнуть в результате возникновения уязвимостей в АС при

использовании ПО имеющего дефекты или направленного на преодоление

установленной защиты;

материальный ущерб от необходимости восстановления нарушенных

защищаемых информационных ресурсов **64**, доступ или целостность которых может

быть нарушена;

22

моральный и материальный ущерб от дезорганизации деятельности

организации, вызванный уменьшением работоспособности пользователей АС;

материальный и моральный ущерб от нарушения международных

отношений, который может возникнуть в связи с нанесенным ущербом деловой

репутации организации.

Источниками угроз являются пользователи АС, реализующие деятельность

по преднамеренной несанкционированной установке ПО с целью извлечения своей

выгоды. Умысел может носить разнообразный характер, как подстраивайте АС

«под себя», так и, намеренная дестабилизация организованной системы защиты и

добывание защищаемой информации.

1.3 Модель вероятного нарушителя информационной безопасности и оценка

его возможностей

По признаку принадлежности к АС организации все нарушители будут

являются внутренними, в силу того, что у них имеются некоторые права допуска в

помещения и доступа к ЭВМ.

Естественно, по отношению к АС организации нарушителем может являться

и внешнее лицо, но, если предполагать, что ИБ организации находится в

устойчивом состоянии до того, как произойдет инцидент с несанкционированной

установкой ПО, тогда можно считать, что для преодоления установленных средств

защиты информации (ЗИ) внешнему нарушителю необходимы вначале действия внутреннего, которые повлекут за собой появление ранее не рассматривавшихся уязвимостей.

Возможности внутреннего нарушителя ограничены действующим в организации, в первую очередь, комплексом действующих в пределах контролируемой зоны режимных и организационно-технических мер [61], направленных, в рамках настоящей тематики работы, на пресечение действий пользователей АС по несанкционированной установке ПО.

23

Однако в зависимости от многих факторов (например, количеству денежных средств, выделяемых на организацию ИБ; наличию и уровню подготовки специалистов по ИБ; размерам штата и ресурсов в организации; добросовестности и уровню подготовки персонала в сфере ИБ; и т.п.), организационно-технические меры могут быть не достаточными и в таком случае у внутреннего нарушителя появляется возможность преодоления установленной политики безопасности.

Исходя из особенностей функционирования различных типов организаций, обрабатывающих информацию, подлежащую защите, к внутренним нарушителям можно отнести следующие общие категории персонала [14]: лица, имеющие непосредственный доступ к АС, но не имеющие доступа к защищаемой информации, при этом действия пользователей, относящихся к данной категории способны привести к дестабилизации комплекса защиты информации и увеличению рисков ИБ в сегменте локальной сети, где они расположены (категория 1);

лица, имеющие непосредственный доступ к АС и имеющие доступ к защищаемой информации, а также удаленный доступ по локальным и распределенным информационным системам, при этом действия пользователей, относящихся к данной категории, в дополнение к вышеописанной первой категории, способны привести к нежелательному воздействию на свойства защищаемой информации (категория 2);

лица, имеющие непосредственный доступ к АС и имеющие полномочия администратора безопасности сегмента информационной системы организации, здесь действия пользователя имеют более весомый характер воздействия на состояние защищенности информации, в дополнение руководящие позиции попадают под ответственность о защите интеллектуальной собственности (категория 3);

лица, имеющие непосредственный доступ к АС и имеющие полномочия системного администратора или администратора безопасности информационной системы организации, здесь следует учитывать, что данная

24

категория пользователей обладает большими правами и возможностями, знаниями относительно технических средств обработки информации, а также используемом программным обеспечением (категория 4);

лица, занимающиеся разработкой, сопровождением и ремонтом технических средств в информационной системе организации, могут являться штатными/внештатными сотрудниками и производить несанкционированные действия по установке ПО (категория 5).

Стоит также отметить, что к защищаемой информации относятся не только сведения, попадающие под указ Президента РФ No188 [7], но и общедоступная информация, ограничение доступа к которой быть не должно [5], а также информация, дестабилизирующее воздействие на которую (уничтожение, модификация, искажение, копирование, блокировка), может привести к моральному или материальному ущербу организации и ее деловой репутации.

Предполагается, что пользователи всех категорий способны реализовывать угрозы ИБ, используя стандартное оборудование, расположенное в локальной сети организации, воздействие через которое способно повлиять на всестороннюю

комплексность ИБ.

Кроме того, предполагается, что лица категории 3 и 4 хорошо знакомы с протоколами и алгоритмами, задействованными в информационной системе организации, и могут располагать специализированным оборудованием и программным обеспечением, позволяющим беспрепятственно производить несанкционированную установку ПО.

Для категорий 1-4 справедлива угроза снижения работоспособности персонала.

Предполагается, что для крупных организаций лица категории 3 и 4 проходят тщательную проверку со стороны руководства организации и являются доверенными. К сожалению, для более малых организаций такой тщательной проверки на квалифицированность и добросовестность не производится, поэтому данные категории наравне с остальными могут относиться вероятным нарушителям.

25

Возможны следующие способы реализации угроз ИБ:

несанкционированный доступ к защищаемой информации за счет использования функциональных особенностей несанкционированно установленного ПО, либо наличия в нем не декларированных возможностей, либо уязвимостей старых/новых версий;

несанкционированный доступ к средствам обеспечения защиты информации;

негативное воздействие на программно-технические компоненты ИС; использование программ игрового, информационного, развлекательного характера или предназначенная для обмена разного вида сообщениями; отправка сведений о нелегальном использовании ПО компетентным органам или не прохождение их аудита.

Таким образом, по причине существующих недостатков современных средств по проверке программного обеспечения, возможностей нарушителя и возможных убытков, необходимо производить более тщательный и своевременный аудит электронных носителей информации на наличие несанкционированно установленного ПО.

Для противодействия описанным угрозам и снижения рисков ИБ, в рамках проведенного исследования были разработаны методы идентификации версий ПО.

Выводы по главе 1

- а) В первой главе поставлена задача обеспечения безопасности автоматизированной системы от потенциальных угроз со стороны несанкционированно установленного программного обеспечения и конкретизирован вид рассматриваемого программного обеспечения.
- б) Произведён анализ проблематики современных подходов к идентификации ПО, обозначена актуальность исследований в области информационной безопасности, описаны основные нормативные документы в области защиты информации и использования ПО.

26

- в) Разработана и представлена модель угроз информационной безопасности при обработке информации конфиденциального характера в информационной системе, где также приводится классификация причин возникновения уязвимостей в автоматизированной системе при несанкционированной установке ПО и потенциальный ущерб организации.
- г) Разработана и представлена модель вероятного нарушителя информационной безопасности и оценка его возможностей, позволяющая сделать вывод о том 139, нарушителем может быть пользователь как первой категории, так и пятой.

27

обеспечения

В задаче идентификации объект $o \in O$ (где O – множество объектов)

описывается при помощи некоторого выделенного признакового пространства

f_1, f_2, \dots, f_n , (где $i = 1, 2, \dots, n$ – число градаций признака и значение для

$f_n \in D_f$ – множеству допустимых значений признака) способного охарактеризовать

основное свойство объекта $p(o)$, т.е. его класс, и записывается в виде вектора

значений признака $o = (f_1(o), f_2(o), \dots, f_n(o))$, именуемое в дальнейшем сигнатурой

программы. Предполагается, что существует функциональная связь между

признаками f_1, f_2, \dots, f_n и основным свойством объекта $p(o)$, позволяющая

определить значение свойства $p(o)$, в данном исследовании в качестве основного

свойства выступает имя программы, по ранее составленной информации о других

объектах o_1, o_2, \dots, o_n (называемых обучающей выборкой), и их свойств

$p(o_1), p(o_2), \dots, p(o_n)$.

Идентификация происходит путем применения выбранного алгоритма

сравнения двух векторов значений признака объектов: o – идентифицируемого и

o' – эталонного.

Различают следующие типы признаков [15]:

количественные признаки – признаки, измеренные в определенных

метрических (числовых) шкалах;

качественные – признаки, измеренные в не метрических шкалах и

используемые для выражения терминов и понятий, не имеющих цифровых

значений [67].

Также в зависимости от множества допустимых значений признака,

выделяют следующие типы [16]:

бинарный – принимающий значение нуля или единицы, $D_f = \{0, 1\}$;

номинальный – где множество допустимых значений признака конечно,

$|D_f| = \text{const}$;

28

порядковый – где множество допустимых значений признака конечно и

упорядочено, $|D_f| = \text{const}$ и $D_f = \{df \mid df_j \leq df_{j+1}, j = 1, \dots, \text{const}-1\}$;

количественный – где в качестве допустимых значений признака

выступает множество действительных чисел $D_f \in \mathbb{R}$.

В свою очередь представление признаков может носить разнообразный

характер, в частности являться перечислением признаков, векторным описанием

файла, или же иметь графическое представление, записанное в древовидной или

графовой модели. В настоящем исследовании применяются номинальные

(номинативные) признаки – это признаки, измеренные в не метрической шкале

наименований, представляющие собой набор числовых характеристик признака.

На сегодняшний день существует ряд подходов к идентификации

программного обеспечения на электронных носителях информации. Процесс

распознавания (или идентификации) программы можно разделить на два этапа:

первый - выделение признакового пространства и второй - применяемый метод

распознавания программы.

Подходы к распознаванию программ могут основаны как на структурной

характеристике файла, описывающей его содержимое, так и на семантической

составляющей, характеризующей смысл анализируемого кода программы. К

последним относится интеллектуальная обработка данных в которую входят [17]:

формальные методы – оперирующие фиксированной моделью объекта и

использующие predetermined алгоритмы обработки;

эвристические методы – способные, в некоторой степени, воспроизвести

процесс познания человека; они направлены на выявление неизвестных объектов.

Методы интеллектуальной обработки данных, как правило, включают

проверку «корректности» исполняемых файлов, вычисление статистического

распределения инструкций процессора или энтропии кода, а также выявление последовательностей инструкций, характерных **21** для вредоносных программ **77**. В связи с чем, использование данной группы методов в целях идентификации не вредоносного программного обеспечения не представляется возможным.

29

Ниже в таблице 1 Таблица 1 – представлены наиболее распространенные подходы в задаче идентификации исполняемого файла.

Таблица 1 – Подходы, используемые для идентификации программ

Способы сбора характеристик программы Методы распознавания программ

По статическим характеристиками

рассматриваемого объекта

По динамическим характеристиками

рассматриваемого объекта

Задействующие встроенные функции

операционной системы или

функциональность внедряемого

программного агента

Задание метрики схожести

Статистический анализ

Машинное обучение

Эвристический анализ

2.1 Методы идентификации, основанные на статическом характере анализа

характеристик программного обеспечения

Статические методы [18,19] оперируют двоичным образом программы,

хранящемся на электронном носителе информации, будь то внешний жесткий диск

или оперативная память компьютера. Данный метод подразумевает под собой

работу с исполняемым файлом или его дампом памяти как с массивом байтовой

последовательности **21**.

Статические методы, по способу сбора характеристик программы, являются

наиболее распространенными **82**. Данный класс методов подразделяется на

несигнатурные (методы проверки целостности): сравнение с полной копией

данных, сравнение контрольных сумм, контроль CRC (Cyclic Redundancy Check),

хеширование, имитовставка, цифровая подпись, и сигнатурные (методы сравнения

признаковых последовательностей): выделение «магического числа», сравнение

характерных последовательностей байтов и др.

30

В средствах защиты информации активно используется сигнатурный

(синтаксический) метод, заключающийся в сравнении байтовых

последовательностей **82** с эталонной сигнатурой.

Идентифицировать файл можно с разной степени конкретности, так самой

простой идентификацией будет являться определение расширения файла по его

сигнатуре, известной также как «магическое число», представляемое в виде строго

определенного набора байт, как правило, внесенных в начало тела файла. Так,

например, для ELF файла эта последовательность задана следующим набором

байт – 7F 45 4C 46; а для EXE файлов – 4D 5A.

Несигнатурные методы также популярны в настоящее время, но сфера их

применения ограничена и ниже будут описаны их достоинства и недостатки

применительно к задаче исследования.

Одним из первых и наиболее простых методов обеспечения целостности

файла разработка метода контрольной суммы. Первоначально под контрольной

суммой понималось некоторое число, занимающее позицию младшего разряда

байта, и вычисляемое при помощи сложения остальных 7 бит. Сегодня под данным

определением находится огромное число разнообразных алгоритмов, некоторые из

которых имеют широкое применение, а некоторые используются только для

специфических задач. Самыми известными и популярными алгоритмами являются: SHA-1, SHA-2, SHA-3, CRC и MD5, при этом стоит отметить, что [139] последние два алгоритма применяются все реже, т.к. не способны обеспечить достаточного уровня стойкости к коллизиям.

Хеширование является простейшим подходом к сравнению файлов [20], [21].

При вычислении хеш-функции в качестве исходных данных принимается исходная двоичная последовательность файла, далее в зависимости от выбранного алгоритма она разбивается на блоки фиксированной длины с последующим вычислением значения хеш-функции.

Достоинством хеш-функции является ее простота, скорость выполнения операции и размер получаемого значения хеш-функции, который всегда имеет одну и ту же длину. Однако, существенным недостатком в задаче сравнения двух

31

файлов является то, что вычисляемое значение хеш-функции существенно зависит от модификации произведенной над файлом, при этом оно ни в коей мере не зависит от степени произведенного вмешательства над оригинальным файлом. Т.е. если сравнивать значение хеш-функции оригинального файла со значением хеш-функции оригинального файла, но с измененным хотя бы одним битом – эти значения будут совершенно различны.

Чтобы снизить влияние модификаций и сравнивать похожие файлы, был предложен подход с использованием контекстно-зависимого кусочного хеширования (Context Triggered Piecewise Hashing [30]), который использует плавающее окно с фиксированным значением шага, таким образом, контекстно-зависимое кусочное хеширование представляет собой конкатенацию значений хеш-функций блоков файла. Данный подход использовался в работах [22], [23], [24]. Модификации алгоритма выбора длины окна, шага и границ образуемых блоков приводят к существенному улучшению при идентификации схожих файлов, однако остается ограничение на количество возможных вносимых изменений в оригинальный файл, показанное в приведенных статьях.

Цифровая подпись в исполняемых файлах используется для верификации программы ее автором, позволяя идентифицировать издателя файла (подлинность файла) и определить, не подвергался ли он изменениям (целостность файла), тем самым гарантируется безопасность программы для АС. Таким образом, если файл содержит некорректную цифровую подпись (или ее нет совсем), то это может означать, что данный файл опубликован ненадежным издателем или был изменен [25], однако существуют способы обхода проверки цифровой подписи и запуска недостоверного программного обеспечения [26].

Решение, представленное в статьях [27], [28] основано на предположении, что разные части одного и того же вредоносного обеспечения имеют схожий порядок кода и областей данных. Каждая такая область файла [95] может характеризоваться не только ее длиной, но также и своей схожестью. Другими словами, файл может быть охарактеризован комплексностью его областей кода.

32

Данный подход состоит в использовании вейвлет-анализа для сегментирования файлов [29] с различными уровнями энтропии и последующим использованием редакционного расстояния между сегментами последовательности для определения степени сходства файлов. Предлагаемое решение имеет ряд преимуществ, помогающих эффективно обнаруживать вредоносные программы на персональных компьютерах. Данный подход не принимает во внимание функциональность анализируемых файлов и основан исключительно на определении сходства в коде программы и данных о позиции области, что делает алгоритм эффективным против множества способов защиты исполняемого кода, с другой стороны, такое сравнение приводит к ложным срабатываниям. Стоит

отметить, что авторы статей рассматривали своей целью распознавание вредоносного ПО, в частности выявление упакованных программ, ориентируясь на характерные признаки их семейства.

В работах [30], [31], также описывается метод идентификации типа файла на основе энтропийного подхода построения статистического профиля файла.

Однако, авторы производят разбиения бинарного кода программы на блоки переменной длины при помощи метода скользящего окна и с учетом границ выявленных по статистическим свойствам. Выбор длины блока основывается на правиле о том, что каждый блок должен представлять собой статистически однородный набор данных произвольного размера [37]. Степень подобия двух файлов вычисляется с помощью описанного в работе показателя степени подобия файлов и задействует редакционное расстояние пары блоков статистического профиля файлов.

Еще один подход по выделению признаков пространства был предложен в работе [32], где авторы представляют исходный файл в виде последовательности цифровых изображений с фиксированной шириной и высотой, определяемой на этапе выбора длины скользящего окна для алгоритма подготовки данных. Также авторами предлагается второй подход к формированию входных данных, основанный на энтропии [17] различных участков кода исполняемого файла. Далее яркость каждого пикселя получаемого изображения принимается равной значению

33

конкретного байта из идентифицируемого объекта. Таким образом, исходный файл представляется в виде набора изображений с фиксированной шириной и высотой.

По причине, возникающей нежелательной, но сильной зависимости получаемых изображений от длины скользящего окна и наличию смещений в идентифицируемых файлах относительно друг друга, авторы обращаются к такому методу классификации объектов как машинное обучение, а именно к неоконгитрону – иерархической многослойной нейронной сети, предназначенной для инвариантного распознавания образов [24]. Однако сами авторы говорят о существенных ограничениях и недостаточной эффективности такого подхода для целей распознавания всех возможных семейств вредоносных файлов.

В работе [33] автор рассматривает идентификацию программ по внутренним характеристикам, при этом выделяет важность использования такого критерия установления степени подобия исходной и исследуемой программы, который был бы независимым от маскировок, вносимых в исходный текст программ нарушителем. Так, ссылаясь на работу [34], автор принимает решение об использовании операторов по тексту программ в качестве используемого признакового пространства. Повторяемость встречаемых операторов в коде программы делает степень подобия программ визуально видимой при построении их гистограмм, что особенно видно для программ имеющих схожее функциональное назначение, поскольку выбор операторов часто определяется целевой направленностью программы. В качестве меры оценки подобия в работе рассматривается корреляционная функция и поиск ее максимума между целевой и оцениваемой программами. Так наиболее схожие файлы будут иметь значение коэффициента корреляции близкое к единице, в свою очередь наименее схожие программы будут иметь значение, приближенное к нулю. Недостатком описанного подхода является то, что частота наиболее популярных операторов программы вносит существенный вклад в значение коэффициента корреляции и увеличивает его, что требует постоянного мониторинга и внесения изменений в пороговое значение коэффициента. Данный подход рассматривался автором в качестве меры по обеспечению безопасности программ, в случае, когда их исходные тексты

34

попадают в руки злоумышленников, желающих внести в код программы разрушающие программные средства до того, как программы подвергнутся компиляции.

В работе [35] для извлечения признаков характеристик программ использовались три способа: составление бинарных профилей, последовательность строк и шестнадцатеричные дампы. Автор применял описанные методы к исполняемым файлам ОС Windows и, в качестве методов распознавания программ, применял алгоритмы машинного обучения, такой как наивный классификатор Байеса. Первый способ извлечения признаковой характеристики заключался в выделении трех типов информации о ресурсе из исполняемых файлов PE формата: списка динамически подключаемых библиотек DLL; функции, вызываемые этими библиотеками; числом отличных системных вызовов каждой из библиотек – и дальнейшем формировании вектора бинарных характеристик. В качестве метода по распознаванию программ выступал алгоритм RIPPER [36]. Второй способ извлечения признаков заключался в использовании UNIX команды strings, которая показывает печатаемые строки в объекте или двоичном файле. Авторы формировали обучающую выборку, рассматривая строки как двоичные атрибуты, которые либо присутствовали в данном исполняемом файле, либо отсутствовали. Третий способ использовал утилиту hexdump, которая представляет исполняемый файл в виде набора шестнадцатеричных чисел. Также, как и для второго способа, автор использовал двухбайтовые последовательности в качестве бинарных атрибутов, которые либо присутствовали в данном исполняемом файле, либо отсутствовали.

В работе [37] авторы рассматривают подход к идентификации вредоносных исполняемых файлов PE формата на основе применения алгоритмов машинного обучения к формируемым n-граммам, представляющим собой последовательность из n элементов. Они используют утилиту hexdump для конвертации всех рассматриваемых исполняемых файлов. Далее формируют n-граммы, состоящие из четырехбайтовых последовательностей. Среди алгоритмов по классификации исполняемых файлов были использованы: метод k-ближайших соседей, наивный

35

Байесовский классификатор, метод опорных векторов (SVM), дерева решений. Также применительно к последним трем представленным алгоритмам использовался бустинг. Авторы решали задачу классификации с двумя классами, первый из которых являлся «исполняемый файл распознан как вредоносный», второй – «исполняемый файл распознан как безвредный». Авторами были достигнуты хорошие результаты, представленные в виде ROC и AUC, достигающие 99,6% покрытия площади под ROC (AUC) для доверительного интервала 0,05. Также авторы произвели классификацию вредоносных исполняемых файлов используя обобщенные три класса, характеризующих ключевые особенности вредоносного ПО: массовая почтовая рассылка (AUC = 0,8986), backdoor (AUC = 0,8704), вирус (AUC = 0,9114). К сожалению, не представляется возможности сравнить предоставленные результаты, так как результаты настоящего исследования используют мульти-классы.

В работе [38] представлено несколько методов дизассемблирования двоичных исполняемых файлов, помогающих сформировать представление исполнения двоичных исполняемых файлов и улучшить разбиение программы на «идиомы» (так называемые последовательности инструкций). Также стоит упомянуть здесь работу [39], в которой предложен метод, основанный на анализе CFG для обработки вредоносных программ с обфускацией. Позже, в [40] был улучшен метод путем включения семантических шаблонов вредоносных характеристик.

В работе [41] исследовалась способность ассемблерных команд обнаруживать вредоносные программы. Его исследование показало, что ассемблерные команды показывают значительные различия между вредоносным и легитимным ПО. А также то, что редко встречаемые команды несут большую долю предсказания, чем часто встречаемые команды.

В [42] изучен подход к идентификации ранее не известного вредоносного ПО

с помощью машинного обучения. Исполняемые файлы рассматривались как непрерывная последовательность ассемблерных команд. Поскольку большинство распространенных ассемблерных команд, которые могут использоваться для

36

вредоносных целей, требуют более одного машинного кода операции, было предложено использовать последовательности ассемблерных команд вместо отдельных команд. Таким образом, происходит добавление синтаксической информации, что, по мнению авторов, улучшает способность разработанного ими подхода к идентификации блоков последовательности команд, описывающих вредоносное поведение исполняемых файлов. В качестве идентификатора программы выступало чередование групп (вектор) из двух значений: частота встречаемости последовательности из n -грамм ассемблерных команд в файле и присвоенный им вес. В эксперименте принимали участие последовательности ассемблерных команд размерностей $n = 1$ и $n = 2$.

Схожий подход рассматривался в работах [43–46], где исследовалась точность классификации исполняемых файлов при изменении размера n -граммы; различном сборе характеристик программы – одновременное использование статических и динамических характеристик; в сфере IoT устройств; определенных семейств вредоносного ПО и т.д.

В работе [47] представляется исследование по классификации программ на основе алгоритмов, по выявлению плагиата. Обнаружение плагиата – это процесс выявления того, что некоторые части исследуемой работы не являются оригинальными наработками автора произведения. Одним из наиболее распространенных применений таких программ является обнаружение плагиата в программах, разработанными студентами. В данной работе считают, что обнаружение плагиата программного обеспечения и кластеризация вредоносных программ связаны друг с другом в том, что они оба пытаются обнаружить некоторую степень сходства между двумя программами среди большого числа исследуемых программ. Однако, из-за уникальности образцов вредоносных программ по сравнению с программами в целом (например, при использовании привилегированных системных ресурсов) и от степени обфускации, обычно применяемой в отношении вредоносных программ [17] в работе был выражен скептический настрой на эффективность рассматриваемых подходов.

37

Работа [48] рассматривает в качестве характеристики файлов – строковый константы, получаемые путем сканирования файла, чтения каждого знака один за другим и созданием списка всех получаемых строк. Длина строки определялась как набор символов. Чем короче длина строки, тем больше вероятность, что характеристика будет иметь усредненную значимость в исследуемой выборке файлов. Однако такая длина предоставляет больше возможностей. В работе отмечается, что для улучшения производительности и избегания переобучения при использовании методов машинного обучения, возникает необходимость в создании фильтра характеристик. В качестве метода идентификации были исследованы три подхода на основе классификатора Байеса: наивный (Naive Bayes), мульти-классификатор (Multi-naive Bayes) и частично-приращенный (Half Increment Bayes). Наилучший результат был достигнут при помощи использования последнего классификатора, однако временная сложность оценивается авторами как $T O (S_2 k_1 S_3 (k_1 k_2) \dots S_n (k_1 k_2 \dots k_g))$, где k_g – характеристики строковых констант, извлеченных из g -ой выборки, S_n – строковые константы, извлеченные из n -ой выборки.

Характеристикой исполняемого файла в работе [49] служит статистический профиль – вектор статистических описаний файла, полученный при движении скользящего окна от начала к концу файла с шагом в один байт, и расчета в каждом положении окна некоторой статистической величины. Ширина окна обычно

подбирается экспериментально. Затем в исследовании применяется метод разбиения файла на блоки фиксированной длины и подсчета их длин, полученных кодированием набора байтов при помощи алгоритма Хаффмана [37]. Также авторами рассматриваются графики энтропий для программы, упакованной с помощью различных упаковщиков.

2.2 Методы идентификации, основанные на динамическом характере сбора характеристик программного обеспечения

Динамический анализ основан на запуске анализируемой программы на исполнение. Отсутствие каких-либо предположений о ходе исполнения программы [48]

38 и проверка её в процессе или сразу после исполнения является значимым преимуществом динамического анализа. Очевидным требованием, предъявляемым при реализации динамического анализа – его проведение должно наименьшим образом влиять на ход исполнения. При определенных условиях на детерминированность программы, динамический анализ позволяет избежать проблемы ложных срабатываний [48].

Большинство современных вредоносных программ использует методы запутывания (обфускации), в частности бинарные упаковщики, чтобы помешать статическому анализу двоичных файлов. Таким образом, динамический анализ таких вредоносных программ зачастую бывает намного эффективнее статического анализа. Мониторинг поведения программ в ходе их исполнения включает, например, сбор выполняемых программой бинарных операций и предлагает потенциально более глубокое понимание самого кода. Пока подходы такого типа имеют свои ограничения, например – может быть трудно вызвать определенное поведение программ, а некоторые из них могут требовать определенные условия окружающей среды – тем не менее его использование считается более эффективным, по сравнению с только лишь статическими подходами. По этой причине, динамическому анализу программ уделяется большее внимание в научном сообществе. Системы анализа, как CWSandbox [50], Anubis [51], BitBlaze [52], AVG [53] и Advanced Threat Protection [54] выполняют запуск образцов вредоносных программ в защищенной среде [55] и отслеживают их поведение для анализа и разработки защитных механизмов. [17]

Так основными считаются подходы по [56]:

выявлению последовательности необычных путей исполнения кода программы или вызываемым системным функциям (API, system calls, Windows native API) [57];

нанесенным повреждениям окружающей среды в изолированной системе [58];

попыткам получения высоких привилегий [59];

39 изменениям параметров вызываемых функций [35];
способам обработки данных программой [60];
трассировке машинных инструкций программы [61].

Данный подход является более эффективным в плане анализа программы, однако для задачи идентификации большого количества версий не вредоносных программ его применение становится не целесообразным, поскольку главный минус динамического анализа связан с большими временными затратами, требуемыми на выполнение и тщательное исследование программы.

2.3 Методы идентификации, основанные на сборе информации путем обращения к встроенным функциям операционной системы или задействования собственного программного агента

В настоящее время множество компаний обращаются к использованию систем управления ИТ-активами (ITAM), представляющие собой комплексные решения, нацеленные на физический учёт, финансовый контроль и соблюдение контрактных обязательств, связанных с ИТ-активами, на протяжении всего их

жизненного цикла [62,63]. Здесь под ИТ-активами подразумеваются все аппаратные и программные элементы ИТ-инфраструктуры, обеспечивающие деятельность бизнес-среды.

В свою очередь ИТАМ подразделяется на управление аппаратными активами (НАМ), охватывающее управление материальными составляющими ИТ-инфраструктуры: пользовательские компьютеры, сервера, телефоны и т.д.; и на управление программными активами (САМ), охватывающее управление не материальными составляющими ИТ-инфраструктуры: программное обеспечение, лицензии, версии, конечные точки инсталляции и т.д.

На рынке услуг представлено немало решений [64], позволяющих идентифицировать программные активы, управлять учетом, а также производить контроль их изменений и др. Из числа наиболее известных программных продуктов можно выделить:

40

Работающие на основе агента:

- 1) Samanage – поддерживает работу с Windows, Linux, Mac OS [65].
- 2) Kaspersky Systems Management – поддерживает работу с Windows, Linux, Android, iOS [66].

Работающие без агента:

- 1) Microsoft Assessment and Planning Toolkit – поддерживает работу с Windows ОС [67].

Работающие как с агентом, так и без него:

- 1) AIDA64 – поддерживает работу с Windows, Linux, Android [68].
- 2) Lansweeper – поддерживает работу с Windows, Linux, Mac OS и устройствами, поддерживающими IP адресацию [69].

Агентом является программный агент, устанавливаемый на компьютер пользователя АС и ведущий скрытую деятельность по выполнению своих функций, в частности инвентаризации ПО. Недостатком данного подхода является необходимость установки агента на каждый исследуемый компьютер, однако преимуществом агента является возможность его установки на корпоративные ноутбуки и смартфоны, используемые пользователями за пределами внутренней сети организации, получают самый высокий уровень детализации, могут быть использованы для вмешательства в компьютер пользователя. Использование подхода без установки агента не требует наличия подходящих для данного агента систем, его обновлений, не повышает производительные затраты ресурсов компьютера, а также имеет удобное централизованное управление. Однако, уровень детализации получаемой информации ниже, чем при задействовании агента, а также необходимо удостовериться, что исследуемые активы доступны для их обнаружения.

Подход к оценке ресурсов с(без) агентом имеет свои преимущества и недостатки [139], особенно стоит учитывать число самих ресурсов в организации.

41

Однако стоит отметить, что большая часть предоставляемой информации собирается посредством встроенных технологий инвентаризации программного и аппаратного окружений, операционной системы:

Windows Management Instrumentation – является одним из базовых инструментариев централизованного управления и слежения за работой различных частей компьютерной инфраструктуры под управлением платформы Windows;

Active Directory Domain Services – служба каталогов Windows, которая централизованно сохраняет все данные и настройки среды в базе данных;

SMS Provider – поставщик инструментария управления Windows, назначающий доступ на чтение и запись к базе данных Configuration Manager на сайте сервера;

lshw – утилита Linux, которая выводит полный список аппаратных

компонентов системы вместе с информацией об устройствах;
dpkg -l – утилита Linux, которая выводит полный список программных
компонентов системы;
и других технологий.

Недостатки такого подхода очевидны. Отсутствие должного уровня
оценивания роли информационной безопасности в бизнес процессах организации
приводит к недооценке кадровой политики со стороны руководства при подборе
ИТ персонала. В то же время, пользователи АС имеют возможность достичь
достаточного уровня компетенции в сфере компьютерных технологий,
позволяющие им обходить установленные политики безопасности. Все это
предоставляет возможность внести изменения в конфигурационные данные
устанавливаемого программного обеспечения.

Ниже на рисунках 2-4 приведены примеры по изменению версии программ
(ABBYY Lingvo и Htop) в операционных системах Windows (10) и Linux (Ubuntu)
соответственно. Из которых становятся очевидны пути обмана средств аудита
программного обеспечения как для установленного ранее программного
обеспечения, так и для устанавливаемого впервые.

42

Рисунок 2 – Изменение версии программы ABBYY Lingvo x5 15.0.511
на ABBYY Lingvo x5 15.0.510

Рисунок 3 – Изменение версии программы ABBYY Lingvo x5 15.0.511
на ABBYY Lingvo x7 15.0.512

Из рисунков 2 и 3 видно, что путем произведенных манипуляций в реестре
операционной системы Windows, версия программы ABBYY Lingvo была изменена
дважды, целью данного вмешательства являлось намеренное влияние на
результаты выдачи запроса по инвентаризации программного обеспечения с
помощью средств Windows Management Instrumentation (WMI) (результат выдачи
представлен на рисунке 2) и распространенного программного средства Lansweeper
(результат выдачи представлен на рисунке 3).

На рисунке 4 в свою очередь представлена установка программы Htop в
операционной системе Ubuntu с заведомо измененной версией, исправленной
путем вмешательства в конфигурационный файл устанавливаемого пакета .deb.

43

Рисунок 4 – Установка измененной версии программы htop 8.1.0-3
вместо htop 2.1.0-3

Сравнение наиболее близких к задаче исследования работ по идентификации
исполняемых файлов представлено ниже в таблице 2.

44

Таблица 2 – Сравнение основных работ по идентификации исполняемых файлов

Источник
Формат
файлов
Используемые
характеристики файла
Метод
идентификации
Число классов
Оценка качества метода
1
Точность
(Ассигасу),
%
Бикубическая
мера качества

кластеризации

1 2 3 4 5 6 7

Data mining methods for

detection of new

malicious

executables 28 , [35]

PE/?

2

Шестнадцатеричный

дамп

Signature

Бинарная

классификация

49,31 -

Используемые DLL RIPPER 83,61 -

функции, вызываемые

DLL

RIPPER 89,36 -

Число отличных

системных вызовов

каждой из DLL

RIPPER 89,07 -

Печатаемые строки Naive Bayes 97,11 -

Набор

шестнадцатеричных

чисел

Voting Naive Bayes 96,88 -

Opcode sequences as

representation of

executables for data-

mining-based unknown

malware detection 51 , [42]

PE

Последовательность

частоты встречаемости

очередности (n = 1)

ассемблерных команд

SVM: Pearson VII 92,92 -

Последовательность

частоты встречаемости

очередности (n = 2)

ассемблерных команд

SVM: normalised

polynomial

95,90 -

Последовательность

частоты встречаемости

очередности (n = 1 и

n = 2) ассемблерных

команд

95,80 -

1 Показана максимально достигаемая величина для используемом способе сбора характеристик файла и методе идентификации

2 Только часть файлов формата PE, для остальной части файлов формат не известен.

Продолжение таблицы 2

1 2 3 4 5 6 7

On Challenges in

Evaluating Malware

Clustering, [47]

?

Последовательность

частоты встречаемости

очередности вызовов

API (n = 3)

Мера Жаккара +

кластеризация на

основе работы [70]

Кластеризация

- 0,81

Последовательность

вызовов API

Сравнение строк

программными

средствами (diff,

CCFinder) +

кластеризация на

основе работы [70]

- 0,78

Identifying almost

identical files using

context triggered

piecewise hashing **28**, [22]

PE

Побайтовый код

программы

Контекстно-

зависимое

кусочное

хеширование

Мульти-

классификация

- 0,29

A Fast Randomness Test

that Preserves Local

Detail, [49]

Вектора для блоков

постоянного размера

побайтового кода

программы

Евклидово

расстояние

- 0,55

Идентификация типа

файла на основе

структурного

анализа, [30]

?

Вектора для блоков

постоянного размера

побайтового кода

программы

Редакционное

расстояние

- 0,69

Identification of

Executable Files

on the basis of Statistical

Criteria, [71]

ELF

Последовательность

частот встречаемости

118 ассемблерных

команд

Статистический

критерий

однородности хи-

квадрат

98,67/

79,55

3

-

Последовательность

частоты встречаемости

одной ассемблерной

команды на неравных

интервалах

Статистический

критерий согласия

Колмогорова

95,09/

50,41

3

-

3 Оценивание точности статистического критерия происходило двумя способами (Confusion matrix и Accuracy), описание которых представлено в разделе 2.6.

46

Продолжение таблицы 2

1 2 3 4 5 6 7

Подход к выбору

информативного

признака в задаче

идентификации

программного

обеспечения, [72]

ELF

Последовательность

частоты встречаемости

одной ассемблерной

команды на равных

интервалах

Статистический

критерий

однородности хи-

квадрат

Мульти-

классификация

98,58/ 78,05

3

-

Разработка способа

идентификации elf-

файлов на основе

нейронной сети, [73]

Нейронная сеть

MLP, с методом

сопряженных

градиентов и 50

нейронами

82,11 -

Алгоритм градиентного

бустинга деревьев

решений в задаче

идентификации

программного

обеспечения **8**, [74]

CatBoost 89,43 0,88

Сравнительный анализ

подходов к

идентификации

программного

обеспечения, [75]

LightGBM 86,99 0,84

XGBoost 95,93 0,96

Повышение точности

идентификации

программного

обеспечения путем

использования

аддитивного критерия

Фишберна, [76]

XGBoost и

Фишберн

99,19 0,99

47

В данной главе был произведен обзор существующих методов

идентификации программного обеспечения, из которого можно выделить ряд

ограничений, не позволяющий их использование в задаче идентификации версий

программного обеспечения:

Для статического сбора характеристик файла:

- 1) Методы основанные на оценке целостности файла не позволяют исследовать схожесть программ, не задействованных в обучающей выборке.
- 2) Использование особенностей PE формата программ ОС Windows, отсутствующие в ELF-файлах.
- 3) Большинство существующих подходов к классификации исполняемых файлов направлены на бинарное разделение тестовой выборки и выделение класса вредоносных программ.
- 4) Так же часть подходов рассматривают мульти-классификацию с ограниченным набором классов.
- 5) Работы же, направленные на мульти-классификацию с неограниченным набором классов, обладают низкой точностью идентификации исполняемых

файлов.

Для динамического сбора характеристик файла:

- 1) Использование данного типа подходов является не целесообразным, по причине задействования огромного числа ресурсов.
- 2) Для выполнения программы бывает необходимо создавать особую среду исполнения.

Для сбора информации путем обращения к встроенным функциям

операционной системы или задействования собственного программного агента:

- 1) Главным недостатком является сам принцип получения информации с АС, где факт наличия способов изменения информации о программе был приведен в соответствующем пункте 1.4.3.

Автором научной работы были проанализированы методы идентификации исполняемых файлов АС, использующие в качестве идентификационных

48 признаков последовательности системных вызовов, байтов, ассемблерных команд, энтропийных характеристик и др. Представлены их ограничения в возможностях идентификации программ.

Приведенный анализ методов идентификации различных версий исполняемых файлов, позволяет сделать вывод, что идентификация на основе статического анализа характеристик программы и использование статистических и машинных алгоритмов классификации, является наиболее перспективным направлением деятельности в рамках поставленной научной задачи исследования.

Таким образом, в современных условиях **139**, не существует решения, позволяющего с достаточной точностью производить идентификацию **7** различных версий ПО в АС в задаче по выявлению факта нарушения пользователем установленных правил по несанкционированной установке ПО. Отсюда вытекает необходимость в проведении настоящего исследования, направленного на достижение максимальной точности идентификации исполняемых файлов, путем совершенствования существующих подходов, поиска новых, ранее не применявшихся, решений в области выделения признакового пространства файлов, методов идентификации, а также способов постобработки результатов.

2.4 Постановка задачи исследования

В качестве идентифицирующей исполняемые файлы информацией – сигнатуры, выступают структурные характеристики программы, считываемые с помощью подхода, основанном на статическом характере сбора. Имеются сигнатуры идентифицируемого исполняемого файла (СИИФ) и эталонные сигнатуры программ (ЭСП).

Идентифицировать исполняемый файл при аудите электронных носителей информации означает распознать программу путем анализа ее структурных характеристик и сравнения сформированной СИИФ с ЭСП некоторых известных нам потенциальных программ, хранимых в архиве, или сделать вывод о том, что **139** программа не является ни одной из нам известных программ.

49

Определение того, является ли идентифицируемый файл той или иной программой (принадлежит определенному классу) из числа возможных потенциальных программ, производится на основании анализа значений выделенного признакового пространства СИИФ и ЭСП, хранимых в архиве сигнатур программ (АЭСП).

Архив может строиться на основе одного из принципов – в него заносятся сигнатуры программ, нахождение которых на электронном носителе информации разрешено; либо в него заносятся сигнатуры программ, нахождение которых на электронном носителе информации запрещено. Выбор одного из принципов зависит от поставленных целей и задач при аудите.

Исходя из описанного выше, задачу идентификации исполняемых файлов по

статическим характеристикам можно разделить на следующие две подзадачи:

1) Формирование АЭСП А путем сбора выбранной характеристики F (feature) различных версий эталонных программ Рэтал.,i (program), и создание на их основе ЭСП S(Рэтал.,i) (signature).

2) Идентификация неизвестного исполняемого файла путем сравнения его СИИФ S(ej) с ЭСП и последующей постобработкой результатов. В случае выдачи положительного результата применяемым методом идентификации, исполняемый файл ej (executable) считается идентифицированным как та или иная программа Рэтал.,i, или же как не файл, не относящийся ни к одной из программ АЭСП.

Задача идентификации ПО в общем случае сводится к решению задачи мульти-классификации.

Для каждого идентифицируемого файла необходимо определить его меру близости с потенциальной эталонной программой из АЭСП, т.е. необходимо отнести исполняемый файл к определенному классу известных нам программ. Имеется множество дизассемблированных представлений потенциальных эталонных программ 32-х и 64-х разрядностей формата ELF в Linux ОС для архитектур процессоров x86 и x86-64, взятых с официальных репозиториях и обработанных программой objdump P = {Рэтал.,1, Рэтал.,2,..., Рэтал.,i} и их имен N = {N1, N2,..., Ni}, при этом функция соответствия программы Рэтал.,i и имени

50

программы Ni – биективна, где i – число программ разных наименований формирующих АЭСП (обучающая выборка). Так же есть множество версий для программ Рэтал.,i = {v1, v2,..., vm}, где m – число версий программы с одним наименованием, которые участвуют в формировании сигнатур S(Рэтал.,i).

В свою очередь для формирования сигнатуры версии программы используется функция $q(v_m, F): X \rightarrow N_0$, которая при помощи заданного алгоритма и выбранной характеристики F формирует частотную последовательность признака версии программы, где X – множество принимаемых значений характеристики F. Затем полученные сигнатуры версий программы формируют наборы ЭСП S(Рэтал.,i) = {(Ni, q(v1,F)), (Ni, q(v2,F)),..., (Ni, q(vm,F))}, где Ni – имя программы Рэтал.,i, которые заносятся в АЭСП: A = {S(Рэтал.,1), S(Рэтал.,2),..., S(Рэтал.,i)}.

При проведении аудита имеется некоторое множество неизвестных исполняемых файлов E = {e1, e2,..., ej} (тестовая выборка), где j – число исследуемых исполняемых файлов. С помощью той же выбранной характеристики F и функции q производится формирование частотной последовательности признака для каждого исследуемого файла ej, СИИФ S(ej) представляет собой результат функции $q(e_j, F): X \rightarrow N_0$.

Требуется построить алгоритм I(S(Рэтал.,i), S(ej)): ej → Ni (identification) способный определить вероятность/ меру сходства идентифицируемого исполняемого файла ej с именем Ni эталонной программы Рэтал.,i, который бы удовлетворял следующему критерию: точность (метрика оценки результатов идентификации: Accuracy или F-measure) должна быть максимальной, при условии ограничений на:

потенциальное количество различных программ, установленных на одном компьютере или на всех компьютерах организации;
количество различных версий для каждой программы, измеряемое от одной до десятков, но в практике не превосходящее сотни;
наличие существенных изменений в различных версиях одной программы;

51

наличие индивидуального характера распределения признака идентифицируемого файла от других программ;
возможность производить дизассемблирование исполняемого файла;
не способность в выявлении закладки вредоносного кода в легитимный

исполняемый файл;

не способность классификации выдать корректный результат для программы, класс которой не был определен на этапе формирования модели классификации (нельзя создать класс «файл не похож ни на одну из программ»). Необходимо провести обучение данного алгоритма по обучающей выборке P и протестировать его на тестовой E . При проведении сравнения, в зависимости от выбранной меры схожести, идентифицируемый исполняемый файл будет считаться программой N_i , при достижении алгоритмом I порогового значения (статистические критерии) или же максимальной вероятности принадлежности к классу (алгоритмы машинного обучения).

При решении задачи идентификации исполняемых файлов на качество идентификации могут оказывать влияние следующие факторы:

используемая характеристика исполняемого файла;

метод сравнения СИИФ и ЭСП.

Таким образом, основными подзадачами идентификации исполняемых файлов, можно выделить следующие:

выделение информативных характеристик ELF-файлов, обладающих максимальной различающей способностью;

разработка метода идентификации исполняемого ELF-файла по его внутренним характеристикам, обеспечивающего максимальную точность идентификации, при заданных ограничениях.

Из приведенных подзадач можно выделить более частные задачи исследования:

52

разработка метода выделения признакового пространства программного обеспечения с последующим созданием сигнатур исполняемых файлов на его основе;

разработка метода сравнения СИИФ с ЭСП, обеспечивающий более

высокую точность идентификации;

разработка методики идентификации 7-й версии программного обеспечения на основе созданного архива сигнатур с использованием статистического подхода и машинного обучения;

проведение вычислительного эксперимента и обоснование применимости разработанного метода по идентификации программного обеспечения.

Выводы по главе 2

а) Проведенный анализ методов идентификации ПО и выявленные их ограничения и недостатки, позволяют сделать вывод о том, что 139 для задачи данного

исследования наиболее подходящим будет разработка подхода к идентификации на основе статического анализа характеристик программ и использование статистических и машинных алгоритмов классификации.

б) Анализ зарубежных и отечественных работ позволяет сделать вывод, что методы машинного обучения дают наиболее высокий показатель точности идентификации ПО.

в) В сфере идентификации ПО большая часть исследований направлена на задачу классификации вредоносного ПО, имеющего дополнительные отличительные черты функциональности, облегчающие задачу его анализа.

Однако перед настоящим исследованием стоит более трудная задача по идентификации не вредоносного ПО.

г) На основании предыдущих пунктов и проведенного анализа литературы можно утверждать, что на данный момент не существует 139 решения, позволяющего с максимальной точностью производить идентификацию не вредоносного ПО на основе его характеристик. Произведенный анализ позволил определить основные

53

задачи в области идентификации ПО и обозначил необходимость проведения исследования, направленного на максимизацию точности идентификации

исполняемых файлов, путем совершенствования существующих подходов, поиска новых, ранее не применявшихся, решений в области выделения признакового пространства файлов, методов идентификации, а также способов постобработки результатов.

54

Глава 3. Разработка методики идентификации установленного программного обеспечения на основе статических характеристик программного кода файла

Разработка методики идентификации исполняемых файлов включает в себя следующие задачи:

исследование исполняемых файлов в соответствии со стандартом представления формата ELF, проведение их структурного и семантического анализа с целью формирования представления о процессе исполнения программ, их расположения в АС;

изучение способов представления исполняемых файлов, выделение списка характеристик, пригодных для использования при формировании сигнатур исполняемых файлов, их анализ и выбор наиболее информативных;

формирование модели представления исполняемых файлов, исходя из результатов, полученных в предыдущих пунктах;

разработка метода формирования сигнатур исполняемых файлов, предоставляющий возможность идентификации на их основе программ, ранее не задействованных на этапе создания АЭСП. Необходимо сформировать признаковое пространство на основе выбранных характеристик исполняемых файлов, проанализировать существующие подходы к определению их информативности;

разработка метода сравнения нескольких сигнатур и предоставления идентификатора программы, основываясь на проведенном обзоре зарубежных и отечественных исследований. Также следует обозначить подход к оценке вычислительной сложности созданных методов;

последней решаемой задачей станет компоновка результатов описанных задач и описание методики идентификации программного обеспечения.

Описание этапов проведения идентификации, а также процесс оценки точности идентификации ПО.

55

3.1 Анализ структуры и характеристик исполняемого файла

Исполняемый файл или программа представляет собой комбинацию компьютерных инструкций и данных, позволяющая аппаратному обеспечению вычислительной системы выполнять вычисления или функции управления [77]. Данные исполняемого файла имеют свой формат (например: EXE, PE, ELF, COFF и т.д.) в зависимости от операционной системы, и состоят из нескольких частей, как правило заголовков, самого тела программы – кода, и остального.

В данной работе исследуются исполняемые ELF-файлы операционной системы Linux. Аббревиатура данного формата расшифровывается как Executable and Linkable Format и в переводе означает формат исполнимых и компоуемых файлов.

ELF-файлы включают в себя program linking (собираение программы) и program execution (запуск программы). Для удобства и эффективности данный формат файла предоставляет параллельное представление содержимого, отражающее различные потребности его активности. На рисунке 5 показана организация ELF-файла, где область слева (Linking View) отображает процесс собирания программы, а область справа (Execution View) – процесс ее исполнения.

ELF заголовок располагается в начале файла и содержит, так называемую, «дорожную карту», которая описывает структурную организацию файла. Также стоит отметить, что первые четыре байта отведены для «магического числа» – идентификатора, имеющего значение 7F 45 4C 46, где 7F – префикс, а остальные

расшифровываются в соответствии с ISO/IEC 8859-1 как 45 = E, 4С = L, 46 = F.

Секции содержат основную часть информации исполняемого файла для сборки программы: инструкции, данные, таблицы символов, информацию о перемещении и т.д. Таблица заголовков программы сообщает системе, как создать образ процесса. Файлы, используемые для создания образа процесса (исполнения программы), должны иметь таблицу заголовков программы; перемещаемые файлы в ней не нужны. Таблица заголовков секций содержит информацию, описывающую секции файла. Каждая секция имеет запись в таблице, которая

56
предоставляет такую информацию, как название секции, размер секции и т.д. Файлы, используемые при сборке программы, должны иметь таблицу заголовка раздела; остальные файлы могут не иметь ее.

Рисунок 5 – Представление ELF-файла

С помощью команды `readelf` можно посмотреть на структуру и состав ELF-файла. Ввод данной команды в окне консоли выдаст следующую информацию:

Команда для отображения ELF заголовка:

```
readelf -file-header htop
```

```
Magi: 7f 45 4c 46 01 01 00 00 00 00 00 00 00 00 00
```

```
Class:
```

```
Data:
```

```
Version:
```

```
OS/ABI:
```

```
ABI Version:
```

```
Type:
```

```
Machine:
```

```
ELF32
```

```
2's complement, little endian
```

```
1 (current)
```

```
UNIX - System V
```

```
0
```

```
EXEC (Executable file 20 )
```

```
Intel 80386 20
```

```
57
```

```
Version:
```

```
Entry point address:
```

```
Start of program headers 20 :
```

```
Start of section headers:
```

```
Flags:
```

```
Size of this header:
```

```
Size of program headers:
```

```
Number of program headers:
```

```
Size of section headers:
```

```
Number of section headers:
```

```
Section header string table index:
```

```
0x1
```

```
0x804ee80
```

```
52 (bytes into file)
```

```
122368 (bytes into file)
```

```
0x0
```

```
52 (bytes)
```

```
32 (bytes)
```

```
9
```

```
40 (bytes 46 )
```

```
28
```

Команда для отображения таблицы заголовков программы:

```
readelf --program-headers htop
```

Elf file type is EXEC (Executable file)

Entry point 0x804ee80

There are 9 program headers, starting at offset 52

Program Headers:

Type Offset VirtAddr PhysAddr FileSiz MemSiz Flg Align

PHDR 0x000034 0x08048034 0x08048034 0x00120 0x00120 R E 46 0x4

INTER

P

0x000154 0x08048154 0x08048154 0x00013 0x00013 86 R 0x1

LOAD 0x000000 0x08048000 0x08048000 0x1c17c 0x1c17c R E 0x1000

LOAD 0x01cef0 0x08065ef0 0x08065ef0 0x00e08 0x019f0 RW 0x1000

DYNAM

IC

0x01cefc 0x08065efc 0x08065efc 0x00100 0x00100 RW 0x4

NOTE 0x000168 0x08048168 0x08048168 0x00044 0x00044 R 0x4

GNU_E

H_FRA

ME

0x0190dc 0x080610dc 0x080610dc 0x0073c 0x0073c R 0x4

GNU_S

TACK

0x000000 0x00000000 0x00000000 0x00000 0x00000 RW 0x4

GNU_R

ELRO

0x01cef0 0x08065ef0 0x08065ef0 0x00110 0x00110 R 0x1

Section to Segment mapping:

Segment Sections...

00

01

02

03

04

05

.interp

.interp .note.ABI-tag.note.gnu.build-id.gnu.hash.dynsym.dynstr

.gnu.version.gnu.version_r.rel.dyn.rel.plt.init.plt.text

.fini.rodata.eh_frame_hdr.eh_frame

.init_array.fini_array.jcr.dynamic.got.got.plt.data.bss

.dynamic

.note.ABI-tag.note.gnu.build-id 20

58

06

07

08

.eh_frame_hdr 20

.init_array.fini_array.jcr.dynamic.got 20

Команда для отображения таблицы заголовков секций:

```
readelf --section-headers htop
```

There are 28 section headers, starting at offset 20 0x1de00:

Section Headers:

[Nr]

Name

Type

Addr

Off

Size

ES

Flg

Lk

Inf

AI 20

[0]

[1]

[2]

[3]

[4]

[5]

[6]

[7]

[8]

[9]

[10]

[11]

[12]

[13]

[14]

[15]

[16]

[17]

[18]

[19]

[20]

[21]

[22]

[23]

[24]

[25]

[26]

[27]

.interp

.note.ABI

-tag

.note.gnu

.build-i

.gnu.hash

.dynsym

.dynstr

.gnu.vers

ion

.gnu.vers

ion_r

.rel.dyn

.rel.plt

.init

.plt

.text
.fini
.rodata
.eh_frame
_hdr
.eh_frame
.init 20 _arr
ay
.fini_arr
ay
.jcr
.dynamic
.got
.got.plt
.data
.bss
.gnu_debu
glink
.shstrtab
NULL
PROGBITS
NOTE
NOTE
GNU_HASH
DYNSYM
STRTAB
VERSYM
VERNEED
REL

REL
PROGBITS
PROGBITS
PROGBITS
PROGBITS
PROGBITS
PROGBITS
PROGBITS 56
INIT_ARR
AY
FINI_ARR
AY
PROGBITS
DYNAMIC
PROGBITS
PROGBITS
PROGBITS
NOBITS
PROGBITS
STRTAB
00000000
08048154
08048168
08048188
080481ac
08048918

08049de8
0804af82
0804b21c
0804b2bc
0804b2f4
0804b6ac
0804b6d0
0804be50
0805eb08
0805eb20
080610dc
8061818
08065ef0
08065ef4
08065ef8
08065efc
08065ffc
08066000
08066200
08066d00
00000000
00000000
000000
000154
000168
000188
0001ac
000918
001de8
002f82
00321c
0032bc
0032f4
0036ac
0036d0
003e50
016b08
016b20
0190dc
019818
01cef0
01cef4
01cef8
01cefc
01cffc
01d000
01d200
01dcf8
01dcf8
01dd04
000000
000013
000020
000024
00076c

0014d0
00119a
00029a
0000a0
000038
0003b8
000024
000780
012cb8
000015
0025bc
00073c
002964
000004
000004
000004
000100
000004
0001e8
000af8
000be0
00000c
0000fc
00
00
00
00
04
10
00
02
00
08
08
00
04
00
00
00
00
00
00
00
00
00
00
00
00
00
08
04
04
00
00
00
00
A
A
A
A

A

A

A

A

A

A

AX

AX

AX

AX

A

A

A

WA

WA

WA

WA

WA

WA

WA

WA

0

0

0

0

5

6

0

5

6

5

5

0

0

0

0

0

0

0

0

0

0

6

0

0

0

0

0

0

0

0

0

0

0

0

1

0
0
2
0
12
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
1
4
4
4
4
4
1
2
4
4
4
4
16
16
4
32
4
4
4
4
4
4
4
4
4
32
32
1
1
59

Key to Flags:
W (write), A (alloc), X (execute), M (merge), S (strings), I (info),
L (link order), O (extra OS processing required), G (group), T (TLS),
C (compressed), x (unknown), o (OS specific), E (exclude **81**), p

(processor specific)

Команда для отображения содержимого секции в шестнадцатеричном формате

представления:

```
readelf --hex-dump .data htop
```

Hex dump of section '.data':

```
0x08066200 00000000 00000000 00000000 00000000 .....
0x08066210 00000000 00000000 00000000 00000000 .....
0x08066220 8cef0508 8cef0508 8cef0508 88eb0508 .....
0x08066230 8cef0508 8cef0508 8feb0508 96eb0508 .....
0x08066240 9deb0508 a4eb0508 00000000 00000000 .....
0x08066250 00000000 00000000 00000000 00000000 .....
0x08066260 8cef0508 8cef0508 8cef0508 8cef0508 .....
0x08066270 abeb0508 b2eb0508 8cef0508 8cef0508 .....
0x08066280 8cef0508 a4eb0508 00000000 00000000 .....
0x08066290 00000000 00000000 00000000 00000000 .....
```

...

3.1.1 Представление исходного кода на языке ассемблера

Структура ELF-файла была подробно рассмотрена в предыдущем пункте, однако любой файл имеет представление и в других форматах, например, в виде бинарной последовательности, декомпилированного или дизассемблированного кода, набора последовательностей выполняемых функций и т.д.

Как было отмечено ранее, в настоящем исследовании в качестве идентификационных характеристик исполняемых файлов были выбраны ассемблерные команды, таким образом, ниже будет приведено представление программы Htop в ассемблерном коде.

Язык ассемблера представляет собой низкоуровневый язык машинно-ориентированного программирования, его команды прямо соответствуют отдельным командам машины или их последовательностям. Синтаксис, в общем виде, представляет собой последовательность из адреса – команды – значения.

60

Представление части программы Htop в дизассемблированном виде, где команды выделены синим цветом:

```
htop: file format elf32-i386
```

Disassembly of section .interp:

```
<.interp>:
2f
6c
69 62 2f 6c 64 2d 6c
69 6e 75 78 2e 73 6f
2e 32 00
das
insb
imul
imul
xor
(%dx),%es:(%edi)
$0x6c2d646c,0x2f(%edx),%esp
$0x6f732e78,0x75(%esi),%ebp
%cs:(%eax),%al
```

Disassembly of section .data:

```
<__data_start>:
```

...

```
8c ef
05 08 8c ef 05
08 8c ef 05 08 88 eb
```

05 08 8c ef 05
08 8c ef 05 08 8f eb
05 08 96 eb 05
08 9d eb 05 08 a4
eb 05
08 00
...
00 00
00 8c ef 05 08 8c ef
05 08 8c ef 05
08 8c ef 05 08 ab eb
05 08 b2 eb 05
08 8c ef 05 08 8c ef
05 08 8c ef 56 05
08 a4 eb 05 08 00 00
...
00 00
8c ef
05 08 8c ef 05
08 8c ef 05 08 8c ef
05 08 8c ef 05
08 8c ef 05 08 8c ef
05 08 8c ef 05
08 8c ef 56 05 08 a4 eb
05 08 00 00 00
...
mov
add
or
add
or
add
or
jmp
or
add
add
add
or
add
or
add
mov
add
or
add
or
add
or
add
mov
add
or
add
or
add
or
add
%gs,%edi
\$0x5ef8c08,%eax
%cl,-0x1477f7fb(%edi,%ebp,8)

```
$0x5ef8c08,%eax
%cl,-0x1470f7fb(%edi,%ebp,8)
$0x5eb9608,%eax
%bl,-0x5bf7fa15(%ebp)
806624c
%al,(%eax)
%al,(%eax)
%cl,-0x1073f7fb(%edi,%ebp,8)
$0x5ef8c08,%eax
%cl,-0x1454f7fb(%edi,%ebp,8)
$0x5ebb208,%eax
%cl,-0x1073f7fb(%edi,%ebp,8)
$0x5ef8c08,%eax
%ah,0x805(%ebx,%ebp,8)
%al,(%eax)
%gs,%edi
$0x5ef8c08,%eax
%cl,-0x1073f7fb(%edi,%ebp,8)
$0x5ef8c08,%eax
%cl,-0x1073f7fb(%edi,%ebp,8)
$0x5ef8c08,%eax
%cl,-0x145bf7fb(%edi,%ebp,8)
$0x8,%eax
```

61

00 8c ef 05 08 8c ef

05 08 8c ef 05

08 8c ef 05 08 8c ef

05 08 8c ef 05

08 8c ef 05 08 8c ef

05 08 8c ef 56 05

08 a4 eb 05 08 00 00

...

add

add

or

add

or

add

or

```
%cl,-0x1073f7fb(%edi,%ebp,8)
```

```
$0x5ef8c08,%eax
```

```
%cl,-0x1073f7fb(%edi,%ebp,8)
```

```
$0x5ef8c08,%eax
```

```
%cl,-0x1073f7fb(%edi,%ebp,8)
```

```
$0x5ef8c08,%eax
```

```
%ah,0x805(%ebx,%ebp,8)
```

3.1.2 Анализ особенностей ассемблерного кода программного обеспечения

Представление программы в том или ином виде (бинарном, ассемблерном, на уровне инструкций высокоуровневого языка программирования, структуры формата файла и других), обладает некоторыми особенностями, которые будут рассмотрены ниже.

Уровень абстракции. Наличие смысловых конструкций, упрощающих понимание структуры данных и операций над ними для человека. Высокий уровень абстракции используется в высокоуровневых языках программирования для описания структур данных и операций над ними, описание которых на машинном

коде очень длинны и сложны для понимания человеком. Здесь присутствует широкий выбор команд, а имена переменных задаются программистом, в дальнейшем такой программный код компилируется и переводится на низкоуровневый язык (например, ассемблер или байт-код). Чем выше уровень абстракции, тем разнообразнее можно реализовать некоторую функцию в программе, необходимо принимать больше усилий в детальном исследовании всех процедур и данных программы. В то же время, представление исполняемого файла на низкоуровневом языке программирования значительно сокращает число возможных используемых значений признака, в том числе таких, чьи функциональные назначения носят схожий характер. Такой код становится легче анализировать и автоматизировать подход к идентификации, не привлекая к работе человека;

Диапазон признака. Количество вариаций некоторой характеристики файла носит ограниченный характер. Известно, что степень свободы в статистике – число

62

количества значений, которые могут варьироваться, играет важную роль в количестве информации, носимой признаковым пространством. Так при рассмотрении признака, имеющего три градации (степень свободы здесь равна двум) количество информации, которое будет носить построенная на нем сигнатура, будет ниже, чем при использовании признака с десятью градациями (степень свободы здесь равна девяти);

Уникальность и информативность значения признака. Значение признака может не носить никакой информативной составляющей, как например, значение байта «00», присутствие которого во всех файлах носит скорее заполняющий характер, нежели информационный. Очевидно, что некоторые значения признака будут обладать более существенными идентифицирующими способностями, таковыми обычно являются редко встречающиеся градации признака, они и обладают наибольшей информативностью. Однако в рамках разрабатываемой методики идентификации следует учитывать присущую разнообразность исполняемым файлам, участвующим в процессе аудита электронных носителей информации. Программы отличаются не только функциональной направленностью, но и объемом, что существенно влияет на вид формируемых сигнатур при использовании разработанного метода, описанного в разделе 2.3. Если в качестве признака выбрать редко встречаемую характеристику файла, то формируемая сигнатура будет содержать большое число нулей, не несущих никакой важной информации. При этом есть большая вероятность, что для исполняемых файлов небольшого объема и сильно отличающихся функциональностью, выбранный признак и вовсе будет давать нулевую сигнатуру, анализ и сравнение которой будет невозможно.

В данной работе производится исследование представления исполняемых файлов в их ассемблерном виде. Так же коды программ рассматриваются как целый объем, так и разделенный на определенные части.

Анализ программы как целого объема данных имеет некоторое преимущество, выраженное в том, что формируемая сигнатура независима от внутреннего расположения блоков данных программы, это актуально для случая

63

выпуска новой версии программы, содержащей одну или несколько частей измененного кода. При разбиении исследуемого кода файла на части и исследование признака на каждой из них, уже не будет обладать данным преимуществом. Однако анализ распределения признака по коду программы является достаточно информативным показателем.

3.2 Подход к представлению программного обеспечения

В исследование приняло участие 578 исполняемых файла ОС Linux, из них обучающей выборке принадлежало 455 файла различных версий и разрядностей (32x и 64x), относящихся к 63 различным программам. К тестовой выборке же

принадлежало 123 файла, относящихся к тем же 63 программам, все они были отличны от задействованных файлов, используемых в обучающей выборке, и имели 32х и 64х разрядности. Весь объем файлов был сформирован формировалась путем скачивания программ с официальных репозиториях ОС Linux для архитектур процессоров x86 и x86-64. Затем, для формирования тестовой выборки она отделялась от обучающей путем извлечения по два исполняемых файла различных версий и разрядностей по каждой программе (кроме одной, представленной только в одной разрядности) из всего объема скаченных файлов.

Стоит отметить, что архив сигнатур не является фиксированным и запрещенным для добавления новых сигнатур, наоборот, в реальных условиях при проведении аудита электронных носителей информации, он должен периодически обновляться и иметь актуальные данные для выполнения поставленных исследователем задач.

Важным также является количество программ и исполняемых файлов, участвующих в процессе аудита. Очевидно, что компьютер пользователя или же совокупность всех компьютеров организации, не смотря на все разнообразие существующих, содержит ограниченное число программ. На основе легитимных из них происходит формирование АЭСП. Также учитывая тот факт, что каждая из программ реализуется в различных версиях, число которых часто не

64 превышает и двух десятков, работа в данной области неизбежно происходит на ограниченном наборе данных, увеличение которого в практических реалиях просто не представляется возможным.

Рисунок 6 – Распределение исполняемых файлов по их размеру

Выше на рисунке 6 представлено распределение размеров участвующих в исследовании программ. Как видно их объем измеряется от десятков байтов, до тысяч. Каждый маркер на рисунке отображает один файл из 578 скаченных файлов. Данные приведены в логарифмической шкале и видно, что наиболее большое число исполняемых файлов имеет размер в пределах 104 – 106 байтов (желтая и оранжевая области линейчатой диаграммы с накоплением), так средним показателем является 7,24·10⁵ или 707 Кб.

Необходимо подчеркнуть, что весь объем исследуемых файлов подчиняется закону Бенфорда, который описывает вероятность появления определенной первой значащей цифры в распределениях величин, взятых из реальной жизни [120] [78,79]. Считается, что компьютерные файлы подчиняются этому закону [80], таким образом можно утверждать о разнообразии и репрезентативности изучаемого

65 объема файлов. В таблице представлены два распределения - Бенфорда и первых значащих цифр исследуемых файлов. Так, на основании критерия однородности хи-квадрат и уровне значимости $p = 0,05$, можно утверждать, что расхождения между теоретическим и экспериментальным распределениями вероятностей первой цифры для десятичной системы счисления, статистически недостоверны.

Таблица 3 – Сравнение двух распределений

Первая значащая цифра 1 2 3 4 5 6 7 8 9

Распределение Бенфорда 0,301 0,176 0,125 0,097 0,079 0,067 0,058 0,051 0,046

Распределение первых

значащих цифр

исследуемых файлов

0,336 0,111 0,123 0,130 0,093 0,050 0,073 0,048 0,036

3.2.1 Модель представления исполняемого файла

Исходя из понятий векторной модели, каждый исполняемый файл можно представить в качестве точки n -мерного пространства действительных чисел, т.е. вектора, где n – количество градаций признака и являющееся одинаковым для всех рассматриваемых исполняемых файлов.

Если ввести в рассмотрение n -мерное признаковое пространство

$F = (f_1, f_2, \dots, f_n)$, где $i = 1, 2, \dots, n$ – число градаций признака, тогда каждый объект o (файл vm или ej) в данном пространстве будет отображаться точкой с координатами $o = (f_1(o), f_2(o), \dots, f_n(o))$, а каждый класс объектов, т.е. имя программы Рэтал., i множеством таких точек.

Любой исполняемый файл может быть представлен как совокупность некоторых количественных значений. Ранее было выделено, что наиболее подходящим способом сбора характеристик ПО в данном исследовании является статический сбор частотных характеристик. Нескольким версиям одной и той же программы свойственны определенные характерные особенности. Измерение частот этих особенностей – градаций признака – позволяет нам представить исполняемый файл, как последовательность частот характеристик и перенести его вектор в пространство. Исполняемому файлу может быть поставлена в

66 соответствие точка в пространстве, координатами которой являются частоты его градаций признака.

Считая, что идентифицируемый файл обладает некоторыми свойственными ему характеристиками, можно сопоставить ему определенные характеристики конкретных программ. Идентификация файла происходит путем применения выбранного алгоритма сравнения к двум векторам значений признака объектов: o – идентифицируемого и o' – эталонного.

3.2.2 Идентификационное признаковое пространство

В качестве значений идентифицирующих признаков выступают статистические характеристики кода программ, такие как частота признака на всем объеме кода программы, частота признака в частях кода программы; также сюда входит такая градация признака, как наименование ассемблерной команды (aaa, aad, ..., xor).

Для проведения аудита электронных носителей информации главной задачей является определение такого набора характеристик ELF-фалов, выступающих в качестве соответствующих ему признаков и впоследствии включаемых в ЭСП. Очевидно, что выявление признаков, позволяющих достичь максимальной точности идентификации, является крайне важным. На качество идентификации могут оказать влияние множество факторов, таких как: размер программы, общее количество проверяемых исполняемых файлов, количество сформированных ЭСП в АЭСП, выбранный признак и его градация, способ построения СП и метод сравнения сигнатур.

В качестве идентификатора сопоставляемого, определенной программе, выступает совокупность выявленных характерных особенностей исполняемых файлов, являющихся версиями данной программы. Важным является тот факт, что 139 последующая версия программы должна носить дополняющий предыдущую версию характер, т.е. если создатель программы решил полностью сменить функциональность программы или внес в нее некоторое критически большое число

67 изменений, такая программа будет считаться новой и для нее необходимо построение отдельной сигнатуры.

Как было представлено выше исполняемый файл является сложной системой, в основе которой лежат последовательности структур данных и операций над ними. Многие из этих функций проявляются в виде конкретных наборов ассемблерных команд, которые могут быть использованы при 139 решении задачи идентификации ПО, однако, 139 в то же время некоторые 139 характеристики являются однотипными, использующиеся во многих различных между собой программах.

Учитывая возможное количество ПО установленного на компьютере пользователя и возможный список легитимного ПО сохраняющегося в АЭСП в работе предлагается использовать сигнатуры, размер которых согласуется с

объемами анализируемых файлов и возможными частотами встречаемости характеристик файлов, а также имеет фиксированную длину для всех сравниваемых исполняемых файлов.

С целью повышения точности идентификации и повышения скорости проведения сравнения сигнатур, в работе предложено и обосновано применение новых методов по формированию СИИФ и ЭСП из признаков, а также построение унифицированных сигнатур (УС), рассмотрение которых будет в пункте 2.3.2.

Возможными характеристиками эталонных программ являются наименования ассемблерных команд $As.com$. Так, в качестве совокупности характеристик версии программы могут выступать:

$vm = (As.com)$, где $As.com = (ac1, ac2, \dots, ac118)$ – последовательность из 118 отобранных ассемблерных команд. Сигнатура представляет собой частотную последовательность 118 ассемблерных команд на всем объеме кода программы;

$vm = (asp, 1, asp, 2, \dots, asp, k)$, где asp, k – n -ая выбранная ассемблерная команда из 118 ассемблерных команд $asp \in As.com$, k – число интервалов, на которые был разделен код программы. Сигнатура представляет собой частотную

68

последовательность выбранной ассемблерной команды в определенных интервалах кода программы.

3.2.3 Модель представления программного обеспечения

Идентифицируемый исполняемый файл обладает характеристиками, свойственными некоторой программе. Таким образом, можно детерминировать некоторую программу $R_{etal,i}$ как совокупность ее выпущенных версий vm , представляемых в виде набора, координатам которого сопоставляется множество характеристик F . Схематическое отображение модели представления программного обеспечения изображено на рисунке 7.

Рисунок 7 – Модель представления программного обеспечения

69

3.3 Метод формирования сигнатур исполняемых файлов, обладающий достаточной гибкостью по распознаванию различных версий программ, а также позволяющий наиболее точно производить идентификацию

Предлагается производить отбор признаков, которые обладают наиболее высокой различающей способностью и обеспечивающие достаточную частоту встречаемости для формирования не нулевых сигнатур исполняемых файлов.

В исследовании рассматривается несколько подходов к идентификации ELF-файлов, каждый из которых обладает своими ограничениями, которые следует учитывать еще на этапе формирования сигнатур файлов. Так, в данном разделе будут описаны следующие подходы к формированию сигнатур:

Формирование унифицированных сигнатур программ, объединяющих в себе несколько частотных характеристик различных файлов версий одной программы. Решение о группировке тех или иных частотных последовательностей принимается на основе метода кластеризации k -средних. Сигнатуры самих исполняемых файлов строятся на одном из следующих подходов:

- 1) подсчет распределения частоты появления 118 ассемблерных команд на всем объеме дизассемблированного кода файла;
- 2) подсчет распределения частоты появления одной ассемблерной команды на k не равных интервалах дизассемблированного кода файла;
- 3) подсчет распределения частоты появления одной ассемблерной команды на k равных интервалах дизассемблированного кода файла.

Формирование сигнатур по каждой версиям программ, которые строятся на подсчете распределения частоты появления одной ассемблерной команды на k равных интервалах дизассемблированного кода файла.

Среди всего многообразия ассемблерных команд были отобраны 118, которые присущи файлам обучающей выборки. Однако, в ходе работы над

исследованием, было решено проанализировать информативность этих команд – данный этап описан в пункте 2.3.1, и произвести отбор идентификационных признаков с целью разработки метода по формированию СП на основе

70

распределения ассемблерной команды в определенных интервалах кода программы.

На рисунке 8 представлены графики распределений частот выбранных ассемблерных команд для одной версии программы Атагок. На оси абсцисс отмечены номера интервалов, на которые был разделен код программы k от 1 до 30, а на оси ординат – частоты встречаемости той или иной ассемблерной команды $асі$.

Рисунок 8 – Частотные распределения ассемблерных команд.

Из рисунка 8 видно, что характеры распределений разных ассемблерных команд различаются между собой, одни команды располагаются в основном в середине программы, другие в начале и в конце.

3.3.1 Критерии и алгоритм отбора наиболее информативных признаков

Информативность признака может выражаться в различных величинах, в зависимости от используемого подхода:

71

энергетический – где информативность оценивается по величине признака. Здесь наиболее информативным признаком будет считаться тот, чье абсолютное значение наибольшее. Из этого становится понятным, что использование энергетического подхода к задаче идентификации ПО не пригодно, ведь в ситуации, когда значение какого-либо признака велико по абсолютной величине и практически не отличается у файлов, относящихся к различным программам, идентификация объекта к какому-либо одному классу является практически невозможной;

информационный – где информативность признака оценивается по величине достоверного различия между классами образов в пространстве признаков.

Наиболее пригодным подходом для оценивания информативности признаков, является информационный подход, позволяющий производить отбор признаков, обладающих максимальной различающей способностью именно на определенной группе программ, задействованных в формировании АЭСП.

Одним из ключевых вопросов является выбор алгоритма вычисления информативности 118 ассемблерных команд. Основываясь на проведенных исследованиях [81], [82], был выбран метод Шеннона, не зависящий от объема выборки наблюдаемых признаков, а также он позволяет вычислить информативность признака для произвольного числа классов. Данный алгоритм оценивает информативность как средневзвешенное количество информации (величину устраненной энтропии), свойственной рассматриваемой градации признака $f \in F$ и предоставляет оценку информативности $I(f)$ в виде нормированной величины, принимающей значения в интервале от нуля (отсутствие информативности) до единицы (высокая степень информативности).

Для оценки информативности признака f пользуются следующей формулой:

$$I(f) = 1 + \sum_{i=1}^U (P_i \cdot \sum_{u=1}^n P_{i,u} \cdot \log U P_{i,u})$$

У

n $u=1$)

$i=1, (1)$

72

где n – число градаций признака; U – число программ, участвующих в

подсчете информативности; P_i – вероятность попадания значения признака в i -ю

градацию, которая вычисляется по формуле:

$P_i =$

$\sum m_{i,u}$

U

u=1

N

.

где $m_{i,u}$ – частота появления значения признака в i-й градации в u-ом классе **87**;

N – общее число наблюдений признака; $P_{i,u}$ – вероятность появления i-й градации признака в k-ой программе, которая вычисляется по формуле:

 $P_{i,u} =$ $m_{i,u}$ $\sum m_{i,u}$

U

u=1

.

Очевидно, что произвести оценку информативности признаков на всем пространстве существующих программ и их версий, является невозможным, поэтому в расчетах принимало ограниченное число исследуемых файлов. Было отобрано 10 различных программ и версий, из которых было сформировано 52 частотных распределения характеристики – ассемблерной команды, путем подсчета каждой из 118 ассемблерных команд на всем объеме дизассемблированных кодов исполняемых файлов.

Решение задачи по отбору наиболее информативных признаков разбивается на два случая:

первый, когда мы принимаем число градаций признака n равным 1, т.е.

появление выбранной ассемблерной команды;

второй, когда мы принимаем число градаций признака n равным 2, т.е.

появление выбранной ассемблерной команды и появление другой, отличной от данной, команды.

Далее, в таблицах 4 и 5 приведены результаты расчета информативности для $n = 1$ и $n = 2$ соответственно. Все ассемблерные команды расположены в порядке возрастания их информативности. Полный перечень 118 ассемблерных команд и их информативности по Шеннону можно найти в приложении А.

73

Таблица 4 – Значения информативности для 118 ассемблерных команд, по методу Шеннона и заданным числом градаций признака $n = 1$

Ассемблерные команды Информативность, I(x)

cmpsb, cmpsw, esc, jc, jcxz, jna, jnae, jnb, jnbe, jnc, jng, jnge, jnl,

jnz, jpe, jpo, jz, lodsb, lodsw, loopnz, loopz, movsb, movsw,

repe, repne, retn, sal, scasb, scasw, stosb, stosw, wait

0

rep, jle, dec, lea, std, cmp, shr, jmp, jg, shl, and, cld, rol, js, jl, add,

pop, mul

0,1408 – 0,1990

jne, push, ret, sub, xor, jge, ror, not, test, mov, div, jbe, jb, loopne,

clc, je, jae, xchg, ja, or, repz, cli, adc, pushf, hlt, sar, sti, sbb, neg,

jns, idiv, in, call, imul, jp, loop, jno, pop, repnz, out, ired, loope, jnp,

xlat, int, rcl, cmc

0,2001 – 0,2972

jo, stc, lahf, retf, rcr, sahf, inc, cwd, popf, cbw 0,3026 – 0,3778

lock 0,4228

daa, les, into, aaa, lds, aam, aas, das, aad 0,6144 – 0,6785

Таблица 5 – Значения информативности для 118 ассемблерных команд, по методу Шеннона и заданным числом градаций признака $n = 2$

Ассемблерные команды Информативность, I(x)

cmpsb, cmpsw, esc, jc, jcxz, jna, jnae, jnb, jnbe, jnc, jng, jnge, jnl,

jnl, jnz, jpe, jpo, jz, lodsb, lodsw, aaa, les, daa, aas, aam, ja, mul,
js, loop, jge, jp, hlt, jl, jns, loopne, ired, das, rep, jae, idiv, loopnz,
loopz, movsb, movsw, repe, repne, ret, sal, scasb, scasw, stosb,
stosw, wait, cwd, cbw, div, neg, into, lds, aad

0,2024

rcl, std, sbb, lahf, ror, ret, jne, cmc, popf, rcr, adc, jno, in, shl, jbe,
jb, pushf, not, clc, rol, int, jle, jg, sub, loope, xlat, cld, sti, shr, jnp,
repnz, cli, dec, repz

0,2025

imul, sar, sahf, jmp, xor, nop, and, jo, je, stc, test, push, retf 0,2026
out, cmp, inc 0,2027

or 0,2028

xchg 0,2029

add, pop 0,203

lea 0,2032

call 0,2036

mov 0,2044

lock 0,2095

Из приведенных результатов можно сделать вывод, что значение числа градаций несет важную роль, интерпретировать которую следует как результат существенного влияния на параметры P_i и P_i, u в формуле (1) дополнительной градации, т.е. не появление рассматриваемого признака, а появление другой, отличной ассемблерной команды. Таким образом, происходит учет отношения числа встречаемости рассматриваемой ассемблерной команды к числу появления

74

всех остальных 117 ассемблерных команд, что позволяет оценить информативность не только на основе различия частот встречаемости команды в различных классах, но и на основе ее доли по отношению к остальным командам. Стоит отметить, что использование наиболее информативных команд из таблицы 4 невозможно, т.к. их частота встречаемости в дизассемблированных кодах программ слишком мала для формирования сигнатур с достаточным числом ненулевых значений. Использование второго же подхода с числом градаций признака $n = 2$ позволяет устранить данный недостаток.

Имена ассемблерных команд, выделенные курсивом, были отобраны как наиболее подходящие для решения задачи данной научной работы. Основываясь на результатах таблиц 4 и 5, примем следующий порядок информативности выбранных 10 ассемблерных команд: *mov*, *call*, *pop*, *push*, *je*, *lea*, *add*, *cmp*, *and*, *jmp*.

3.3.2 Формирование унифицированных сигнатур программ (УСП), используемого в дальнейшем в статистическом методе сравнения сигнатур и с использованием искусственной нейронной сети

Формирование УСП, объединяющих в себе несколько частотных характеристик различных файлов версий одной программы. Решение о группировке тех или иных частотных последовательностей принимается на основе метода кластеризации k -средних. Среди всех версий одной программы выделяются кластеры, на основе которых происходит формирование одного варианта унифицированной ЭСП.

Подход на основе 118 ассемблерных команд.

Первым подходом к формированию сигнатур исполняемых файлов рассмотрим использование в качестве характеристики F файла ассемблерный код *As.com*, а именно его команды $ac_1, ac_2, \dots, ac_{118}$ в качестве градаций признака. Таким образом, здесь $F = (f_1, f_2, \dots, f_n) = As.com = (ac_1, ac_2, \dots, ac_n)$, $n = 118$ – является признаковым пространством, где значение признака $L(ac_n) \in N_0$ и показывает количественное значение для градации признака *asp*.

75

На основе отобранных 118 градаций признака происходит формирование как

сигнатур программ обучающей выборки, так и сигнатур идентифицируемых файлов тестовой выборки. Однако, ЭСП будут создаваться путем унификации сигнатур исполняемых файлов обучающей выборки и добавления наименования известной программы.

Пусть обучающей выборке программ $P = \{P_{\text{Этал.}}, 1, P_{\text{Этал.}}, 2, \dots, P_{\text{Этал.}}, i\}$ присвоены имена этих программ $N = \{N_1, N_2, \dots, N_i\}$. Каждая отдельная программа $P_{\text{Этал.}}, i$ состоит из ряда различных версий v_1, v_2, \dots, v_m , тогда версия v_m будет описываться n градациями признака $F = As.com = (ac_1, ac_2, \dots, ac_{118})$.

Процесс создания сигнатур всех программ обучающей выборки будет состоять в следующих шагах (описание дается для каждой рассматриваемой программы в отдельности):

1) Каждая версия программы v_m дизассемблируется.

2) Из полученного представления кода выделяется частотная

характеристика выбранного признака F (118 ассемблерных команд), являющаяся основой для формирования сигнатуры одного исполняемого файла

$S(v_m) = q(v_m, As.com) = (ac_1, ac_2, \dots, ac_{118})$, где $q()$ – функция $q(v_m, F): As.com \rightarrow N_0$,

определяющая алгоритм формирования количественного значения для ac_n

ассемблерной команды, где $As.com$ – множество принимаемых значений

характеристики F . Сигнатура версии программы принимает вид:

$S(v_m) = (L(ac_1), L(ac_2), \dots, L(ac_n))$, $n = 118$,

где $L(ac_n)$ – значение частоты встречаемости градации признака ac_n в

дизассемблированной коде файла v_m .

Однако, многообразие версий ПО, растущие функциональные возможности

компьютерных систем, различные временные промежутки выпуска новой версии

ПО и т.п. эти факторы вносят существенные коррективы в содержательную часть

кода программ. Так, выпуск новых версий некоторой программы может содержать

в себе существенные изменения в отличии от выпущенных ранее версий. Следует

учитывать данный нюанс и при последующем формировании унифицированных

76

сигнатур из сигнатур версий $S(v_m)$ производить оценку отклонения значения

частоты ассемблерной команды $L(ac_n)$ файла v_m от значения средней частоты

ассемблерной команды $L(ac_n)$ программы $P_{\text{Этал.}}, i$ по формулам:

а) Расчет средних частот $L(ac_n)$ для каждой из ассемблерных команд по всем

$S(v_m)$

$L(ac_n) =$

1

m

$\sum m S(v_m)[ac_n]$

1,

получаем строку со 118-ю значениями: $S'(P_{\text{Этал.}}, i) = (L(ac_1), L(ac_2), \dots, L(ac_n))$.

б) Расчет промежуточной унифицированной сигнатуры

$S_{\text{униф.}}(P_{\text{Этал.}}, i) = (p(ac_1), p(ac_2), \dots, p(ac_n))$

$p(ac_n) =$

{

$L(ac_n)$, если число не нулевых элементов $Q(v_m)[ac_n] \geq 0,85 \cdot m$

0, если число не нулевых элементов $Q(v_m)[ac_n] < 0,85 \cdot m$

, (2)

где $Q(v_m) = (q_1, q_2, \dots, q_n)$ множество элементов

$q_n = \{$

$L(ac_n), |L(ac_n) - L(ac_n)| \leq 0,2 \cdot L(ac_n)$

0, $|L(ac_n) - L(ac_n)| > 0,2 \cdot L(ac_n)$

,

таких же, как и в сигнатуре версии программы $S'(v_m)$, но обнуляемые при

невыполнении условия о попадании частоты ассемблерной команды $L(ac_n)$ в

заданный интервал, ограниченный величиной допустимого отклонения от среднего

показателя частоты ассемблерной команды для программы L (асп).

Смысл формулы (2) состоит в следующем: если количество ненулевых значений q_n будет не менее 85% от общего числа файлов vm в выборке программы Рэтал, i , то в промежуточную унифицируемую сигнатуру $S_{\text{униф.}}(\text{Рэтал},i)$ записывается средняя частота встречаемости каждой из ассемблерных команд для всех версий $S(vm)$.

Коэффициенты 0,2 и 0,85 в формулах были вычислены экспериментальным путем.

3) Окончательно унифицированная сигнатура программы имеет вид:

77

$S_{\text{униф.}}(\text{Рэтал},i) = (N_i, S_{\text{униф.}}(\text{Рэтал},i)) = (N_i, p(ас1), p(ас2), \dots, p(асn))$.

Все сформированные таким образом унифицированные сигнатуры эталонных программ помещаются в АЭСР, для дальнейшего обращения к нему либо в процессе идентификации, либо при необходимости модификации сигнатур.

В соответствии с описанными выше этапом под пунктом 2) сигнатура идентифицируемого файла имеет вид:

$S(e_j) = (L(ас1), L(ас2), \dots, L(асn)), n = 118$.

Подход на основе одной ассемблерной команды на неравных интервалах.

В разделе 2.3 было приведено обоснование цели исследования подхода к формированию сигнатур программ на основе одной ассемблерной команды. Здесь и в следующем подпункте мы рассмотрим данный подход.

В отличие от способа построения УСП по частотам встречаемости в дизассемблированном коде 118 различных ассемблерных команд, следующий подход к исследованию характеристик исполняемых файлов основан на анализе одной ассемблерной команды, выбранной в качестве признака для формирования частотных распределений. Таким образом, здесь F представляет собой распределение ассемблерной команды на интервалах ассемблерного кода и $F = (f_1, f_2, \dots, f_n) = асп = (асп,1, асп,2, \dots, асп,k)$, асп – ассемблерная команда $асп \in As.com$, k – число интервалов разбиения дизассемблированного кода программы Рэтал, i , где значение признака $L(асп) \in N_0$ и показывает количественное значение для признака асп.

На основе отобранного признака происходит формирование как сигнатур программ обучающей выборки, так и сигнатур идентифицируемых файлов тестовой выборки. Однако, ЭСП будут создаваться путем унификации сигнатур исполняемых файлов обучающей выборки и добавления наименования известной программы.

Пусть обучающей выборке программ $P = \{\text{Рэтал},1, \text{Рэтал},2, \dots, \text{Рэтал},i\}$ присвоены имена этих программ $N = \{N_1, N_2, \dots, N_i\}$. Каждая отдельная программа Рэтал, i

78

состоит из ряда различных версий v_1, v_2, \dots, v_m , тогда версия vm будет описываться одной градацией признака $As.com$ на n интервалах разбиения дизассемблированного кода программы $F = асп = (асп,1, асп,2, \dots, асп,k)$.

Опишем процесс формирования сигнатур файлов обучающей выборки:

- 1) Первоначально фиксируется один из файлов vm_0 версии программы и происходит его дизассемблирование.
- 2) Полученный ассемблерный код программы разбивается на интервалы неравных длин, но с фиксированной частотой появления в них заданного признака асп (ассемблерной команды). В этом случае за длину интервала принимается количество различных ассемблерных команд из $As.com$ на промежутке фиксированной частоты признака $L(асп)$. Таким образом, формируется частотная характеристика фиксированного файла vm_0 , имеющая вид $S(vm_0) = ((асп, h_1), (асп, h_2), \dots, (асп, h_k))$, где асп – заданная частота признака, равная константе; h_i – длина интервала разбиения; k – количество полученных интервалов разбиения.

Формирование $S'(vm_0)$ необходимо для того, чтобы стало возможным построение сигнатуры отдельной программы vm с фиксированным количеством интервалов разбиения k .

3) К остальным файлам версий ПО из обучающей выборки применяется следующая формула:

$$S'(vm) = ((acn,1, h'1), (acn,2, h'2), \dots, (acn,k, h'k)),$$

где acn,k – получаемая частота признака; $h'k =$

$$hk \cdot lm$$

,

l

– длина интервала

разбиения, l – длина vm_0 -го файла, hk – длина k -го интервала разбиения vm_0 -го файла,

$l'm$ – длина m -го файла vm ; k – количество полученных интервалов разбиения.

Таким образом, получаются распределения фиксированного файла $S'(vm_0)$ и остальных файлов $S'(vm)$ обучающей выборки, с одинаковой длиной k по отдельной программе Рэтал, i , из которых затем, на основе средних значений частот признака

79

$$L(acn) =$$

1

m

$$(S'(vm_0)[acn] + \sum S'(vm)[acn]$$

m

1

)

4) формируется УСП:

$$\text{Суниф}(P\text{этал},i) = (N_i, S'(P\text{этал},i)) =$$

$$(N_i, L(ac_1, h_1), L(ac_2, h_2), \dots, L(ac_k, h_k)).$$

Все сформированные таким образом сигнатуры помещают в АЭСП.

Особенностью сигнатур идентифицируемых файлов является то, что их длины должны совпадать с длиной сигнатуры соответствующей программы из архива, поэтому сигнатура идентифицируемого файла представляет собой следующую структуру в соответствии с описанными выше этапом под пунктом 3):

$$S(e_j) = ((acn,1, h'1), (acn,2, h'2), \dots, (acn,k, h'k))$$

и является распределением частот ассемблерной команды в интервалах, заданных в сигнатуре программы из архива, где acn – частота признака; $h'k =$

$hk \cdot l'$

l

– длина интервала разбиения, l' – сумма длин интервалов hk , заданных в сигнатуре соответствующей программы из архива Суниф(Pэтал, i), hk – длина k -го интервала разбиения сигнатуры программы из архива Суниф(Pэтал, i), l' – длина идентифицируемого файла e_j .

Заметим, что предложенный особый метод формирования сигнатур программ имеет некоторые нюансы:

во-первых, каждая сигнатура в архиве является статистической оценкой (аналогом) равномерного распределения частот некоторой ассемблерной команды, которая осуществлена на основе выборки, отличной от выборки частот идентифицируемого файла;

во-вторых, для каждого идентифицируемого файла преобразуется его длина относительно сравниваемой с ним программы, что приводит к одинаковой размерности сравниваемых сигнатур;

в-третьих, для каждой ЭСП из архива необходимо стоять свою СИИФ.

80

Подход на основе одной ассемблерной команды на равных интервалах.

Процедура формирования сигнатур исполняемых файлов схожа с выше

описанной, с тем лишь отличием, что частотное значение признака строится на равных интервалах ассемблерного кода программ. Градацией признака здесь являются интервалы разбиения и его присутствие в них. Таким образом, F представляет собой набор выделенного числа ассемблерных команд и

$$F = (f_1, f_2, \dots, f_n) = (ac_1, ac_2, \dots, ac_n), ac_n - \text{ассемблерная команда } ac_n \in As.com.$$

Формирование ЭСП происходит в соответствии со следующими шагами:

1) Каждый файл vm дизассемблируется и разбивается на интервалы равных длин.

При этом вводится фиксированный множитель для формирования длины шага, корректирующий количество получаемых интервалов k для файлов различного объема, где за длину интервала принимается количество различных ассемблерных команд на промежутке одного шага, k – число интервалов разбиения дизассемблированного кода программы $R_{\text{этл},i}$, где значение признака $L(ac_n) \in N_0$ и показывает количественное значение для признака ac_n . В исследовании коэффициент k был подобран таким образом, чтобы число интервалов равнялось тридцати.

2) Далее из полученного представления кода выделяется частотная характеристика выбранной градации признака ac_n (одной ассемблерной команды), являющаяся основой для формирования сигнатуры одного исполняемого файла $S'(vm) = (L(ac_n, 1), L(ac_n, 2), \dots, L(ac_n, k))$, где ac_n, k – значение частоты появления признака ac_n в k -ом интервале.

3) На основе сигнатур версий vm программы $R_{\text{этл},i}$ строится промежуточная УСП, представляющая собой $S'_{\text{униф.}}(R_{\text{этл},i}) = (p(ac_n, 1), p(ac_n, 2), \dots, p(ac_n, k))$, где $p(ac_n, k) = \{$

$$L(ac_n, k), | \min(S'(vm)[L(ac_n, k)]) - L(ac_n, k) | \text{ AND } | \max(S'(vm)[L(ac_n, k)]) - L(ac_n, k) | \leq \text{ratio} * L(ac_n)$$

$$0, | \min(S'(vm)[L(ac_n, k)]) - L(ac_n, k) | \text{ OR } | \max(S'(vm)[L(ac_n, k)]) - L(ac_n, k) | > \text{ratio} * L(ac_n)$$

,

81

где $L(ac_n, k) =$

1

m

$$\sum m S'(vm)[L(ac_n, k)]$$

1 – среднее значение для $L(ac_n)$ всех vm в

$R_{\text{этл},i}$; ratio – изменяемый коэффициент; AND – логическое «И»; OR – логическое «ИЛИ».

Смысл формулы (3) заключается в аннулировании тех значений частот, где абсолютная разница минимального или максимального значения элемента $S'(vm)[L(ac_n, k)]$ от среднего $L(ac_n, k)$ между файлами vm программы $R_{\text{этл},i}$ в k -ом интервале превышает установленную коэффициентом ratio долю от среднего показателя частоты признака в данном интервале для всех версий программы $L(ac_n, k)$.

4) Таким образом, унифицированная сигнатура программы примет вид:

$$S'_{\text{униф.}}(R_{\text{этл},i}) = (N_i, S'_{\text{униф.}}(R_{\text{этл},i})) = (N_i, p(ac_n, 1), p(ac_n, 2), \dots, p(ac_n, k)).$$

А сигнатура идентифицируемого файла, в соответствии с расчетами под пунктом 2) имеет вид:

$$S(e_j) = (L(ac_n, 1), L(ac_n, 2), \dots, L(ac_n, k)).$$

3.3.3 Формирование сигнатур программ (СП), используемого в дальнейшем в методе сравнения сигнатур на основе градиентного бустинга деревьев решений

Процесс формирования сигнатур обучающей и тестовой выборки для входных данных для модели обучения градиентного бустинга деревьев решений состоит из тех же этапов, что и подход на основе одной ассемблерной команды, анализируемой на равных интервалах разбиения дизассемблированного кода версий программ. Однако, в данном подходе не происходит формирования

унифицированных сигнатур, по причине того, что рассматриваемые алгоритмы лучше работают на большом наборе данных и их классифицирующая способность показала лучшие результаты именно на не унифицированных сигнатурах программ, а на сигнатурах версий этих программ.

82

Уточним, что сигнатурами файлов обучающей выборки являются структуры вида:

$$S(\text{Рэтэл}, i) = \{(N_i, L(\text{асп}, 1), L(\text{асп}, 2), \dots, L(\text{асп}, k)), \dots, \{N_i, L(\text{асп}, 1), L(\text{асп}, 2), \dots, L(\text{асп}, k)\}\},$$

а СИИФ:

$$S(e_j) = (L(\text{асп}, 1), L(\text{асп}, 2), \dots, L(\text{асп}, k)).$$

Входными параметрами обучаемых моделей, кроме наименований классов, могут быть лишь количественные меры, однако один из алгоритмов градиентного бустинга деревьев решений способен обрабатывать и категориальные значения.

3.4 Метод сравнения сигнатуры идентифицируемого исполняемого файла

(СИИФ) с эталонной сигнатурой программы (ЭСП)

Целью разработки данного метода является повышение точности

идентификации исполняемых файлов в условиях:

мульти-классификации. Каждое наименование программы – это

отдельный класс;

отсутствие класса вредоносных программ как такового. Считается, что

программы не обладают ярко выраженными признаками присущими вредоносному

ПО. Каждая из программ представляет собой легальный продукт, наличие в АС

которого является несанкционированным только в рамках установленных правил

действующей политики безопасности;

ограниченного набора проверяемых программ при проведении аудита. Из

приведенного выше материала видно, что АЭСП строится на основе файлов

обучающей выборки и распознавание идентифицируемого файла возможно лишь в

рамках количества изученных наименований программ;

постоянного выпуска новых версий ПО разработчиками. Производители

поддерживают выпускаемые ими программы, исправляют ошибки, оптимизируют

код и добавляют или убирают функционал – это приводит к изменению кода

83

программ, а с ним и значения выделяемых характеристик. Как уже было отмечено,

подходы, основанные на целостности файла, не предоставляют возможности

сравнивать две версии одной и той же программы. Присутствует необходимость в

постоянном мониторинге программных средств, установленных на АС, и

формировании актуальной базы таких показателей.

Из приведенных условий следует потребность в решении задачи по созданию

метода сравнения СИИФ с ЭСП, позволяющий производить с высокой точностью

идентификацию исполняемых файлов.

Задача сравнения двух сигнатур может рассматриваться, как:

задача проверки статистической гипотезы о не различности двух

распределений (сигнатур) или принадлежности известному закону распределения;

задача мульти-классификации.

Метод сравнения сигнатур, в зависимости от используемого подхода, может

включать в себя этап построения и обучения модели классификатора, а затем

определение меры близости того, что идентифицируемый файл является той или

иной программой.

3.4.1 Статистические методы сравнения сигнатур на основе критерия

однородности хи-квадрат и критерия Колмогорова

Одним из достаточно простых и подходящих в использовании

статистических критериев является критерий однородности хи-квадрат,

используемый как для дискретных, так и для непрерывных распределений и

способный сравнивать распределения характеристик файлов, представленных в

любой шкале, начиная от шкалы наименований [83], [84]. Однако, необходимо учитывать ограничение на равное число групп для сравниваемых распределений.

Так, если в результате эксперимента были получены две независимые выборки объемов n_1 и n_2 – представление двух сравниваемых файлов в ассемблерном виде, при этом по рассматриваемому признаку выборки

84
структурируются в k классов с частотами m_1, m_2, \dots, m_k и m'_1, m'_2, \dots, m'_k , то при расчете эмпирического значения χ^2 -критерия однородности используют формулу:

$$\chi^2 = n_1 \cdot n_2 \sum_{i=1}^k \frac{(m_i - m_i')^2}{m_i + m_i'}$$

где $m_1 + m_2 + \dots + m_k = n_1$ и $m'_1 + m'_2 + \dots + m'_k = n_2$.

Доказано, что эта статистика при больших значениях n_1 и n_2 распределена по закону χ^2 с $k - 1$ степенями свободы [83] и имеет правостороннюю критическую область, поэтому нет оснований отвергнуть гипотезу об однородности распределений если на заданном уровне значимости p выполняется неравенство $\chi^2 < \chi^2_p$.

Стоит заметить, что проверку гипотезы об однородности распределений по χ^2 -критерию можно применять не только к двум, но и к нескольким распределениям [84].

Еще одним интересным статистическим критерием является критерий согласия Колмогорова, предназначенного для проверки гипотезы о том, что значения исследуемой выборки подчиняются некоторому известному закону распределения. Процесс идентификации заключается в сравнении сигнатуры программы из АЭСП Суниф(Рэтал,1), в которой распределение частот ассемблерной команды является «эталонным» и имеет равномерный закон распределения за счет того, что частота подсчитывалась на интервалах не равных длин, с сигатурой идентифицируемого файла $S(ej)$ по критерию согласия Колмогорова.

В отличие от χ^2 -критерия, по которому сопоставлялись частоты двух распределений по каждому интервалу разбиения, здесь сопоставление происходит по накопленным частотам, т.е. сначала сопоставляются значения частот встречаемости ассемблерной команды по первому интервалу, затем по сумме первого и второго интервалов, далее по сумме первого, второго и третьего интервалов и т.д. Таким образом, критерий способен найти точку, в которой сумма накопленных расхождений между двумя распределениями является максимальной, и оценить достоверность данного расхождения [84].

85
Однако критерий Колмогорова-Смирнова, как и все статистические критерии, имеет ряд ограничений, наиболее существенным из которых, в контексте решаемой задачи по идентификации исполняемых файлов, является требование к формированию разрядов значений признака, которые должны быть упорядочены по нарастанию или убыванию, иначе накопление частот будет отражать лишь элемент случайного соседства разрядов [85].

Чтобы учесть данное требование, был разработан особый подход к формированию сигнатур файлов, основанный на формировании интервалов различных длин, но таких, что в каждом из них частота рассматриваемого признака была одинакова и имела вид равномерного распределения.

Статистикой критерия Колмогорова является величина, представляющая собой модуль максимального отклонения эмпирической функции распределения $F^*(x)$ от теоретической $F(x)$, которая при фиксированном объеме выборки n записывается следующим образом [84]:

$$D_n$$

$$D_n = \max_{-\infty < x < \infty} |F^*(x) - F(x)|$$

Учитывая, что критерий также имеет правостороннюю критическую область, то нет оснований отвергнуть гипотезу о согласии двух распределений, если на уровне значимости p выполняется следующее неравенство D_n

$$D_n < d_p/n$$

Отметим также, что при достаточно большом количестве независимых испытаний Колмогоровым была доказана независимость от закона распределения случайной величины, а вероятность $P(\lambda)$ неравенства $\lambda > \lambda$ приближенно выражается следующим образом и называется функцией обеспеченности [86]:

$$P(\lambda) = P(\lambda > \lambda) \approx 1 - \sum_{k=-\infty}^{\infty} (-1)^k e^{-2k^2 \lambda^2}$$

∞

$k=-\infty$

Случайная величина $\lambda = D_n$

$\sqrt{n} D_n$ может служить критерием согласия между

эмпирическим и теоретическим законами распределения. Если уровень значимости p принять равным 5%, то $k_p \approx 1,36$, если 1%, то $k_p \approx 1,63$.

86

3.4.2 Метод сравнения сигнатур на основе машинного обучения

Аппаратная или программная реализация некоторой математической модели, построенной на принципе организации и функционирования биологических нейронных сетей, называется искусственной нейронной сетью, по аналогии с сетями нервных клеток живого организма. Нейронная сеть является частью класса методов искусственного интеллекта и представляет собой обучаемую систему, т.е. ее функционирование происходит не только в соответствии с заданными алгоритмами и функциями, но и на основании полученного опыта.

Структура нейронной сети состоит из нейронов – элементов, представляющих собой простые процессоры и способные быть охарактеризованы каким-либо состоянием, по аналогии с клетками головного мозга. Нейрон обладает группой синапсов – однонаправленных входных связей, соединенных с выходами других нейронов. Также нейрон имеет аксон – выходную связь данного нейрона, с которой сигнал поступает на синапсы следующих нейронов [87]. На рисунке 9 показана модель многослойного перцептрона, структура которого была выбрана для задачи идентификации ПО экспериментальным путем [73].

87

Скрытые слои

...

...

Выходные

данные

Входные

данные

...

Синапсы Аксоны

Рисунок 9 – Модель многослойного перцептрона с прямой передачей

сигнала

Синапсы обладают двумя параметрами – задаваемой величиной

синоптической связи x_i и вычисляемым значением веса w_i , который соответствует

«силе» одной синоптической связи. Синапс способен оказывать на нейрон как

возбуждающее, так и тормозящее воздействие. Именно благодаря весу, происходит

изменение входного сигнала при его передаче к другому нейрону. Обучение

искусственной нейронной сети происходит путем подачи значений входных

переменных во входные элементы модели. Затем последовательно отработываются

нейроны промежуточных и выходного слоев, где каждый из них вычисляет

свое значение активации, беря взвешенную сумму $S = \sum_{i=1}^n x_i \cdot w_i$

n

$i=1$ выходов

элементов предыдущего слоя и вычитая из нее пороговое значение w_0 , где n – число

нейронов предыдущего слоя. Затем значение активации преобразуются с помощью

функции активации, и в результате получается выход нейрона. После проведения

всех итераций, выходные значения элементов выходного слоя принимаются за

выход всей сети в целом [87].

88

Другим популярным алгоритмом машинного обучения является деревья

решений, которые предлагают способ представления правил в иерархической

структуре, представляющей собой древовидный граф. Они содержат вершины, в

которые записываются проверяемые условия, ребра или ветви, являющиеся

переходами из одного узла дерева в другой, а также листья, в которые

записываются конечные значения, например, один из классов при решении задачи

классификации [89]. На рисунке 10 представлено схематическая модель дерева

решений.

Рисунок 10 – Модель дерева решений

Градиентный бустинг является общим методом для повышения

производительности, который может быть использован для любого алгоритма

машинного обучения. Суть его алгоритма заключается в обучении каждой

последующей модели с использованием данных об ошибках в предыдущих

моделях и дальнейшего снижения данных ошибок:

89

$$\text{obj}(\theta) = \sum_{i=1}^n l(y_i, y_i)$$

$$+ \sum_{k=1}^K \Omega(f_k) \rightarrow$$

$$k=1 \min,$$

где x – значения признаков обучаемой выборки; K – число деревьев; f –

функция в функциональном пространстве F . В данном исследовании для

программы XGBoost это функция Softmax:

$$P(y = j | x) =$$

$$\frac{\exp(\sum_{i=1}^n w_i)}{\sum_{i=1}^n \exp(w_i)}$$

$$\sum_{i=1}^n \exp(w_i)$$

$$i=1$$

коэффициент прогнозирования (prediction scores) для ансамбля деревьев

$$y_i = \sum_{k=1}^K f_k(x_i),$$

$$K$$

$$k=1$$

$$f_k \in F$$

w – весовые коэффициенты признаков; F – набор всех возможных CART

(Classification and Regression tree); n – число классов; l – функция потерь обучения

(training loss function); Ω – регуляризация (regularization term).

Считается, что градиентный бустинг, можно применить в отношении любого

слабого алгоритма в целях снижения его ошибки обучения [90]. Градиентный

бустинг деревьев решений позволяет строить аддитивную функцию в виде суммы

деревьев решений итерационно, по аналогии с методом градиентного спуска [91]. 93

Так на рисунке 11 схематично представлен ансамбль из k деревьев и принцип минимизации ошибки $obj(\theta)$.

Рисунок 11 – Минимизация ошибки обучения при использовании алгоритма градиентного бустинга деревьев решений

90

В проведенном исследовании были рассмотрены три программных реализации градиентного бустинга на деревьях решений: LightGBM, XGBoost и CatBoost.

LightGBM – распределенная, быстрая и высокопроизводительная реализация градиентного бустинга, основанная на алгоритмах деревьев решений и применяется для решения задач машинного обучения, таких как ранжирование данных, классификации и многих других. Данная платформа является разработкой компании Microsoft [92].

XGBoost – одна из самых популярных и эффективных реализаций алгоритма градиентного бустинга на деревьях, так же является очень гибкой и портативной [93]. Написанный код работает на большинстве существующих платформ, например, Hadoop, SGE, MPI и способен решать миллиарды примеров [94].

CatBoost – отечественная библиотека, представленная в 2017 г. Особенности алгоритма являются: построение симметричных деревьев, возможность работы с категориальными признаками, кроме того, он позволяет обучаться на относительно небольшом количестве неоднородных данных [95]. Сравнительный анализ ключевых возможностей представленных реализаций, представлен в таблице 6.

Таблица 6 – Сравнение возможностей описанных платформ

LightGBM	XGBoost	CatBoost
Скорость обучения модели высокая	высокая	низкая
Поддержка CPU/GPU да	да	да
Точность высокая	высокая	высокая
Поддержка параллельного вычисления	да	да
Обработка крупных массивов данных	да	да
Поддерживаемые языки программирования	Python, R	Comand Line, Python, R, Julia, JVM языки
	Comand Line,	Python, R

91

Стоит отметить, что недостатком всех классифицирующих методов является отсутствие уникального класса, типа «не похоже ни на один из классов», т.е. идентифицируемая сигнатура всегда будет принадлежать одному из классов, сформированных на обучающей выборке. Однако можно произвести дополнительную проверку имени предсказанного класса и имени исследуемого файла на электронном носителе информации. В случае их несовпадения, проверяющему стоит задуматься о наличии нелегитимной программы на исследуемом объекте.

3.4.3 Оценка вычислительной сложности различных методов сравнения

СИИФ с ЭСП

Для сравнения вычислительных сложностей двух описанных подходов – на основе статистических критериев и на основе машинного обучения, обратимся к тем этапам построения ЭСП и их сравнения с СИИФ, которые различаются. Так, статистический подход, в отличие машинного, дополняется этапом создания унифицированных сигнатур, который включает в себя этап кластеризации с применением метода k-средних и формирования для каждого кластера единого варианта ЭСП. Таким образом, при создании АЭСП вычислительная сложность формирования УСП составляет $O(nk)$, где n – число объектов, k – число кластеров, l – число итераций. Метод же на основе машинного обучения не требует создания УСП, отсюда можно сказать, что вычислительная сложность данного этапа постоянна $O(1)$ [96].

Далее, при рассмотрении этапа непосредственной идентификации исполняемых файлов, для подхода, основанного на статистических критериях, вычислительная сложность попарного сравнения всех сигнатур тестовой выборки с сигнатурами в АЭСП составляет $O(n')$, где n' – число объектов (унифицированных ЭСП) в архиве сигнатур, получившихся после этапа кластеризации, где очевидно следующее соотношение $n \geq n' \geq m$ (m – количество различных программ). Для подхода же на основе применения алгоритмов

92
машинного обучения – $O(n)$, где n остается неизменным в виду отсутствия этапа унификации сигнатур.

Из проведенного анализа следует, что для подхода на основе идентификации СИИФ с УСП посредством статистических методов, вычислительная сложность составляет – $O(nk) + O(n')$, тогда как для подхода идентификации СИИФ с ЭСП на основе машинного обучения – $O(1) + O(n)$.

3.5 Методика идентификации программного обеспечения на основе статических характеристик программного кода файла

На основе разработанных методов формирования и сравнения ЭСП и СИИФ предложена методика идентификации исполняемого файла, схема которой представлена на рисунке 12.

Методика содержит в себе обязательный подготовительный этап, на котором производится формирование АЭСП, включающем в себя ЭСП или унифицированные ЭСП. И этап идентификации исполняемого файла, состоящий из формирования СИИФ, их непосредственного сравнения с ЭСП из АЭСП, постобработку с применением аддитивного критерия и принятие итогового решения о подлинности предъявленного идентификатора.

В качестве решающего правила выступает выбор наибольшей вероятности принадлежности сигнатуры идентифицируемого исполняемого файла к эталонной сигнатуре программы из архива, по десяти ассемблерным командам:

$$\text{Name}(e_j) = \max \{ \text{rej}(N_i) = \sum \lambda_n \cdot \text{rej} \}$$

$$10 \text{ asp}(N_i)$$

$$n=1 \}$$

93
Рисунок 12 – Блок-схема методики идентификации исполняемых файлов

94
3.5.1 Предварительный этап методики идентификации программного обеспечения. Сбор и формирование обучающей выборки и составление

АЭСП, содержащий ЭСП или унифицированные ЭСП

Идентификация исполняемых файлов возможна только при наличии сформированного заранее архива (АЭСП), содержащего в себе ЭСП или унифицированные ЭСП, по причине того, что в процессе идентификации итерационно происходит сравнение последовательности измеренных характеристик идентифицируемого файла с такими же ранее собранными характеристиками версий программного обеспечения, хранящимися в АЭСП, для

которых известна их принадлежность к той или иной программе.

Первым шагом является сбор, обработка и структурирование версий программ. На данном этапе формируются множества N и P , а также различные версии $\{v_1, v_2, \dots, v_m\}$ по каждой $P_{\text{этал},i}$, представляющие собой обучающую выборку для того, чтобы в последствии имелась возможность отнесения идентифицируемого файла к одному из имеющихся наименований программ. Список всех рассматриваемых программ и их версий обязательно должен быть обновляемым и дополняемым на протяжении всего эксплуатационного периода защищаемой ИС.

Формирование ЭСП происходит путем дизассемблирования файлов обучающей выборки и подсчета частот встречаемости выбранного признака (рисунок 13). Которые в конце данного этапа формируют АЭСП.

Предлагается следующий алгоритм сбора версий программ и сохранения их ЭСП:

1) На основании действующих правил организации формируется эталонный список разрешенного или запрещенного ПО, регулярный пересмотр и обновление которого, является важным требованием.

Формирование такого списка может происходить несколькими путями, например, информация об эталонных программах может быть получена с электронных носителей информации, содержащих только легитимное ПО или

наоборот запрещенное; информацию можно получить, выгрузив дистрибутивы программ с официальных репозиториях ОС Linux и другими подходящими для задачи выделения определенного пула программ способами.

Автоматизация сбора программ легко реализуема путем написания скрипта на каком-либо языке программирования.

Рисунок 13 – Сбор характеристик файлов

2) После сбора ELF-файлов, необходимо произвести их обработку и структурирование.

Она заключается в формировании двух множеств: непосредственно программ $P = \{P_{\text{этал},1}, P_{\text{этал},2}, \dots, P_{\text{этал},i}\}$ с их различными версиями $P_{\text{этал},i} = \{v_1, v_2, \dots, v_m\}$ и множества имен этих программ $N = \{N_1, N_2, \dots, N_i\}$.

3) Следующий шаг состоит в формировании СП, для чего происходит дизассемблирование каждого ELF-файла и выделение из всего полученного кода последовательностей ассемблерных команд.

В зависимости от выбранного метода сбора характеристик файла и способа построения сигнатур на основе признакового пространства, производится подсчет встречаемости градации признака (ассемблерной команды) по каждому ELF-файлу в отдельности. После чего формируются промежуточные сигнатуры версий программ $S(v_m)$ и, на их основе, создаются СП $S(P_{\text{этал},i})$ сохраняемые в АЭСП. При этом, для методов идентификации, основанных на статистических критериях, производится дополнительный шаг, заключающийся в кластеризации СП $S(P_{\text{этал},i})$ и формировании унифицированных ЭСП, записываемых в АЭСП унифицированных сигнатур.

Стоит отметить, что в данном исследовании были выбраны десять ассемблерных команд в качестве признаков (для тех методов формирования сигнатур, которые основаны на подсчете частоты встречаемости градации признака на интервалах ассемблерного кода программ), анализ которых в совокупности позволяет достичь более высоких показателей точности идентификации, чем их использование по отдельности или всех вместе за раз. Использование аддитивного критерия Фишберна будет приведено в пункте 3.2.3.

3.5.2 Основные этапы методики идентификации исполняемого файла

Идентификация ELF-файлов при проведении мероприятий по аудиту ПО на электронных носителях информации состоит из двух подэтапов: формирования

СИИФ исследуемых файлов и процесса сравнения сигнатур и присвоения идентификатора (имени программы).

Алгоритм формирования СИИФ состоит в следующем:

97

- 1) Как и в предварительном этапе методики производится автоматический сбор всех ELF-файлов с анализируемого электронного носителя информации.
- 2) Формирование СИИФ происходит путем дизассемблирования файлов тестовой выборки и подсчета частот встречаемости выбранного признака (также, как и на рисунке 12).

Для каждого идентифицируемого файла выполняется итерационный процесс формирования СИИФ по всем десяти ассемблерным командам в точности, как и при формировании ЭСП. После чего все СИИФ формируют некоторую базу исследуемых сигнатур.

Алгоритм сравнения СИИФ с ЭСП из АЭСП состоит в следующем:

- 1) В зависимости от выбранного метода идентификации: при помощи статистического критерия или машинного обучения, происходит не только выбор способа формирования сигнатур, произведенный ранее в методике, но и, в соответствии с ним, способ сравнения сигнатур, а также выбор АЭСП.
- 2) Процесс сравнения сигнатур происходит путем попарного сравнения двух последовательностей признака – у идентифицируемого файла и у эталонной программы.

Так при формировании унифицированных ЭСП и сравнении их при помощи статистического критерия или нейронной сети, необходимо обратиться к АЭСП (унифицированных), при формировании ЭСП и сравнении их при помощи градиентного бустинга, необходимо обратиться к АЭСП.

Сравнение сигнатур для идентифицируемого файла также происходит по всем десяти ассемблерным командам, после чего следует этап постобработки результатов.

- 3) Постобработка результатов состоит в формировании конечного решения о принадлежности идентифицируемого файла к той или иной программе и заключается в применении коэффициентов важности аддитивного критерия Фишберна к результатам идентификации по десяти ассемблерным командам.

98

Коэффициенты должны быть рассчитаны заранее. Автором данные коэффициенты были рассчитаны при проведении эксперимента на результатах идентификации тестовой выборки.

В итоге, после произведения сравнения всех СИИФ с ЭСП, получается список имен программ, находящихся на исследуемом электронном носителе информации.

3.6 Оценка точности идентификации программного обеспечения

Оценка качества или эффективности методики, а также разработанных методов в частности, подразумевает оценку точности идентификации исполняемых файлов. В исследовании было использовано три подхода к оценке классификации: матрица ошибок (Confusion Matrix), точность (Accuracy), бикубическая мера (F-measure).

Оценка точности идентификации, на основе использования статистических критериев, оценивалась при помощи матрицы ошибок, позволяющей рассмотреть все возможные результаты, ее структура приведена в таблице 7.

Таблица 7 – Матрица ошибок (Confusion Matrix)

Генеральная совокупность

(Total population)

Реальное состояние

(True condition)

Положительное

состояние
(Condition positive)

Отрицательное

состояние
(Condition negative)

Прогнозируемое

состояние
(Predicted
condition)

Прогнозируемое состояние

положительное
(Predicted condition
positive)

Верно

положительное
(True positive)

Не верно

положительное
(False positive)

Прогнозируемое состояние

отрицательное
(Predicted condition
negative)

Не верно

отрицательное
(False negative)

Верно отрицательное

(True negative)

За точность идентификации Accuracy выступает отношение числа верно идентифицированных исполняемых файлов TP (true positive results) к общему объему исполняемых файлов, участвовавших в идентификации J (числу файлов

99

тестовой выборки). Показатель Accuracy считается по данной формуле и представляется в процентах:

Accuracy =

TP

J

· 100%.

Точность классификатора можно рассчитать при помощи F-measure, которая представляет собой среднее гармоничное между точностью (precision) и полнотой (recall). Так для мульти-классификатора необходимо произвести расчет F-measure по каждому объекту (программы тестовой выборки) по формуле:

Fmeasure(Prest.,i) =

1

α·

1

precision

+(1-α)·

1

recall

,

где коэффициент α позволяет взвешивать соотношение точности (precision) и полноты (recall), и принимает значение от нуля до единицы включительно.

Бикубическая мера для всех объектов рассчитывается как среднее значение

F-measure для каждого из i объектов:

Fmeasure =

$\sum_i Fmeasure(P_{тест,i})$

1

i

.

Разнообразие в подходе оценки точности вызвано следующими факторами:

для учета всех результатов в подходе идентификации на основе статистических критериев необходимо рассматривать матрицу ошибок (таблица 7), которая включает в себя все возможные результаты, получаемые при проведении статистического анализа;

для учета всех результатов в подходе идентификации на основе машинного обучения, применительно к задаче мульти-классификации, матрица ошибок уже не подходит, т.к. на выходе присутствуют результаты только двух видов – true positive (файл верно идентифицируется как программа) и false positive (файл не верно идентифицируется как программа), таким образом для оценки точности классификатора подходит вычисление Accisasy;

100

оценка бикубической меры производилась в целях сравнения получаемых результатов точности с существующими методами идентификации других исследователей.

3.7 Ограничения методики

Важно понимать, что методике идентификации ПО присущи ограничения, посредством влияния ряда факторов, таких как:

- потенциальное количество различных программ, установленных на одном компьютере или на всех компьютерах организации;
- количество различных версий для каждой программы, измеряемое от одной до десятков, но в практике не превосходящее сотни;
- наличие существенных изменений в различных версиях одной программы;
- наличие индивидуального характера распределения признака идентифицируемого файла от других программ;
- возможность производить дизассемблирование исполняемого файла;
- не способность в выявлении закладки вредоносного кода в легитимный исполняемый файл;
- классификация не способна выдать корректный результат для программы, класс которой не был определен на этапе формирования модели классификации (нельзя создать класс «файл не похож ни на одну из программ»);
- и многих других.

Выводы по главе 3

а) В третьей главе был произведен анализ структуры ELF-файла, его представление в ассемблерном коде и выявлены характеристики пригодные для использования при формировании сигнатур исполняемых файлов.

101

б) Представлен подход к формированию признакового пространства с выделением наиболее информативных признаков и их дальнейшего использования.

в) Сформирована модель представления программного обеспечения, заключающая в себе наименования и программные реализации кодов, версии программ и их внутренние характеристики, представляющие собой последовательность появления выбранного признака.

г) Разработаны методы по формированию ЭСП и унифицированных ЭСП, создание АЭСП, на основе которого происходит дальнейшая идентификация исполняемых файлов.

д) Предложено несколько методов по проведению сравнения нескольких сигнатур ЭСП и СИИФ, основанных как на проведении статистического анализа с

применением критерия однородности хи-квадрат и критерием согласия Колмогорова, так и с использованием алгоритмов машинного обучения, в частности искусственной нейронной сети и градиентного бустинга деревьев решений.

е) Представлен анализ вычислительной сложности методов сравнения сигнатур. Так подход на основе машинного обучения является менее затратным по вычислительным ресурсам и простым в исполнении.

ж) На основе разработанных методов формирования сигнатур и их сравнения была предложена методика идентификации исполняемых файлов.

з) Необходимо произвести ряд экспериментов, определяющих точность идентификации файлов с помощью предложенной методики. А также анализ методов, описанных в ней и выделение наиболее эффективного с точки зрения качества классификации.

102

Глава 4. Экспериментальная проверка точности идентификации программного обеспечения и анализ полученных результатов

В данной главе описывается серия проведенных экспериментов, целью которых являлась проверка достигаемой точности идентификации исполняемых файлов, получаемых при использовании разработанной методики, а также выявление наиболее результативного метода по формированию и сравнению сигнатур ELF-файлов.

Эксперименты производились на файлах, скаченных с официального репозитория Linux в 2016 году, формирующих из собой базу из 63 программ, 578 исполняемых файлов, более детальное описание было описано в разделе 2.2. Отсюда соотношение обучающей к тестовой выборке составляет 79% к 21%.

4.1 Входные данные экспериментов. Обучающая и тестовая выборки

База программ была разделена на две выборки: обучающую с 455 файлами и тестовую со 123 файлами. По разработанной методике идентификации, файлы были предварительно разделены по версиям различных программ, для обучающей выборки из них были выделены исследуемые характеристики, формирующие ЭСП или унифицированные ЭСП, и помещены в АЭСП.

Помимо того, что одной из ключевых особенностей рассматриваемой задачи является ограниченное число идентифицируемых файлов, с этим связана еще проблема ограниченного числа версий по каждой программе. Реализуемая методика идентификации должна позволять производить идентификацию исполняемых файлов в условиях дисбаланса числа версий для различных программ.

Для приближения экспериментов к реальной ситуации, обучающая выборка была составлена таким образом, что для различных программ число используемых различных версий для них было неравномерно.

В представленных ниже экспериментах, оценка точности производилась в соответствии с описанными критериями точности в разделе 2.6.

103

4.2 Определение точности идентификации при использовании разработанных методов сравнения СИИФ с ЭСП

Цель эксперимента – оценка сравнительного анализа точности идентификации ПО с использованием разработанных методов по формированию сигнатур исполняемых файлов и их сравнение с различными существующими методами, а также с показателями точности идентификации, полученными авторами других работ.

В ходе теоретического исследования были разработаны несколько подходов к сбору характеристик исполняемых файлов и их сравнению, чьи показатели точности были исследованы ниже.

полученные результаты. А также, для наглядности, рассмотрим на примерах описанные ранее в разделе 2.4 методы сравнения СИИФ с ЭСП.

4.2.1 Точность идентификации при использовании статистических методов сравнения сигнатур

Подход к построению унифицированных ЭСП и СИИФ на основе распределения 118 ассемблерных команд и идентификации при помощи критерия однородности хи-квадрат.

Пусть идентифицируемым файлом является программа `bacula-console` версии `5.2.5-0ubuntu6`. Сравнение выбранного файла будет происходить с программами: `bacula-console` и `baobab`. Таким образом, были сформированы сигнатура идентифицируемого файла и заранее на обучающей выборке две сигнатуры выбранных программ.

На рисунках 14 и 15 показаны распределения частот признака в сигнатуре идентифицируемого файла `bacula-console_5.2.5-0ubuntu6` и соответственно в сигнатурах программ `bacula-console` и `baobab` из АЭСП. По оси абсцисс отложены градации признака с ненулевыми значениями в сигнатуре, по оси ординат – частоты ненулевых ассемблерных команд из сигнатуры программы.

104

Рисунок 14 – Сигнатуры идентифицируемого файла является `bacula-console_5.2.5-0ubuntu6` и программы `bacula-console`

При проверке статистической гипотезы о сходстве распределения частот градаций признака в сигнатуре идентифицируемого файла `bacula-console_5.2.5-0ubuntu6` и в сигнатуре программы `bacula-console_i386` из АЭСП эмпирическое значение критерия однородности $\chi^2 = 6,76$, что ниже критического значения $\chi^2_{p=0,05} = 26,30$, поэтому на уровне значимости $p = 0,05$ можно утверждать, что идентифицируемый файл является программой из архива (что действительно так). Действительно, на рисунке 14 видно, что гистограмма частот встречаемости градаций признака идентифицируемого файла `bacula-console_5.2.5-0ubuntu6` и кривая распределения частот для программы `bacula-console` различаются незначительно.

В свою очередь при проверке статистической гипотезы о сходстве распределения частот градаций признака в сигнатуре идентифицируемого файла `bacula-console_5.2.5-0ubuntu6` и в сигнатуре программы `baobab_i386` из АЭСП получено эмпирическое значение критерия однородности $\chi^2 = 301,903$, превышающее критическое значение $\chi^2_{p=0,05} = 35,17$ на уровне значимости $p = 0,05$. Следовательно, на данном уровне значимости можно утверждать, что

105

идентифицируемый файл не является программой из архива (что действительно так).

Рисунок 15 – Сигнатуры идентифицируемого файла является `bacula-console_5.2.5-0ubuntu6` и программы `baobab`

Полученный вывод подтверждается на рисунке 15, где графически показано, что различия в распределениях частот встречаемости градаций признака идентифицируемого файла `bacula-console_5.2.5-0ubuntu6` и программы `baobab_i386` из архива значительны.

Общие результаты эксперимента с применением критерия однородности хи-квадрат для распределения 118 ассемблерных команд представлены в матрице ошибок (таблица 8), где в первой и четвертой строках содержится показатель (в процентах) корректных результатов, в третьей и четвертой строках – показатель ошибки первого и второго рода соответственно.

В первом и втором столбцах таблицы описываются все возможные исходы идентификации сигнатур, а в третьем столбце отображены полученные результаты идентификации для этих исходов, выраженные в процентных долях.

Таблица 8 – Матрица ошибок для подхода на основе критерия однородности хи-

квадрат

106

Отношение между ЭСП из

АЭСП и СИИФ

Основная гипотеза

о сходстве

сигнатур

Результат эксперимента, %

Уровень

значимости

$p = 0,05$

Уровень

значимости

$p = 0,01$

Сигнатуры относятся к одной и

той же программе

Принимается 0,36 0,42

Отвергается 1,23 1,17

Сигнатуры относятся к разным

программам

Принимается 0,10 0,16

Отвергается 98,32 98,25

Из таблицы 8 следует, что показатель точности на основе Confusion Matrix

для идентификации сигнатур ELF-файлов, построенных на основе частот 118

ассемблерных команд, составляет 98,68% для уровня значимости $p = 0,05$ и 98,67%

для $p = 0,01$.

Подход к построению унифицированных ЭСП и СИИФ на основе

распределения одной ассемблерной команды на равных интервалах и

идентификации при помощи критерия однородности хи-квадрат.

Пусть идентифицируемым файлом является программа amarak версии 2.3.0-

0ubuntu4. Сравнение выбранного файла будет происходить с программами: amarak

и apasorn. Таким образом, были сформированы сигнатура идентифицируемого

файла и заранее на обучающей выборке две сигнатуры выбранных программ.

На рисунке 16 изображены контурные графики, где для построения

трехмерной поверхности использовалась абсолютная разница между десятью

распределениями ассемблерных команд (сигнатурами СИИФ и ЭСП),

построенными для идентифицируемого файла amarak_2.3.0-0ubuntu4 и десятью

унифицированными сигнатурами из АЭСП для двух программ apasorn и amarak.

По оси абсцисс отмечены номера интервалов разбиения ассемблерной

последовательности эталонной программы k от 1 до 30, а по оси ординат –

порядковый номер ассемблерной команды, используемой для построения

сигнатур.

107

а) б)

Рисунок 16 – Абсолютная разница между СИИФ amarak_2.3.0-0ubuntu4

и унифицированных ЭСП а) amarak (левый) и б) apasorn (правый)

Видно, что для сигнатур, относящихся к одной и той же программе (рисунок

16 а)), высота пиков не превышает значения 150, а для многих ассемблерных

команд и вовсе минимальна, у различных же программ (рисунок 16 б)) поверхность

имеет большее число возвышений, а их максимальное значение почти в два раза

больше 139, чем 139 для одинаковых программ.

На рисунке 17 представлено влияние изменения значения коэффициента

(ratio) в формуле (3) на результат идентификации. Эксперимент проводился для

ассемблерной команды call, а значение коэффициента (ratio) изменялось от 0,3 до 1

с шагом 0,05.

На графиках рисунка 17 а) изображены корректные результаты идентификации, т.е. когда верно принята основная гипотеза H_0 (верхний) или верно принята альтернативная гипотеза H_1 (нижний). На графике рисунка 17 б) изображены результаты возникающих ошибок первого и второго рода, когда неверно отвергнута основная гипотеза H_0 и когда неверно принята основная гипотеза H_0 .

108

а) б)

Рисунок 17 – Зависимость результатов идентификации от коэффициента (ratio)

Как видно, оптимальным значением коэффициента (ratio) является 0,6, именно в этой точке график ошибки второго рода замедляет свой рост.

Общие результаты эксперимента с применением критерия однородности хи-квадрат для распределения одной из десяти ассемблерных команд на равных интервалах представлены в матрице ошибок (таблица 9). Принятие или отклонение основной гипотезы H_0 принималось на уровне значимости $p = 0,05$.

Таблица 9 – Матрица ошибок для подхода на основе критерия однородности хи-квадрат

Отношение между

ЭСП из АЭСП и

СИИФ

Сигнатуры относятся к одной и

той же программе

Сигнатуры относятся к разным

программам

Основная гипотеза о

сходстве сигнатур

Принимается

Отвергается

Принимается

Отвергается

Результат эксперимента, %

Ассемблерная команда

add 0,16 1,44 0,01 98,4

and 0,3 1,29 1,37 97,04

call 0,34 1,25 2,19 96,22

cmp 0,44 1,15 6,34 92,07

je 0,45 1,14 3,43 94,98

jmp 0,33 1,26 1,15 97,26

lea 0,47 1,12 11,44 86,97

mov 0,22 1,37 0,69 97,71

pop 0,33 1,26 2,34 96,07

push 0,26 1,33 0,74 97,67

109

Из таблицы 9 следует, что наибольший показатель точности на основе Confusion Matrix для идентификации сигнатур ELF-файлов, построенных на основе частот выбранной ассемблерной команды, составляет 98,56% для ассемблерной команды add.

Подход к построению унифицированных ЭСП и СИИФ на основе распределения одной ассемблерной команды на неравных интервалах и идентификации при помощи критерия согласия Колмогорова.

В качестве признака для формирования частотных распределений была выбрана ассемблерная команда cmp.

Пусть идентифицируемым файлом является программа bacula-console версии

5.2.5-0ubuntu6. Сравнение выбранного файла будет происходить с программами:

bacula-console и baobab. Таким образом, были сформированы сигнатура идентифицируемого файла и заранее на обучающей выборке две сигнатуры выбранных программ.

На рисунках 18 и 19 изображены гистограммы накоплений относительных частот для сигнатуры идентифицируемого файла bacula-console_5.2.5-0ubuntu6 и кривые соответственно для сигнатур программ bacula-console и baobab из архива.

По оси абсцисс отложена последовательность интервалов разбиения дизассемблированного кода программ с ненулевой частотой встречаемости признака, по оси ординат – накопления частот встречаемости ассемблерной команды str на этих интервалах.

Как видно из рисунка 18, накопление относительных частот в сигнатуре идентифицируемого файла bacula-console_5.2.5-0ubuntu6 незначительно отклоняется от накопления относительных частот в сигнатуре программы bacula-console и накопления для равномерного закона распределения. Максимальное расхождение составляет d_n

$\ast = 0,073$, которое меньше критического значения

$d_{p;n} = 0,106$, поэтому на уровне значимости $p = 0,05$ можно утверждать, что идентифицируемый файл является программой из АЭСП (что действительно так).

110

Рисунок 18 – Сигнатуры идентифицируемого файла является bacula-console_5.2.5-0ubuntu6 и программы bacula-console

Рисунок 19 – Сигнатуры идентифицируемого файла является bacula-console_5.2.5-0ubuntu6 и программы baobab

Из рисунка 19 следует, что расхождения между накоплениями относительных частот в сигнатуре идентифицируемого файла bacula-console_5.2.5-0ubuntu6 и в сигнатуре программы baobab_i386 значительны, при этом максимальное расхождение составляет d_n

$\ast = 0,362$, которое больше критического

значения $d_{p;n} = 0,106$, поэтому на уровне значимости $p = 0,05$ можно утверждать,

111

что идентифицируемый файл не является программой из АЭСП (что действительно так).

Общие результаты эксперимента с применением критерия согласия Колмогорова представлены в матрице ошибок (таблица 10).

Таблица 10 – Матрица ошибок для подхода на основе критерия согласия

Колмогорова

Отношение между ЭСП из

АЭСП и СИИФ

Основная

гипотеза

о сходстве

сигнатур

Результат эксперимента, %

Уровень

значимости

$p = 0,05$

Уровень

значимости

$p = 0,01$

Сигнатуры относятся к одной и

той же программе

Принимается 1,86 1,94

Отвергается 0,61 0,53

Сигнатуры относятся к разным

программам

Принимается 4,30 6,39

Отвергается 93,23 91,14

Из таблицы 10 следует, что показатель Accuracy на основе Confusion Matrix для идентификации сигнатур ELF-файлов, построенных на основе частот выбранной ассемблерной команды, составляет 95,09% для уровня значимости $p = 0,05$ и 93,08% для $p = 0,01$.

4.2.2 Точность идентификации при использовании метода сравнения сигнатур на основе машинного обучения

Подход к построению унифицированных ЭСП и СИИФ на основе распределения одной ассемблерной команды на равных интервалах и идентификации при помощи искусственной нейронной сети.

Эксперимент проводился для ассемблерной команды jmp и с использованием программного средства STATISTICA, которая позволяет настроить пользовательскую нейронную сеть и предоставляет удобный интерфейс взаимодействия с пользователем.

Обучение и тестирование нейронной сети предполагает определенное формирование входных данных. На рисунке 20 представлен скриншот окна программного средства STATISTICA, где в первом столбце приведены числовые

112

идентификаторы имен различных программ (их соответствие представлено в приложении 2), далее следуют унифицированные ЭСП встречаемости признака (одной ассемблерной команды на 30 равных интервалах разбиения), затем СИИФ, а в последнем столбце стоит отметка о принадлежности каждой из сигнатур к обучающей или тестовой выборке.

Рисунок 20 – Инициализация искусственной нейронной сети в программе STATISTICA

В пользовательской нейронной сети предлагается выбор типа архитектуры нейронной сети, число нейронов и итераций алгоритм обучения и другое.

По результатам ряда экспериментов было установлено, что наиболее точные показания идентификации достигаются при установке следующих параметров обучения искусственной нейронной сети (hbseujr 21):

тип сети: многослойный перцептрон (MLP);

алгоритм обучения: метод сопряженных градиентов (Conjugate gradient);

число итераций: 1000;

число скрытых нейронов: 50.

113

Рисунок 21 – Задание параметров нейронной сети

В ходе эксперимента было построено 5 моделей нейронных сетей (Модель 1- Модель 5). Результаты правильной классификации которых показаны в таблице 11.

Таблица 11 – Показатели Accuracy для пяти моделей нейронной сети по ассемблерной команде jmp

Accuracy

Модель 1 75,61

Модель 2 76,42

Модель 3 64,23

Модель 4 74,80

Модель 5 78,86

Среднее значение Accuracy 73,98

На рисунке 22 представлен график, отображающий распределение результатов произведенной классификации СИИФ из тестовой выборки. На оси абсцисс отмечены спрогнозированные классы (имена программ) для идентифицируемых файлов, на оси ординат – реальные классы (имена программ)

этих файлов. Так, в случае совпадения прогнозируемого результата с реальным, на графике маркер ставится в точке с координатой (x, y), где x = y, т.е. на диагонали.

114

Рисунок 22 – График зависимости результатов идентификации от истинной принадлежности файлов классам

Общие результаты эксперимента оценивались при помощи точности

Ассигасу по десяти ассемблерным командам с применением искусственной нейронной сети представлены в таблице 12.

Таблица 12 – Ассигасу для подхода на основе искусственной нейронной сети

Ассемблерная команда Среднее значение Ассигасу по пяти моделям

add 76,42

and 69,11

call 74,8

cmp 82,11

je 78,05

jmp 73,98

lea 56,1

mov 74,8

pop 69,1

push 53,66

Из таблицы 12 следует, что наибольший показатель Ассигасу для идентификации сигнатур ELF-файлов, построенных на основе частот выбранной ассемблерной команды, составляет 82,11% для ассемблерной команды `cmp`.

115

Подход к построению ЭСП и СИИФ на основе распределения одной ассемблерной команды на равных интервалах и идентификации при помощи градиентного бустинга деревьев решений.

Для данного подхода были выбраны три библиотеки, реализующие алгоритм градиентного бустинга деревьев решений LightGBM, CatBoost и XGBoost.

Формирование обучаемой модели для решения задачи мульти-классификации [89] происходило по нескольким параметрам, значения для которых представлены в таблице 13. Остальные значения параметров моделей были установлены по умолчанию.

Таблица 13 – Параметры обучения моделей

Параметры обучения

Описание параметров XGBoost LightGBM CatBoost

booster/

boosting/

boosting_type

тип бустинга gbtree gbdt plain

eta/

learning_rate

скорость обучения,

используемая для

уменьшения шага

градиентного спуска

0,3 0,3 0,7

-/- l2_leaf_reg

коэффициент

регуляризации L2,

используемый для

расчета значения

листов

-- 1

max_depth/
depth
глубина дерева 2 7 2
num_round/
num_iterations/

iterations
число итераций

бустинга
1000 1000 1000

objective/
loss_function

метрика оценки
(функция потерь),
используемая в
обучении модели

Multi:SoftMax MultiClass MultiClass

Общие результаты эксперимента оценивались при помощи точности

Accuracy по десяти ассемблерным командам с применением градиентного бустинга
деревьев решений представлены в таблице 14.

116

Таблица 14 – Accuracy для подхода на основе градиентного бустинга деревьев
решений

Ассемблерная
команда

Accuracy для LightGBM Accuracy для CatBoost Accuracy для XGBoost

add 82,93 85,37 92,68

and 79,67 86,18 92,68

call 81,30 84,55 94,31

cmp 78,05 85,37 84,55

je 86,99 89,43 91,06

jmp 84,55 83,74 95,93

lea 74,80 76,42 91,87

mov 74,80 85,37 94,31

pop 78,86 83,74 88,62

push 74,80 82,93 89,43

Очевидно, что для почти всех ассемблерных команд число верно
идентифицированных программ выше при использовании библиотеки XGBoost и
достигает максимального значения в 95,93% для команды jmp (118 правильно
идентифицированных исполняемых файла из 123).

При этом из рисунка 23 видно значительное преимущество LightGBM во
времени, затрачиваемого на обучение модели и идентификацию всех исполняемых
файлов тестовой выборки.

Рисунок 23 – Отношение времени обучения модели классификации и
идентификации файлов тестовой выборки

XGBoost

0,09

LightGBM

0,03

CatBoost

0,88

117

4.2.3 Повышение точности идентификации путем использования

аддитивного критерия Фишберна

В задачах, когда эффективность принимаемых решений оценивается не

одним, а несколькими критериями, является центральной проблемой теории принятия решений. В работе [97] даются строгие определения равенства и неравенства критериев по важности. Классификация методов определения коэффициентов важности критериев представлена в работе [98], среди которых выделяется две группы методов: первые – определяют коэффициенты важности критериев, использование которых возможно в обобщенных свертках; и вторые – определяют весовые коэффициенты, использование которых в обобщенных свертках не рекомендуется.

Аддитивный критерий Фишберна [83], [100] как раз относится к первой группе методов, а в частности к методам аддитивной свертки. Данный критерий обладает рядом достоинств, в частности, отсутствует необходимость в опросе мнений экспертов, отсутствуют ограничения на условия реализации, простота расчетов и т.д.

Методы аддитивной свертки можно использовать, если функция полезности $\Phi(f(x))$ представима в аддитивной форме:

$$\Phi(f_1(x), \dots, f_n(x)) = \sum_{i=1}^n \lambda_i \cdot \Phi_i$$

n

$i=1$

$\cdot \Phi_i(x)$

где λ_i – коэффициенты относительной важности критериев, которые определяются по формуле:

$$\lambda_i =$$

$$2 \cdot (n - i + 1)$$

$$n \cdot (n + 1)$$

$$, i = 1, \dots, n$$

при упорядоченных критериях $f_1(x) \geq \dots \geq f_n(x)$.

При этом, если значения нескольких критериев

равны $f_i(x) = f_{i+1}(x) = \dots = f_{i+m}(x)$, тогда и коэффициенты относительной важности

для них одинаковы и равны среднему значению

$$1$$

$$\sum_{j=1}^{i+m} \lambda_j$$

$$\cdot \sum_{j=1}^{i+m} \lambda_j$$

$$i+m$$

$$j=1, \text{ как если бы они были}$$

118

рассчитаны для не равных критериев с сохранением их порядка

$$f_i(x) > f_{i+1}(x) > \dots > f_{i+m}(x).$$

По результатам проведенных экспериментов можно сделать вывод о том, что **139**

наиболее **139** точные результаты идентификации были получены при построении

сигнатур, основанных на частоте одной ассемблерной команде на равных

интервалах разбиения, с использованием метода их сравнения на основе

градиентного бустинга деревьев решений в реализации библиотеки XGBoost. В

таблице 13 представлены результаты по десяти ассемблерным командам

выделенные цветовой шкалой, так наибольшая точность идентификации

исполняемых файлов достигается при использовании ассемблерной команды `jmp`,

а наихудшая при – `str`.

Таблица 15 – Упорядоченные цветом ассемблерные команды по достигаемой

точности идентификации

`xgboost`

`add and call cmp je jmp lea mov pop push`

113 113 115 104 112 118 112 116 109 110

Произведем расчет коэффициентов относительной важности критериев

(ассемблерных команд) для нашей задачи.

`jmp mov call add and`

$\lambda_1 =$

2·(10-1+1)
10·(10+1)
λ2=
2·(10-2+1)
10·(10+1)
λ3=
2·(10-3+1)
10·(10+1)
λ4=
2·(10-4+1)
10·(10+1)
λ5=
2·(10-5+1)
10·(10+1)
λ1=0,182 λ2=0,164 λ3=0,145 λ4=0,118 λ5=0,118

je lea push pop cmp

λ6=
2·(10-6+1)
10·(10+1)
λ7=
2·(10-7+1)
10·(10+1)
λ8=
2·(10-8+1)
10·(10+1)
λ9=
2·(10-9+1)
10·(10+1)
λ10=
2·(10-10+1)
10·(10+1)
λ6=0,082 λ7=0,082 λ8=0,055 λ9=0,036 λ10=0,018

119

Таким образом, мы получаем аддитивную функцию вида:

$\Phi(\text{jmp}, \dots, \text{cmp}) = 0,182 \cdot \text{jmp} + 0,164 \cdot \text{mov} + 0,145 \cdot \text{call} + 0,118 \cdot \text{add} + 0,118 \cdot \text{and} + 0,082 \cdot \text{je} + 0,082 \cdot \text{lea} + 0,055 \cdot \text{push} + 0,036 \cdot \text{pop} + 0,018 \cdot \text{cmp}$

применив которую к результатам классификации XGBoost получим вероятности принадлежности исполняемых файлов тестовой выборки к определенным программам на основании совокупности десяти ассемблерных команд и коэффициентов их важности. Так на рисунке 24 отображены результаты идентификации с постобработкой результатов классификации XGBoost по всем 123 исполняемым файлам тестовой выборки, отмеченным на оси ординат, по оси абсцисс отмечается вероятность принадлежности исполняемого файла к классу (программе).

а) б)

Рисунок 24 – Результаты постобработки результатов классификации

а) для всех файлов; б) для 20 файлов для ближайшего рассмотрения

Автором, для более понятной визуализации, специально была выбрана линейчатая диаграмма с накоплением, где первым прямоугольником отображается вероятность принадлежности исследуемого файла к истинной программе (которой он и является). Далее следуют вероятности принадлежности к другим, не истинным, программам. Так, после применения аддитивного критерия Фишберна,

120

точность идентификации исполняемых файлов возросла до 99,19% (только один файл из 123 был неправильно идентифицирован) [76].

4.3 Определение итоговой точности идентификации программного

обеспечения на основе предложенной методики

Цель эксперимента – определить точность идентификации исполняемых файлов по характеристикам их ассемблерного кода с использованием разработанной методики в условиях ограниченного количества версий и программ разнообразной функциональности, а также провести сравнение с результатами, других исследовательских работ.

На унифицированных сигнатурах программ наилучшим методом идентификации, при подходе на основе оценки матрицы ошибок (Confusion Matrix) и уровне значимости $p = 0,05$, показал себя критерий однородности хи-квадрат с построением сигнатур на основе 118 ассемблерных команд, чья точность достигла 98,68%, что выше на 0,12% и на 3,59%, чем при использовании критерия однородности хи-квадрат и согласия Колмогорова, соответственно, с построением сигнатур на основе частоты встречаемости одной ассемблерной команды.

Однако, при расчете точности при помощи Ассигасу, достигаемые результаты равны лишь 79,55%, 78,05% и 50,41% соответственно. Точность же с использованием искусственной нейронной сети достигает 82,11% и является наиболее эффективной с точки зрения числа верно идентифицированных исполняемых файлов, при построении унифицированных сигнатур.

На не унифицированных сигнатурах методы градиентного бустинга деревьев решений, при подходе к оценке через Ассигасу, показали наиболее высокую достигаемую точность со средним показателем равным 84,98%, чем с использованием статистических критериев однородности хи-квадрат и Колмогорова, а также с использованием нейронной сети, чей средний показатель равен 66,07%. Среди подходов к реализации градиентного бустинга, наиболее результативным является XGBoost со средним показателем точности 91,29%.

121

Среди ассемблерных команд наилучший результат получается при использовании в качестве признака команды je с 82,93%, наихудший при использовании push с 71,54%. Так разница между достигаемыми показателями точности для этих команд составляет 11,39%.

Максимально достигаемая величина значения точности идентификации, через оценку Ассигасу, равняется 95,95% и достигается за счет применения метода идентификации на основе градиентного бустинга деревьев решений в реализации XGBoost к сигнатурам, сформированным на основе распределения одной ассемблерной команды jmp на равных интервалах дизассемблированного кода программы.

Таблица 16 – Показатель Ассигасу по каждой ассемблерной команде, в процентах и F-measure, (в скобках)

Ассемблерные команды

Статистические критерии,

$p = 0,05$

Машинное обучение

однородности хи-

квадрат

согласия

Колмогорова

однородности хи-

квадрат

Нейронная сеть

MLP

Градиентный бустинг деревьев

решений

LightGBM

CatBoost

XGBoost
add
79,55
- 40,65 76,42 82,92 85,37 92,68
and - 60,16 69,11 78,87 86,18 91,87
call - 65,85 74,80 81,3 84,55 93,50
cmp 50,41 70,73 82,11 78,04 85,37 84,55

je - 69,11 78,05

86,99

(0,84)

89,43

(0,88)

91,05

jmp - 65,85 73,98 83,74 83,74

95,93

(0,96)

lea - 78,05 56,10 73,99 76,42 91,06

mov - 47,97 74,80 73,99 85,37 94,30

pop - 60,16 69,11 78,86 83,74 88,60

push - 56,91 53,66 74,79 82,93 89,40

Фишберн --- 84,93 87,81 91,87

99,19

(0,99)

122

После применения подхода постобработки результатов, на основе аддитивного критерия Фишберна и сочетания десяти ассемблерных команд, точность идентификации исполняемых файлов решений в реализации XGBoost возрастает на 3,26% и достигает показателя 99,19%.

Оценка точности по различным ассемблерным командам приведена в таблице 16. Ассигасу рассчитана для каждой ассемблерной команды, F-measure только для наилучших результатов.

Оценка точности идентификации на основе Ассигасу и F-measure, позволяет сделать вывод о том, что [139] разработанная методика идентификации ELF-файлов позволяет получить более высокий показатель точности с использованием всех рассмотренных методов классификации в сравнении с существующими методами идентификации программ, описанными в рассмотренных работах отечественных и зарубежных авторов [139]. Так разработанная автором методика, не только позволяет решить совершенно другую область распознавания ПО, но и превосходит другие наиболее результативные подходы (таблица 17).

Таблица 17 – Точность идентификации ПО с использованием разработанной методики и различных современных методов

Признаковое пространство и метод

идентификации

Число классов

Оценка качества метода

Точность

(Ассигасу),

%

Бикубическая

мера качества

кластеризации

Печатаемые строки + Naive Bayes, [35]

Бинарная

классификация

97,11 -

Последовательность частоты

встречаемости очередности ($n = 2$)

ассемблерных команд + SVM:

normalised polynomial, [42]

Бинарная

классификация

95,9 -

Вектора для блоков постоянного

размера побайтового кода программы

+ Редакционное расстояние, [30]

Мульти-

классификация

- 0,69

Последовательность частоты

встречаемости одной ассемблерной

команды на равных интервалах +

XGBoost, [75]

Мульти-

классификация

95,93 0,96

Последовательность частоты

встречаемости одной ассемблерной

команды на равных интервалах +

XGBoost + Фишберн, [76]

Мульти-

классификация

99,19 0,99

123

4.4 Использование результатов исследования для повышения

информационной безопасности автоматизированных систем

Методика идентификации ПО может быть применена для повышения

безопасности информационных систем организации, путем внедрения ее в качестве

технической меры по аудиту ПО на электронных носителях информации.

Контроль установленных программ представляет собой важную

составляющую в организации защищенного функционирования бизнес-процессов.

Методика позволяет осуществлять автоматизированную идентификацию

исполняемых файлов формата ELF в соответствии с формируемых архивом

легитимных или не легитимных программ.

До применения методики реализация угрозы со стороны злоумышленника

может быть произведена, как показано на рисунке 25 [101].

Рисунок 25 – Схема возможного снижения уровня безопасности

системы и воздействие злоумышленника на ИС организации

Группа пользователей АС взаимодействует с СВТ и другими электронными

носителями информации, присутствующими в организации и соединенными с

124

главным сервером (или серверами, хранящими и обрабатывающими защищаемую

информацию). В условиях соблюдения установленной политики безопасности,

службе защиты информации известны присущие ИС уязвимости и сформирована

стратегия по минимизации возможных рисков при инцидентах ИБ. Однако, при

появлении несанкционированно установленного ПО, присутствует вероятность

возникновения новой уязвимости в системе, существование которой не будет

известно службе безопасности до тех пор, пока не будет **139** проведен повторный

мониторинг АС для выявления уязвимостей.

В этот промежуток времени злоумышленник, обладая необходимыми

знаниями и средствами, способен произвести атаку на сервер организации путем эксплуатации возникшей уязвимости из-за несанкционированной установки ПО.

Рисунок 26 – Схема возможного снижения уровня безопасности системы и воздействие злоумышленника на ИС организации, включающая дополнительную меру по аудиту ПО

Применение разработанной методики в ИС организации способна предотвратить описанный сценарий путем проведения автоматизированного

125

аудита установленного ПО на СВТ и электронных носителях информации (рисунок 26).

Так располагая АЭСП на стороне сервера и проводя сравнения СИИФ всех присутствующих программ с каждого СВТ и подключаемых к ним съемных электронных носителей информации можно в короткие сроки уведомить службу ИБ о нарушении правил безопасности и принять своевременные защитные меры.

Разработанные методы и методика могут быть применены для выявления факта нарушения установленных организационных мер о запрете на установку определенных программ путем распознавания установленного исполняемого файла как ту или иную программу в АЭСП или не входящую в данный архив.

Можно выделить смежные задачи, решаемые разработанной методикой:

выявление не добросовестных пользователей, многократно устанавливающих запрещенное ПО;

выявление пользователей, пытающихся обойти меры по обеспечению ИБ;

анализ конкретной программы, на основе версий других, схожих по функционалу программ, для выявления степени принадлежности их области назначения.

Стоит отметить, что методика эффективно показывает себя на ограниченном наборе данных и способна идентифицировать исполняемый файл, не задействованный при создании АЭСП, имеющий новую версию или внесенные в него изменения, а также имеющий исправленные, или не имеющий их вообще, метаданные.

Методика может быть реализована в программном комплексе и использоваться как при периодически проводимых мероприятиях, так и быть установлена на компьютере пользователя для проведения идентификации всего устанавливаемого ПО в автоматическом режиме.

Методика позволяет производить идентификацию непосредственно исполняемого файла, а не его метаданных, записанных в конфигурационных файлах, и, таким образом, может быть использована в компьютерной криминалистике специальными службами.

126

Выводы по главе 4

а) Для выполнения цели экспериментов при подготовки настоящего диссертационного исследования был произведен сбор и анализ экспериментальных данных, включающий разнообразные виды ПО Linux ОС и их различных версий.

б) Анализ разработанных методов идентификации ПО позволил выявить наилучший из них, основанный на формировании сигнатур программ с подсчетом частоты встречаемости выбранных десяти ассемблерных команд на равных интервалах разбиения дизассемблированного кода исполняемых файлов и использующий алгоритм градиентного бустинга деревьев решений с последующей постобработкой получаемых результатов классификации.

в) Сравнение точности производимой идентификации исполняемых файлов при использовании разработанного метода по формированию ЭСП и СИИФ с наилучшими результатами, получаемыми другими исследователями, позволяет сделать вывод о наличии **139** более высокой точности идентификации исполняемых файлов перед методами, основанными на бинарной классификации на 2,08% измеренными в Ассигасу, и перед методами, основанными на мульти-

классификации на 18% измеренными в F-measure.

г) Подход на основе машинного обучения показал наиболее низкую вычислительную сложность, перед подходом с использованием статистических критериев, что подтверждается проведенными экспериментами.

д) Анализ результатов применения градиентного бустинга деревьев решений на сформированных ЭСП и СИИФ позволяет сделать вывод о том, что данный класс алгоритмов машинного обучения способен достигать наиболее точных результатов идентификации ПО. Применение методики идентификации с данным алгоритмом позволяет в среднем достигать на 4,10% более точные результаты, рассчитанные при помощи Assurasu, и на 11,42% более точные результаты, рассчитанные при помощи Assurasu при использовании библиотеки XGBoost.

127

е) Постобработка результатов классификации для алгоритмов машинного обучения показывает, что использование аддитивного критерия Фишберна с расчетом коэффициентов важности для десяти ассемблерных команд, позволяет в среднем повысить результаты классификации на 4,17%, а для библиотеки XGBoost достигать почти идеальной идентификации с Assurasu равным 99,19%.

128

Заключение

В диссертационной работе решена научная задача увеличения точности идентификации исполняемых файлов, устанавливаемых на средства вычислительной техники, в условиях наличия различных версий ПО, большого числа различных наименований программ и ограниченности числа объектов обучающей выборки, обусловленных реальным состоянием выпускаемых версий того или иного ПО. Использование разработанной методики позволяет увеличить точность идентификации исполняемых файлов, что предоставляет возможность производить более тщательный анализ электронных носителей информации на предмет наличия на них нелегитимно установленного ПО. Методика реализуется как часть технических мер по обеспечению ИБ и предотвращению потенциального возникновения новых уязвимостей ИС, способных повлечь за собой негативные воздействия на защищаемую информацию и персонал АС. В том числе получены следующие научные результаты, составляющие итоги исследования:

- 1) Произведено исследование и анализ существующих подходов и методов идентификации программного обеспечения, используемых отечественными и зарубежными исследователями, выявлены достоинства и недостатки данных методов, а также проанализированы условия и ограничения их использования.
- 2) Рассмотрены ELF формат файлов Linux систем, представление программы в виде ассемблерного кода и распределение ассемблерных команд. Проанализированы характеристики ассемблерных команд, их частота и информативность, а также потенциал их применения при формировании сигнатур исполняемых файлов и их последующего сравнения.
- 3) Разработан метод представления идентификатора исполняемого файла в виде сигнатуры, позволяющий реализовывать гибкий подход к идентификации версий ПО.
- 4) Разработан метод сравнения сигнатур программ на основе использования библиотеки градиентного бустинга деревьев решений с последующей постобработкой результатов классификации, на основе аддитивного критерия.

129

5) Разработана методика идентификации исполняемых файлов, позволяющая повысить точность идентификации в условиях большого числа классов и малого числа элементов обучающей выборки, соответствующих реальным условиям эксплуатации ПО в организации. Достигается значение Assurasu в 99,19%, что на 2,08% превышает бинарный подход классификации на

основе использования UNIX команды strings, показывающая печатаемые строки в объекте или двоичном файле, и сравнения получаемых сигнатур при помощи наивного Байесовского классификатора. Достижимые результаты также превышают измеренный показатель F-measure на 30% чем у мульти-классификационного подхода на основе использования вектора для блоков постоянного размера побайтового кода программы и сравнения получаемых сигнатур при помощи редакционного расстояния.

Рекомендации по применению результатов работы для идентификации исполняемых файлов в автоматизированных системах включают в себя указания по формированию архива сигнатур; применению метода и методики идентификации исполняемых файлов, позволяющих обнаруживать нарушения установленных мер политики безопасности в плане запрета на несанкционированную установку программного обеспечения. Также рекомендации включают разработку подходов, направленных на повышение эффективности процесса аудита электронных носителей информации и его автоматизации, позволяющие идентифицировать исполняемые файлы, не основываясь на методах оценки их целостности и анализе метаданных встроенными средствами операционных систем.

В качестве перспектив дальнейшей разработки тематики можно выделить исследования ⁹ связанные с разработкой динамического подхода к идентификации программного обеспечения в близких к системе реального времени ограничениях на скорость реагирования при загрузке ELF-файлов на компьютер пользователя, до момента окончательного скачивания полной версии файла. Такой подход позволит производить превентивную меру по защите автоматизированной системы от несанкционированной установки запрещенного программного обеспечения.

130

Положения, выносимые на защиту, соотнесены с пунктами паспорта специальности 05.13.19 — «⁹ Методы и системы защиты информации, информационная безопасность ¹³⁹»: «б. ²⁷ Модели и методы формирования комплексов средств противодействия угрозам хищения (разрушения, модификации) информации и нарушения информационной безопасности для различного вида объектов защиты вне зависимости от области их функционирования» (¹³⁸ результаты 4-5); «13. Принципы и решения (технические, математические, организационные и др.) по созданию новых и совершенствованию существующих средств защиты информации и обеспечения информационной безопасности» (результаты 2-3); «14. Модели, методы и средства обеспечения внутреннего аудита и мониторинга состояния объекта, находящегося под воздействием угроз нарушения его информационной безопасности» (результат 5). ⁹

131

Список ⁹ сокращений и условных обозначений
API – программный интерфейс приложения (application programming interface ²⁷)

AUC – площадь под roc-кривой (area under roc curve)

CRC – циклический избыточный код (cyclic redundancy check)

SVM – метод опорных векторов (support vector machine)

DLL – динамически подключаемая библиотека (dynamic link library)

ELF – формат исполнимых и компокуемых файлов (Executable and Linkable Format)

HAM – управление аппаратными активами (hardware asset management)

IoT – интернет вещей (internet of things)

ITAM – управление активами информационных технологий (information technology asset management)

ITAM – управление ИТ-активами (IT asset management)

ROC – рабочая характеристика приёмника (receiver operating characteristic)

SAM – управление активами программного обеспечения (software asset management)

SAM – управление программными активами (software asset management)

АС – автоматизированная система

АЭСП – архив эталонных сигнатур программ

ЗИ – защита информации

ИБ – информационная безопасность

ИС – информационная система

ОС – операционная система

ПО – программное обеспечение

СИИФ – сигнатуры идентифицируемого исполняемого файла

СП – сигнатура программы

УК РФ – уголовный кодекс российской федерации

УСП – унифицированных сигнатур программ

132

ЭВМ – электронно-вычислительная машина

ЭСП – эталонной сигнатурой программы

133

Список использованной литературы

1. Tool Interface Standards (TIS). Executable and Linking Format (ELF)

Specification Version 1.2. TIS Commit. 1995. 106 с.

2. ГОСТ Р ИСО/МЭК 27002-2012 Информационная технология (22 ИТ). 49 Методы и средства обеспечения безопасности. Свод норм и правил менеджмента информационной безопасности. М.: 22 Стандартинформ, 2012.

3. 45 ГОСТ Р ИСО/МЭК 27001-2006 Информационная технология (22 ИТ). 49 Методы и средства обеспечения безопасности. Системы менеджмента информационной безопасности. Требования. М.: 22 Стандартинформ, 2006.

4. ГОСТ 39 Р 50922-2006. Защита информации. Основные термины и определения. М.: Стандартинформ, 2006.

5. ФЗ от 27 июля 2006 г. No 149- ФЗ «Об информации, информационных технологиях и о защите информации (с 40 изменениями 13 на 18 декабря 2018 года)». Собрание законодательства Российской Федерации 139 , N 31 (ч.I), 31.07.2006, ст.3448, 2006.

6. Закон РФ от 21 июля 1993 г. No5485-1 «О государственной тайне (с изменениями на 29 июля 2018 года)». Собрание законодательства Российской Федерации 139 , N 41, 13.10.1997, ст.4673, 1993.

7. Указ Президента РФ от 6 марта 1997 г 139 . No188 « Об утверждении перечня сведений конфиденциального характера (с 63 изменениями на 13 июля 2015 года)». Собрание законодательства Российской Федерации 139 , N 10, 10.03.97, ст.1127, 1997.

8. ФЗ от 20 февраля 1995 г. No 24-ФЗ « Об информации, информатизации и защите информации (с изменениями на 10 января 2003 года) (утратил силу с 09.08.2006 на основании Федерального закона от 27.07.2006 N 149-ФЗ)». 138 Собрание законодательства Российской Федерации 139 N 8, 20.02.95, ст.609. ст.2 с.

9. Петренко С.А., Курбатов В.А. Политики безопасности компании при работе в Интернет. 2011. 396 с.

134

10. Boukhtouta A. и др. Graph-theoretic characterization of cyber-threat infrastructures // Digit. Investig. 2015. Т. 14, No 1. С. 3–15.

11. Бондарев В. Введение в информационную безопасность автоматизированных систем. 2018. 252 с.

12. Software Management: Security Imperative, Business Opportunity [Электронный ресурс]. 2018. С. 24. URL: <https://gss.bsa.org/wp->

content/uploads/2018/05/2018_BSA_GSS_Report_en.pdf. Режим доступа:

свободный, дата обращения: 18.03.2018.

13. 69 Отчет о видах наказания по наиболее тяжкому преступлению (без учета

сложения) [94 Электронный ресурс]. 2017. URL:

<http://www.cdep.ru/index.php?id=79&item=69> =4572. Режим доступа: свободный,

дата обращения: 11.12.2018.

14. Базовая модель угроз безопасности персональных данных при их обработке в

информационных системах персональных данных (выписка). ФСТЭК России,

2008 год. 2008.

15. 92 Быкова В.В., Катаева А.В. Методы и средства анализа информативности

признаков при обработке медицинских данных. // Журнал Программные

продукты и системы / Softw. Syst. 2016. Т. 2, No 114. С. 172–178.

16. 68 Hastie T., Tibshirani R., Friedman J. The elements of statistical learning: Data

Mining, Inference, and Prediction. // Springer 105 . 2017. 745 с.

17. Еремеев М.А., Кравчук А.В. Анализ методов распознавания вредоносных

программ. // Журнал Вопросы защиты информации. 2014. Т. 3. С. 44–51.

18. Giesen D. Philosophy and practical implementation of static analyzer tools // Softw.

Dev. J. 1998. No 7. С. 12–31.

19. Christopher C.N. Evaluating Static Analysis Frameworks. 2003. 213 с.

20. Seo K. и др. Detecting Similar Files Based on Hash and Statistical Analysis for

Digital Forensic Investigation // Digital Forensic Investigation. In: Computer

Science and its Applications (CSA 2009). 2009. С. 1–6.

21. Breitinger F. B.H. Similarity Preserving Hashing: Eligible Properties and a New

Algorithm MRSH-v2 // Digital Forensics and Cyber Crime. ICDF2C 2012. Lecture 111

135

Notes of the Institute for Computer Sciences, Social Informatics and

Telecommunications Engineering. 2012. С. vol 111 114.

22. Kornblum J.D. Identifying almost identical files using context triggered piecewise

hashing 28 // Digit. Investig. 2006. Т. 3. С. 91–97.

23. Long C., Guoyin W. An Efficient Piecewise Hashing Method for Computer

Forensics // First International Workshop on Knowledge Discovery and Data

Mining (WKDD 2008). 2008. С. 635–638.

24. Baier H., Frank B. Security Aspects of Piecewise Hashing in Computer Forensics

// 2011 Sixth International Conference on IT Security Incident Management and IT

Forensics. С. 21–36.

25. Цифровые подписи в Windows 7 [Электронный ресурс]. 2012. URL:

<http://windowsnotes.ru/windows-7/cifrovye-podpisi-v-windows-7>. Режим

доступа: свободный, дата обращения: 23.04.2018.

26. Niemela J. It's Signed, therefore it's Clean, right? [Электронный ресурс] //

Protecting the irreplaceable, f-secure.com. 2010. С. 1–30. URL: [https://www.f-](https://www.f-secure.com/weblog/archives/Jarno_Niemela_its_signed.pdf)

[secure.com/weblog/archives/Jarno_Niemela_its_signed.pdf](https://www.f-secure.com/weblog/archives/Jarno_Niemela_its_signed.pdf). Режим доступа:

свободный, дата обращения: 23.04.2018.

27. Sorokin I. Comparing files using structural entropy // J. Comput. Virol. Hacking

Tech. 2011. Т. 7. С. 259–265.

28. Копыльцов А.В., Сорокин И.В. Вейвлет-анализ структурной энтропии файлов

// Известия Российского государственного педагогического университета им.

А.И. Герцена. 2011. Т. 138. С. 7–15.

29. Копыльцов А.В., Сорокин И.В. Алгоритм выравнивания последовательностей

сегментов файлов // Вестник ИНЖЭКОНа. Серия: Технические науки. 2011.

Т. 8. 11 С. 31–36.

30. Антонов А.Е., Федулов А.С. Идентификация типа файла на основе

структурного анализа 29 // ПРИКЛАДНАЯ ИНФОРМАТИКА. 2013. Т. 2, No 44.

С. 68–77.

31. Сорокин И.В., Копыльцов А.В. Модели распознавания вредоносных

программ на основе энтропийных характеристик // Региональная

информатика и информационная безопасность. 2015. С. 313–316.

32. Сорокин И.В., Копыльцов А.В. Использование неокогнитрона для распознавания вредоносных файлов // Известия СПбГЭТУ ЛЭТИ. 2013. Т. 3. 11 С. 45–51.
33. Казарин О.В. Безопасность программного обеспечения компьютерных систем. М.: МГУЛ, 2003. 212 с.
34. Гроувер Д. и др. Защита программного обеспечения: Пер. с англ / под ред. Д.Гроувера. 1992. 288 с.
35. Schultz M.G. и др. Data mining methods for detection of new malicious executables // Proceedings of the IEEE Symposium on Security and Privacy. Los Alamitos 29, CA, 2001. С. 38–49.
36. Cohen W.W. Fast effective rule induction // Proceedings of the Twelfth International Conference on Machine Learning 115. San Francisco, CA 132, 1995. С. 115–123.
37. Kolter J., Maloof M.. Learning to detect and classify malicious executables in the wild // J. Mach. Learn. Res 51. 2006. Т. 7. С. 2721–2744.
38. Bergeron J. и др. Static analysis of binary code to isolate malicious behavior // 1999 Workshop on Enabling Technologies on Infrastructure for Collaborative Enterprises. 1999. С. 184–189.
39. Christodorescu, Mihai Jha S. Static Analysis of Executables to Detect Malicious Patterns 24 // 12th USENIX Security Symposium. 2003. С. 169–186.
40. Christodorescu M. и др. Semantics-aware malware detection // 2005 IEEE Symposium on Security and Privacy. 2005. С. 32–46.
41. Bilal D. Opcodes as predictor for malware // Int. J. Electron. Secur. Digit. Forensics. 2007. Т. 1. С. 156–168.
42. Santos I. и др. Opcode sequences as representation of executables for data-mining-based unknown malware detection 51 // Inf. Sci. (Ny). 2013. Т. 231. С. 64–82.
43. Shabtai A. и др. Detecting unknown malicious code by applying classification techniques on opcode patterns 33 // Secur. Inform. 2012. Т. 1, No 1. С. 1–22.
44. Santos I. и др. OPEM: A Static-Dynamic Approach for Machine-Learning-Based 51 Malware Detection 51 // International Joint Conference CISIS'12-ICEUTE 12-SOCO 12 Special Sessions. Berlin, Heidelberg, 2013. С. 271–280.
45. Darabian H. и др. An opcode-based technique for polymorphic Internet of Things malware detection // Concurr. Comput. Pract. Exp. 2019. С. 1–14.
46. Zhang H. и др. Classification of ransomware families with machine learning based on N-gram of opcodes // Futur. Gener. Comput. Syst. 2019. Т. 90. С. 211–221.
47. Li P. и др. On Challenges in Evaluating Malware Clustering // International Workshop on Recent Advances in Intrusion Detection 80. 2010. С. 238–255.
48. Ying-xu. L., Zeng-hui. L. Unknown Malicious Identification // Adv. Electr. Eng. Comput. Sci. Lect. Notes Electr. Eng. 2009. Т. 39. С. 301–312.
49. Ebringer T., Sun L., Boztas S. A Fast Randomness Test that Preserves Local Detail 24 // 18th Virus Bulletin International Conference. United Kingdom, 2008. С. 34–42.
50. Willems C., Holz T., Freiling F. Toward automated dynamic malware analysis using cwsandbox // 2007 IEEE 130 Symposium on Security and Privacy (S&P 2007). 2007. С. 32–39.
51. Bayer U., Kruegel C., Kirda E. Ttanalyze: A tool for analyzing malware // 15th European Institute for Computer Antivirus Research (EICAR 2006). 2006. С. 180–192.
52. Song D. и др. Bitblaze: A new approach to computer security via binary analysis // 4th International Conference on Information Systems Security. 2008. С. 1–25.
53. AVG [Электронный ресурс]. URL: https://www.avg.com/en-ww/avg-and-norman-business-products?_ga=2.1955633.1678056726.1552380509-211382991.1552380509.

54. Corporation S. Advanced Threat Protection [Электронный ресурс]. URL: <https://www.symantec.com/products/threat-protection>.
55. Lin C.H., Pao H.K., Liao J.W. Efficient dynamic malware analysis using virtual time control mechanics // Comput. Secur. 2018. Т. 73. С. 359–373.
56. Egele M. и др. A Survey on Automated Dynamic Malware-Analysis Techniques and Tools // ACM Comput. Surv. 51. 2012. Т. 44, No 2. С. 42.
57. Gray-Donald T.A., Price M.W. Date and time simulation for time-sensitive applications: пат. US 8,418,151 B2 USA. USA, 2013. С. 9.
58. Bulazel A., Bülent Y. A Survey On Automated Dynamic Malware Analysis Evasion and Counter-Evasion: Pc, mobile, and web. // the 1st Reversing and Offensive-oriented Trends Symposium. 2017. С. 1–21.
59. Windows Privilege Escalation Fundamentals [Электронный ресурс]. 2014. URL: <http://www.fuzzysecurity.com/tutorials/16.html>.
60. Вартанов С.П., Герасимов А.Ю. Динамический анализ программ с целью поиска ошибок и уязвимостей при помощи целенаправленной генерации входных данных 55 // Труды ИСП РАН. 2014. Т. 26, No 1. С. 375–394.
61. Мельников П.В., Горюнов М.Н., Анисимов Д.В. Подход к проведению динамического анализа исходных текстов программ // Вопросы кибербезопасности. 2016. Т. 3, No 16. С. 33–39.
62. Klosterboer L. Implementing ITIL configuration management. 2007. 264 с.
63. England R. Owning ITIL. 2009. 178 с.
64. Bird J. DevOpsSec. 2016.
65. Samanage. Samanage [Электронный ресурс]. 2019. URL: <https://www.samanage.com> (дата обращения: 28.03.2019).
66. Kaspersky. Kaspersky Systems Management [Электронный ресурс]. 2019. URL: <https://www.kaspersky.ru> (дата обращения: 28.03.2019).
67. Microsoft. Microsoft Assessment and Planning Toolkit [Электронный ресурс]. 2018. URL: <https://www.microsoft.com/en-us/Download/details.aspx?id=135> =7826 (дата обращения: 28.03.2019).
68. FinalWire. AIDA64 [Электронный ресурс]. URL: <https://www.aida64.com> (дата обращения: 28.03.2019).
69. Lansweeper. Lansweeper [Электронный ресурс]. 2019. URL: <https://www.lansweeper.com> (дата обращения: 28.03.2019).
70. Bayer U. и др. Scalable, behavior-based malware clustering // Proceedings of the Network and Distributed System Security Symposium 29. 2009. С. 8–11.
71. Krivtsova I.E., Lebedev I.S., Salakhutdinova K.I. Identification of Executable Files on the basis of Statistical Criteria // Proceedings of the 20th Conference of Open Innovations Association FRUCT, IET 7. 2017. С. 202–208.
72. Салахутдинова К.И., Лебедев И.С., Кривцова И.Е. Подход к выбору информативного признака в задаче идентификации программного обеспечения // Научно-технический вестник информационных технологий, механики и оптики 139. 2018. Т. 18, No 2. С. 278–285.
73. Рудина Т.Д. Разработка способа идентификации elf-файлов на основе нейронной сети. 2018. 52 с.
74. Салахутдинова, К.И. Лебедев И.С., Кривцова И.Е. Алгоритм градиентного бустинга деревьев решений в задаче идентификации программного обеспечения 8 // Научно-технический вестник информационных технологий, механики и оптики 139. 2018. Т. 18, No 6.
75. Салахутдинова, К.И. Малков В.В., Кривцова И.Е. Сравнительный анализ подходов к идентификации программного обеспечения // Безопасность информационных технологий. 2019. Т. 26, No 2. С. 58–66.
76. Салахутдинова К.И. Повышение точности идентификации программного обеспечения путем использования аддитивного критерия Фишберна //

- Информационные технологии. 2019. Т. 25, No 10. С. 609–614.
77. Standard I. 24765-2017-I.I. Systems and software engineering--Vocabulary. 2017.
78. Benford F. The Law of Anomalous Numbers // Proceedings of the American Philosophical Society, 78. 1938. С. 551–572.
79. Fewster R.M. A simple explanation of Benford's Law // Am. Stat. 2009. Т. 63, No 1. С. 26–32.
80. TESTING BENFORD'S LAW. File sizes in the Linux 2.6.39.2 source tree [Электронный ресурс]. URL: <http://testingbenfordslaw.com/linux-filesizes>.
81. Салахутдинова К.И., Лебедев И.С., Кривцова И.Е. Подход к выбору информативного признака в задаче идентификации программного обеспечения // Научно-технический вестник информационных технологий, механики и оптики 139 . 2018. Т. 18, No 2. С. 278–285.
82. Salakhutdinova K.I. и др. Studying the Effect of Selection of the Sign and Ratio in the Formation of a Signature in a Program Identification Problem // Autom. Control 140 Comput. Sci. 2018. Т. 52, No 8. С. 1101–1104.
83. Смирнов Н.В., Дунин-Барковский И.. Курс теории вероятностей и математической статистики для технических приложений. М.: Наука, 1969. 511 с.
84. 96 Куликов Е.И. Прикладной статистический анализ. М.: Горячая линия-Телеком, 2008. 463 с.
85. Сидоренко Е.В. Методы математической обработки в психологии. СПб.: Речь, 2010. 349 с.
86. Семенов В.А. Теория вероятностей и математическая статистика 139 : Учебное пособие. Стандарт третьего поколения. СПб.: Питер, 2013. 192 с.
87. Каллан Р. Основные концепции нейронных сетей. Пер. с англ.— М.: Издательский дом «Вильямс», 2001. 291 с.
88. Электронный учебник по статистике [Электронный ресурс]. URL: <http://statsoft.ru/home/textbook 125 /modules/stneunet.html>.
89. Кафтаников И.Л., Парасич А.В. Особенности применения деревьев решений в задачах классификации // Вестник ЮУрГУ. Серия «Компьютерные технологии, управление, радиоэлектроника». 2015. Т. 3, No 15. С. 26–32.
90. Freund Y., Schapire R. Experiments with a New Boosting Algorithm // Machine Learning: Proceedings of the Thirteenth International Conference 54 . 1996. С. 148–156.
91. Дружков П.Н., Золотых Н.Ю., Половинкин 52 А.Н. Реализация параллельного алгоритма предсказания в методе градиентного бустинга деревьев решений 52 // Вестник ЮУрГУ. 2011. Т. 37, No 254. С. 82–89.
92. LightGBM GitHub [Электронный ресурс].
93. XGBoost GitHub [Электронный ресурс]. URL: <https://github.com/dmlc/xgboost>.
94. Китов В.В. Исследование точности метода градиентного бустинга со случайными поворотами // Статистика и экономика. 2016. Т. 4. С. 22–26.
95. CatBoost [Электронный ресурс]. URL: <https://catboost.ai/docs/>.
96. Коварцев А.Н., Даниленко А.Н. Алгоритмы и анализ сложности: учебник. Изд-во Сам. Самара, 2018. 128 с.
- 141
97. Подиновский В.В. Аксиоматическое решение проблемы оценки важности критериев в многокритериальных задачах // Современное состояние теории исследования операций 70 . 1979. С. 117–149.
98. Анохин А.М. и др. Методы определения коэффициентов важности критериев // Автомат. и телемех. 1997. No 8. С. 3–35.
99. Фишберн П. Теория полезности для принятия решений. 1978. 352 с.
100. Фишберн П. Методы оценки аддитивных ценностей // Статистическое измерение качественных характеристик 70 . 1972. С. 8–34.
101. Салахутдинова К.И. и др. Идентификации программного обеспечения в

Приложение 1. Значения информативности для 118 ассемблерных команд

I(x)

Ассемблерная

команда

I(x)

Ассемблерная

команда

I(x)

Ассемблерная

команда

0,2024 cmpsb 0,2024 aaa 0,2025 rcl

0,2024 cmpsw 0,2024 les 0,2025 std

0,2024 esc 0,2024 daa 0,2025 sbb

0,2024 jc 0,2024 aas 0,2025 lahf

0,2024 jcz 0,2024 aam 0,2025 ror

0,2024 jna 0,2024 ja 0,2025 ret

0,2024 jnae 0,2024 mul 0,2025 jne

0,2024 jnb 0,2024 js 0,2025 cmc

0,2024 jnbe 0,2024 loop 0,2025 popf

0,2024 jnc 0,2024 jge 0,2025 rcr

0,2024 jng 0,2024 jp 0,2025 adc

0,2024 jnge 0,2024 hlt 0,2025 jno

0,2024 jnl 0,2024 jl 0,2025 in

0,2024 jnle 0,2024 jns 0,2025 shl

0,2024 jnz 0,2024 loopne 0,2026 imul

0,2024 jpe 0,2024 ired 0,2026 sar

0,2024 jpo 0,2024 das 0,2026 sahf

0,2024 jz 0,2024 rep 0,2026 jmp

0,2024 lodsb 0,2024 jae 0,2026 xor

0,2024 lodsw 0,2024 idiv 0,2026 nop

0,2024 loopnz 0,2025 jbe 0,2026 and

0,2024 loopz 0,2025 jb 0,2026 jo

0,2024 movsb 0,2025 pushf 0,2026 je

0,2024 movsw 0,2025 not 0,2026 stc

0,2024 repe 0,2025 clc 0,2026 test

0,2024 repne 0,2025 rol 0,2026 push

0,2024 retn 0,2025 int 0,2026 retf

0,2024 sal 0,2025 jle 0,2027 out

0,2024 scasb 0,2025 jg 0,2027 cmp

0,2024 scasw 0,2025 sub 0,2027 inc

0,2024 stosb 0,2025 loope 0,2028 or

0,2024 stosw 0,2025 xlat 0,2029 xchg

0,2024 wait 0,2025 cld 0,203 pop

0,2024 cwd 0,2025 sti 0,203 add

0,2024 cbw 0,2025 shr 0,2032 lea

0,2024 div 0,2025 jnp 0,2036 call

0,2024 neg 0,2025 repnz 0,2044 mov

0,2024 into 0,2025 cli 0,2095 lock

0,2024 lds 0,2025 dec

0,2024 aad 0,2025 repz

Приложение 2. Числовые идентификаторы эталонных программ

Имя про-

граммы

Числовой

идентификатор

Имя про-

граммы

Числовой

идентификатор

Имя про-

граммы

Числовой

идентификатор

Имя про-

граммы

Числовой

идентификатор

acpid 1 autogen 15 beanstalkd 28

acpid 1 avahi-autoipd 16 beanstalkd 28

acpid 1 avahi-autoipd 16 bison 29

acpid 1 avahi-autoipd 16 bison 29

activity-log-manager 2 avahi-autoipd 16 bison 29

activity-log-manager 2 avahi-autoipd 16 bogofilter-bdb 30

activity-log-manager 2 avahi-daemon 17 bogofilter-bdb 30

aide 3 avahi-daemon 17 bonnie++ 31

aide 3 avahi-daemon 17 brasero 32

aide 3 avahi-daemon 17 brasero 32

aide 3 avahi-daemon 17 brasero 32

amarok 4 b43-fwcutter 18 brasero 32

anacron 5 b43-fwcutter 18 brasero 32

anacron 5 b43-fwcutter 18 bsd-mailx 33

anacron 5 bacula-console 19 bsd-mailx 33

apmd 6 bacula-console 19 bsd-mailx 33

apmd 6 bacula-console 19 bsd-mailx 33

appstream-util 7 bacula-console 19 bsd-mailx 33

appstream-util 7 bacula-console-qt 20 bsd-mailx 33

apt 8 bacula-console-qt 20 ccache 34

apt 8 bacula-console-qt 20 ccache 34

apt 8 bacula-fd 21 ccache 34

aptitude 9 bacula-fd 21 ccache 34

aptitude 9 bacula-fd 21 ccache 34

ark 10 bacula-fd 21 ccache 34

ark 10 bacula-sd 22 cdparanoia 35

aspell 11 bacula-sd 22 cdparanoia 35

aspell 11 bacula-sd 22 cdrdao 36

aspell 11 bacula-sd+dfsg 23 cdrdao 36

aspell 11 bacula-sd+dfsg 23 ceph-fuse 37

at 12 bacula-sd+dfsg 23 ceph-fuse 37

at 12 bacula-traymonitor 24 ceph-fuse 37

at 12 bacula-traymonitor 24 checkbox-gui 38

at 12 baobab 25 checkbox-gui 38

attr 13 baobab 25 cheese 39

attr 13 baobab 25 cheese 39

authbind 14 bc 26 cheese 39

authbind 14 bc 26 cheese 39

authbind 14 bc 26 choqok 40

autogen 15 bdfresize 27 choqok 40

autogen 15 bdfresize 27 chrpath 41

autogen 15 autogen 15 chrpath 41

144

Продолжение таблицы

Имя про-

граммы

Числовой

иденти-

фикатор

Имя про-

граммы

Числовой

иденти-

фикатор

Имя про-

граммы

Числовой

иденти-

фикатор

clamdsan 42 curl 50 wodim 62

clamdsan 42 curl 50 wodim 62

сmake 43 curl 50 xbrlapi 63

сmake 43 curl 50 xbrlapi 63

сmake 43 curl 50 xbrlapi 63

сmake 43 curl 50 xbrlapi 63

сmake 43 cvs 51

сmake 43 cvs 51

conntrack 44 cvs 51

conntrack 44 cvs 51

conntrackd 45 cvs 51

conntrackd 45 cvsps 52

corosync 46 cvsps 52

corosync 46 cvsps 52

corosync 46 cvsps 52

corosync 46 dc 53

corosync 46 dc 53

corosync 46 dc 53

crash 47 dc 53

crash 47 finger 54

crash 47 finger 54

crash 47 gceph 55

crash 47 gceph 55

crash 47 genisoimage 56

cron 48 genisoimage 56

cron 48 genisoimage 56

cron 48 genisoimage 56

cron 48 iasl 57

cron 48 iasl 57

cron 48 jsvc 58

cups-browsed 49 jsvc 58

cups-browsed 49 nfct 59

cups-browsed 49 nfct 59

cups-browsed 49 radosgw 60

cups-browsed 49 radosgw 60

cups-browsed 49 radosgw 60

curl 50 radosgw 60

curl 50 radosgw 60

curl 50 radosgw 60

curl 50 rbd-fuse 61

curl 50 rbd-fuse 61

curl 50 rbd-fuse 61

curl 50 wodim 62

curl 50 wodim 62

145

Приложение 3. Коды сравнения сигнатур с помощью библиотек градиентного

бустинга деревьев решений

Таблица В.1 – LightGBM

Задаем входные данные

```
import pandas
```

```
import lightgbm
```

```
import numpy
```

```
traindata = pandas.read_csv("tradd.csv",header=None,sep=";")
```

```
testdata = pandas.read_csv("tsadd.csv",header=None,sep=";")
```

```
ytrain=traindata[0]
```

```
ytest=testdata[0]
```

```
xtrain=traindata.drop(0, axis=1)
```

```
xtest=testdata.drop(0, axis=1)
```

```
train = lightgbm.Dataset(xtrain,ytrain)
```

Использование библиотеки lightgbm и расчет Accuracy

```
param = {
```

```
'boosting_type': 'gbdt',
```

```
'objective': 'multiclass',
```

```
'metric': 'multi_logloss',
```

```
'num_class': 64,
```

```
'num_threads': 4,
```

```
'learning_rate': 0.3,
```

```
'max_depth': 7,
```

```
'verbosity': 2
```

```
}
```

```
num_round = 1000
```

```
for learnrate in [0.3]:
```

```
for i in [4]:
```

```
param['learning_rate'] = learnrate
```

146

```
param['max_depth'] = i
```

```
model = lightgbm.train(param,train,num_round)
```

```
preds = model.predict(xtest)
```

```
result = []
```

```
for j in range (len(preds)):
```

```
result.append(preds[j].argmax())
```

```
error_rate = numpy.sum(result != ytest)/ytest.shape[0]
```

```
print("learnrate",learnrate,"max_depth",i," accuracy: ",1-error_rate)
```

147

Таблица В.2 – CatBoost

Задаем входные данные

```
train_data = []
```

```
test_data = []
```

```
train_labels = []
```

```

test_labels = []

with open("Iden/Train_data_ne_unif/sig_ob","r") as f:
for line in f:
train_data.append([int(x) for x in line.split()])

with open("Iden/Train_data/sig","r") as f:
for line in f:
test_data.append([int(x) for x in line.split()])

for x in range(len(test_data)):
test_labels.append(test_data[x][0])

del test_data[x][0]

for x in range(len(train_data)):
train_labels.append(train_data[x][0])

del train_data[x][0]

```

Использование библиотеки CatBoostClassifier и расчет Accuracy

```

from catboost import Pool, CatBoostClassifier

pool = Pool(train_data, train_labels)
eval_pool = Pool(test_data, test_labels)

# Initialize CatBoostClassifier
model = CatBoostClassifier(
iterations=1000,
learning_rate=0.7,
depth=2,
loss_function='MultiClass',
l2_leaf_reg = 1,
)

# Fit model
model.fit(pool, plot=True)

# Get predicted classes
preds_class = model.predict(test_data)

# Get predicted probabilities for each class
preds_proba = model.predict_proba(test_data)

# Get predicted RawFormulaVal

148
preds_raw = model1.predict(test_data, prediction_type='RawFormulaVal')

```

Рассчитываем Accuracy

```

tp = 0
fp = 0

for t in range(len(preds_class)):
if test_labels[t] == preds_class[t]:
tp += 1
else:
fp += 1

print("Accuracy = ", tp/(tp+fp))

```

149

Таблица В.3 – XGBoost

Задаем входные данные

```

from __future__ import division

```

```

import numpy as np

```

```

import xgboost as xgb 65

```

```

import pandas as pd

```

```

#get data fo further training and test

```

```

#label need to be 0 to num_class -1

```

```

traindata = np.genfromtxt('trje.csv',dtype='int', delimiter=';', converters=(29: lambda x:int(x), 30:

```

```

lambda x:int(x)-1))

```

```

testdata = np.loadtxt('tsje.csv', dtype='int', delimiter=',', converters=(29: lambda x:int(x), 30: lambda
x:int(x)-1))
sztr = traindata.shape
sztst = testdata.shape
train_X = traindata[:, :]
train_Y = traindata[:, 30]
test_X = testdata[:, :]
test_Y = testdata[:, 30]

```

```

xg_train = xgb.DMatrix(train_X, label=train_Y)

```

```

xg_test = xgb.DMatrix(test_X, label=test_Y)

```

Использование библиотеки xgb и расчет Accuracy

```

# setup parameters for xgboost

```

```

param = {}

```

```

# use softmax multi-class classification

```

```

param['objective'] = 'multi:softmax'

```

```

# scale weight of positive examples

```

```

param['eta'] = 0.3

```

```

param['max_depth'] = 2

```

```

param['nthread'] = 4

```

```

param['num_class'] = 63

```

```

param['silent'] = 1

```

```

watchlist = [(xg_train, 'train'), (xg_test, 'test')]

```

```

num_round = 1000

```

```

bst = xgb.train(param, xg_train, num_round, watchlist)

```

```

pred = bst.predict(xg_test)

```

```

error_rate = (1-np.sum(pred != test_Y) / test_Y.shape[0])*100

```

```

print("eta=", param['eta'], 'accuracy = {}'.format(error_rate), '%')

```

150

Приложение 4. Свидетельство о регистрации программы для ЭВМ

151

152

153

154

155

156

Приложение 5. Копии актов внедрения

157

158

159

Приложение 6. Публикации соискателя по теме Диссертации

В научных журналах, рекомендованных ВАК:

1. Салахутдинова К.И., Лебедев И.С., Кривцова И.Е. Подход к выбору

информативного признака в задаче идентификации программного обеспечения //

Научно-технический вестник информационных технологий, механики и оптики **139** -

2018. - Т. 18. - No 2(114). - С. 278–285

2. Салахутдинова, К.И. Лебедев И.С., Кривцова И.Е., Сухопаров М.Е.

Исследование влияния выбора признака и коэффициента (ratio) при формировании

сигнатуры в задаче по идентификации программ // Проблемы информационной

безопасности. Компьютерные системы **139** . 2018. No 1. С. 136–141.

3. Салахутдинова, К.И. Лебедев И.С., Кривцова И.Е. Алгоритм

градиентного бустинга деревьев решений в задаче идентификации программного

обеспечения **8** // Научно-технический вестник информационных технологий,

механики и оптики **139** . 2018. Т. 18, No 6.

4. Салахутдинова К.И., Малков В.В., Кривцова И.Е. Сравнительный анализ подходов к идентификации программного обеспечения // Безопасность информационных технологий - 2019. - Т. 26. - No 2. - С. 58-66

5. Салахутдинова К.И., Лебедев И.С., Кривцова И.Е., Анисимов А.С. Идентификации программного обеспечения в задаче аудита электронных носителей информации // Авиакосмическое приборостроение - 2019. - No 9. - С. 20-28

6. Салахутдинова К.И. Повышение точности идентификации программного обеспечения путем использования аддитивного критерия Фишберна // Информационные технологии -2019. - Т. 25. -No 10. - С. 609-614

7. Бажаев Н., Давыдов А.Е., Кривцова И.Е., Лебедев И.С., Салахутдинова К.И. Подход к анализу состояния информационной безопасности беспроводной сети // Прикладная информатика - 2016. - Т. 11. - No 6(66). - С. 121-128

8. Кривцова И.Е., Салахутдинова К.И., Юрин И.В. Метод идентификации исполняемых файлов по их сигнатурам // Вестник Государственного университета **91** **160**

морского и речного флота имени адмирала С.О. Макарова - 2016. - No 1(35). - **91** С. 215-224

В изданиях, индексируемых в международных базах цитирования Web of Science, Scopus **68** :

9. Lebedev I.S., Korzhuk V., Krivtsova I., Salakhutdinova K., Sukhoparov M.E., Tikhonov D. Using Preventive Measures for the Purpose of Assuring Information Security of Wireless Communication Channels // Proceedings of the 18th Conference of Open Innovations Association FRUCT - 2016, pp **16** . 167-173

10. Bazhayev N., Lebedev I.S., Krivtsova I.E., Sukhoparov M.E., Salakhutdinova K., Davydov A.E., Shaparenko I.M. Evaluation of the available wireless remote devices subject to the information impact // 10 th IEEE International Conference on Application of Information and Communication Technologies, AICT 2016 - Conference Proceedings **16** (Azerbaijan, Baku, 12-14 October 2016) - 2016, pp. 1-6

11. Lebedev I.S., Krivtsova I.E., Korzhuk V., Bazhayev N., Sukhoparov M.E., Pecherkin S., Salakhutdinova K. The Analysis of Abnormal Behavior of the System Local Segment on the Basis of Statistical Data Obtained from the Network Infrastructure Monitoring // Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) - 2016, Vol. 9870, pp **16** . 503-511

12. Krivtsova I.E., Lebedev I.S., Salakhutdinova K.I. Identification of Executable Files on the basis of Statistical Criteria// Proceedings of the 20th Conference of Open Innovations Association FRUCT, IET - 2017, pp **7** . 202-208

13. Salakhutdinova K., Lebedev I.S., Krivtsova I.E., Bazhayev N., Sukhoparov M.E., Smirnov P.I., Markelov D.V., Davydov A.E., Pecherkin S., Kolcherin D.V., Shaparenko I.M., Iskanderov Y. A Frequency Approach to Creation of Executable File Signatures for their Identification//11 th IEEE International Conference on Application of Information and Communication Technologies, AICT 2017 - Conference Proceedings **16** (Moscow, 20-22 september 2017), IET - 2017, pp. 261-267.

14. Salakhutdinova, K.I. Krivtsova I.E., Lebedev I.S., Sukhoparov M.E. An Approach to Selecting an Informative Feature in Software Identification **8** // Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics **124**). 2018. С. 318-327.

15. Salakhutdinova K.I., Lebedev I.S., Krivtsova I.E., Sukhoparov M.E. Studying the Effect of Selection of the Sign and Ratio in the Formation of a Signature in a Program Identification Problem // Automatic Control and Computer Sciences - 2018,

16. Semenov, Viktor V., Ilya S. Lebedev, Mikhail E. Sukhoparov and Kseniya

I. Salakhutdinova. Application of an Autonomous Object Behavior Model to Classify the Cybersecurity State. *Internet of Things, Smart Spaces, and Next Generation Networks and Systems* 8. NEW2AN 2019, ruSMART 2019. Lecture Notes in Computer Science - 2019, Vol. 11660, pp. 104-112

В других научных журналах и изданиях:

17. Салахутдинова К.И., Овсянникова В.В., Трофимов А.А., Бессонова Е.Е., Ефремов А.А., Настека А.В. Анализ защищенности систем "Умный дом" // *Региональная информатика (РИ-2014). XIV Санкт-Петербургская международная конференция «Региональная информатика (РИ-2014)». Санкт-Петербург, 29-31 октября 2014 г.: Материалы конференции* 36. \ СПОИСУ. – СПб 45, 2014. - 2014. - С. 124

18. Ефремов А.А., Трофимов А.А., Настека А.В., Салахутдинова К.И., Овсянникова В.В. Защита управляющих сигналов в системе «Умный дом» // *Сборник тезисов докладов конгресса молодых ученых. Электронное издание. – СПб: Университет ИТМО* 16, 2015. – 2015

19. Овсянникова В.В., Настека А.В., Ефремов А.А., Салахутдинова К.И., Трофимов А.А. Защита системы «Умный дом» от программных сбоев // *Сборник тезисов докладов конгресса молодых ученых. Электронное издание. – СПб: Университет ИТМО* 16, 2015. – 2015

162

20. Druzhinin N.K., Salakhutdinova K.I. Identification of executable file by dint of individual feature // *ISPIT-2015. International Conference on Information Security and Protection of Information Technology* 25. St. Petersburg, Russia, November 5-6, 2015, IET - 2015, pp. 45-47

21. Кривцова И.Е., Салахутдинова К.И. Применение х2-критерия для идентификации elf-файлов. V Всероссийский конгресс молодых ученых // *Сборник ВКМУ – 2016*

22. Салахутдинова К.И., Кривцова И.Е. Условия применения критерия Пирсона для идентификации исполняемых файлов. *Региональная информатика "РИ-2016" - 2016*

23. Салахутдинова К.И., Дружинин Н.К. Идентификация исполняемых файлов по их ассемблерным командам // *Научные работы участников конкурса «Молодые ученые Университета ИТМО» 2016 года* 36. 2017.

24. Салахутдинова К.И., Лебедев И.С. Применение методов статистического анализа для идентификаций версий программного обеспечения удаленных автономных объектов транспортных систем // *Информационная безопасность регионов России (ИБРР-2017). 2017.*

25. Салахутдинова К.И., Лебедев И.С. Использование градиентного бустинга в задаче сравнения сигнатур программ // *Сборник тезисов докладов конгресса молодых ученых. 2018. 136-141 с.*

26. Салахутдинова, К.И. Обзор существующих подходов по аудиту электронных носителей информации // *Сборник тезисов докладов конгресса молодых ученых. Электронное издание* 7. 2019.

Свидетельства о государственной регистрации программы для ЭВМ 139 :

27. Ильин А.Г., Салахутдинова К.И., Юшковский А.В., Катаева В.А., Первушин А.О., Павлов К.С. Программа мониторинга Google Apps for Business. No2016614206, 18.04.2016.

163

28. Ильин А.Г., Салахутдинова К.И., Юшковский А.В., Катаева В.А., Первушин А.О., Павлов К.С. Программа для стеганографического сокрытия информации в медиа файлах. No2016614207, 18.04.2016.

29. Ильин А.Г., Салахутдинова К.И., Юшковский А.В., Катаева В.А., Первушин А.О., Павлов К.С. Программа для поиска и анализа криптоконтейнеров

с носителя информации. No2016611975, 15.02.2016.

30. Ильин А.Г., Салахутдинова К.И., Юшковский А.В., Катаева В.А.,

Первушин А.О., Павлов К.С. Программа для криминалистического анализа IM ICQ

и Jabber. No2016614048, 13.04.2016.

31. Ильин А.Г., Салахутдинова К.И., Юшковский А.В., Катаева В.А.,

Первушин А.О., Павлов К.С. Программа для аудита событий информационной

безопасности на основе модели MapReduce. No2016614208, 18.04.2016.

32. Салахутдинова К.И. Сравнение сигнатур исполняемых файлов.

Свидетельство о государственной регистрации программы для ЭВМ **139** No2019619363,

16.07.2019.