

На правах рукописи



Ушаков Игорь Александрович

**ОБНАРУЖЕНИЕ ИНСАЙДЕРОВ В КОМПЬЮТЕРНЫХ СЕТЯХ
НА ОСНОВЕ КОМБИНИРОВАНИЯ ЭКСПЕРТНЫХ ПРАВИЛ, МЕТОДОВ
МАШИННОГО ОБУЧЕНИЯ И ОБРАБОТКИ БОЛЬШИХ ДАННЫХ**

**Специальность 05.13.19 – Методы и системы защиты информации,
информационная безопасность**

АВТОРЕФЕРАТ

диссертации на соискание ученой степени
кандидата технических наук

Санкт-Петербург – 2020

Работа выполнена в Федеральном государственном бюджетном учреждении науки Санкт-Петербургском институте информатики и автоматизации Российской академии наук.

Научный руководитель: доктор технических наук, профессор
Котенко Игорь Витальевич
СПИИРАН,
главный научный сотрудник
лаборатории проблем компьютерной
безопасности

Официальные оппоненты: **Синещук Юрий Иванович**
доктор технических наук, профессор
Федеральное государственное казенное образова-
тельное учреждение высшего образования
«Санкт-Петербургский университет Министер-
ства внутренних дел Российской Федерации»,
профессор кафедры Специальных информацион-
ных технологий

Ефимов Вячеслав Викторович
кандидат технических наук, доцент
Акционерное общество «Научно-
исследовательский институт «Масштаб»,
советник генерального директора

Ведущая организация: **АО «Научно-исследовательский институт
«Рубин»**

Защита состоится 28 апреля 2020 г. в 14:00 часов на заседании диссертационного совета по защите диссертаций на соискание ученой степени кандидата наук, на соискание ученой степени доктора наук Д 002.199.01, созданного на базе Федерального государственного бюджетного учреждения науки Санкт-Петербургского института информатики и автоматизации Российской академии наук (СПИИРАН) по адресу: 199178, Санкт-Петербург, 14-я линия В.О., 39.

Факс: (812)-328-44-50, тел.: (812)-328-34-11.

С диссертацией можно ознакомиться в библиотеке Федерального государственного бюджетного учреждения науки Санкт-Петербургского института информатики и автоматизации Российской академии наук по адресу: 199178, Санкт-Петербург, В.О., 14 линия, д. 39 и на сайте <http://www.spiiiras.nw.ru/dissovet/>

Автореферат разослан « ___ » _____ 2020 года.

Ученый секретарь совета
диссертационного совета Д 002.199.01,
кандидат технических наук



Зайцева
Александра
Алексеевна

Общая характеристика работы

Актуальность темы диссертации

Современную жизнь сложно представить без информационного взаимодействия, затрагивающего как отдельных членов общества, так и крупные организации, в том числе реализующие интересы целого государства. Помимо очевидных получаемых преимуществ такое взаимодействие несет также ряд существенных недостатков. Так, передача информации по сети подвергает ее триаде угроз информационной безопасности: конфиденциальности, целостности и доступности. При этом безопасность информации должна быть обеспечена как при её передаче через открытые сети, так и внутри компьютерной сети (КС), под которой прежде всего понимается корпоративная компьютерная сеть. Однако к информации в КС, в особенности носящей критически важное значение, могут иметь доступ внутренние сотрудники, часть из которых изначально обладает такими полномочиями, входящими в круг их должностных обязанностей. Таким образом, возникает проблема противодействия атакам на КС, как случайным, так и злонамеренным, производимым в том числе внутренними сотрудниками организации.

Существуют различные способы противодействия инсайдерской деятельности на разных этапах – до самой атаки, во время ее проведения и после атаки. У каждого из способов есть свои достоинства и недостатки, однако важным является тот факт, что информация может устаревать и ее ценность, соответственно, уменьшаться. Следовательно, оказание позднего противодействия инсайдерским атакам может оказаться бессмысленным, поскольку информация к этому времени уже будет скомпрометирована и использована третьими лицами. Так, в случае нарушения целостности данных или предоставления неправомерного доступа к данным факт обнаружения одного из подобных нарушений будет иметь для организации существенно меньший эффект, чем недопущение инсайдерской атаки в целом. Следовательно, востребованным является именно недопущение инсайдерской атаки, что может быть достигнуто путем обнаружения инсайдеров до момента проведения самой атаки. После обнаружения инсайдеров, естественно, предполагается их нейтрализация. Нейтрализация инсайдеров может производиться либо автоматически – программными средствами, либо вручную – экспертами по информационной безопасности.

Согласно последним исследованиям, весь IP-трафик и число устройств, подключенных к сети Интернет, утроятся за следующие 5 лет. Считается, что это произойдет вследствие развития сервисов и услуг, предоставляемых телекоммуникационными компаниями. При этом особую популярность набирают: облачные сервисы в виде Platform-as-a-Service (PaaS) и Software-as-a-Service (SaaS); решения для хранения данных; аналитические системы; решения для ведения бизнеса и прогнозирования рисков; рекомендательные системы. Расширение областей применения сетевых технологий означает децентрализацию сетевой инфраструктуры в целом как в плане хранения данных, так и в плане получаемого доступа к этой инфраструктуре. Это усложняет решение задач, стоящих перед специалистами информационной безопасности, поскольку становится труднее контролировать все аспекты сетевой безопасности при защите критически важных данных от угроз, исходящих как из внешней сети, так и изнутри, от самих участников сети.

Таким образом, основная сложность обнаружения инсайдеров в КС напрямую следует из современных тенденций развития информационных технологий, неразрывно связанных с постоянным увеличением параметров сетевого трафика: его объе-

ма; скорости генерации; количества источников и получателей трафика; количества логических потоков, не связанных со своими целями и задачами; увеличения уровня гетерогенности данных и др.

Все это приводит к существенному усложнению анализаторов трафика, поскольку далеко не все существующие системы способны справляться с такими большими объемами и сложностью, в то время как инсайдеры скрывают свои действия в общем потоке действий законных пользователей. Кроме того, современные инсайдерские атаки являются комплексными и используют множество способов реализации и множество векторов для получения несанкционированного доступа и компрометации информационных объектов во внутренней КС.

Таким образом, *основное противоречие предметной области* заключается в следующем: с одной стороны, необходимо повышение точности обнаружения инсайдеров, поскольку их атаки постоянно усложняются и комплексизируются, сетевой трафик атак становится менее различим из-за роста объема всего трафика в КС, а сами инсайдеры маскируют свои действия под законные; с другой стороны, существующие модели, методики и алгоритмы обнаружения инсайдеров не обладают необходимой эффективностью работы, поскольку или имеют высокий риск пропуска инсайдера (ошибка II рода), или, наоборот, – риск отнесения к инсайдеру законного пользователя (ошибка I рода). Возможной причиной порождения данного противоречия является некоторая субъективность, присущая всем вводимым критериям инсайдерской деятельности. Так, например, часть пользователей, определенных как инсайдеры, могли просто выполнять ряд ошибочных действий: неверно ввести свой пароль, ошибочно скачать документ или отправить документ на неверный адрес, подключить чужое устройство и т.п.

Разрешение указанного противоречия может лежать в плоскости применения высокоэффективных специализированных технологий обработки сетевого трафика для сферы информационной безопасности, а также в сочетании существующих и новых способов анализа и обнаружения инсайдерской деятельности. Все это может быть достигнуто следующим образом.

Во-первых, тенденция роста популярности появления решений для работы с большими данными позволяет предположить гипотетическую востребованность данной технологии для разрешения выявленного выше основного противоречия предметной области. Так, с появлением инструментов для разработки систем, использующих концепцию больших данных, встает вопрос об использовании технологий обработки больших данных для информационной безопасности и, в частности, систем мониторинга безопасности. Становится все сложнее обнаруживать потенциальные угрозы безопасности. Пропускная способность современных систем мониторинга и предупреждения сетевых атак перестает удовлетворять требованиям постоянно разрастающихся сетей: в связи с большим количеством поступающего трафика и низкой скоростью его обработки результаты такого анализа получаются неактуальными и не отражают реального состояния сети. Используя новые и эффективные технологии для агрегации и хранения больших объемов данных, а также для организации работы системы обнаружения злоумышленника, можно добиться нужных результатов, а именно получить достаточный уровень контроля над ситуацией в КС.

Также важно учитывать, что не все модели представления данных в достаточной степени адаптированы к своевременной обработке больших объемов информации и событий. Специфика задач кибербезопасности заключается в необходимости при-

менения новых моделей баз данных и использовании методов обработки больших данных для анализа трафика компьютерных сетей.

Во-вторых, существующим и хорошо зарекомендовавшим себя подходом к обнаружению инсайдеров (учитывая сложность строгой формализации критериев обнаружения последних и их возможности к сокрытию своих действий) является использование алгоритмов на основе правил, составленных экспертами с учетом собственного накопленного опыта и существующих «лучших практик» (экспертных правил).

И, в-третьих, учет комплексности проводимых инсайдерами атак, а также их распределенности, в том числе по сети (например, атака на целый ряд не связанных хостов), по объектам (например, попытка доступа к частям документа с целью сбора общей критической массы конфиденциальной информации), по времени (например, последовательность событий, связанных длительным промежутком времени), дает возможность предположить востребованность применения методов машинного обучения, позволяющих учитывать множество, на первый взгляд, трудно связанных друг с другом параметров.

Все вышесказанное предполагает целесообразность применения для обнаружения инсайдеров в КС подхода, основанного на использовании экспертных правил, методов машинного обучения и обработки больших данных. Этим обуславливается актуальность темы диссертационного исследования.

Степень разработанности темы

Проблеме существования инсайдерской деятельности в КС было посвящено большое количество работ как отечественных ученых (П.Д. Зегжды, И.В. Котенко, А.В. Лукацкого, А.А. Молдовяна, В.Ю. Осипова, И.Б. Саенко и др.), так и зарубежных (S. Bellovin, C. Cheh, M. Collins, F. Kammüller, Y. Shuang-Hua, X. Wang и др.). Однако, несмотря на сделанный учеными существенный задел, проблема обнаружения инсайдеров в КС не может считаться разрешенной и требует проведения новых исследований, что и осуществлено в данной работе.

Цели и задачи. Основной целью диссертационной работы является повышение защищенности КС за счет усовершенствования моделей, алгоритмов и методики обнаружения инсайдеров в КС с использованием комбинирования экспертных правил, методов машинного обучения и способов обработки больших данных.

Для достижения данной цели в диссертационной работе поставлены и решены следующие задачи:

- 1) анализ существующих подходов к обнаружению инсайдеров в КС, моделей, методик и алгоритмов обнаружения инсайдеров в КС на основе методов машинного обучения и обработки больших данных;
- 2) разработка модели представления больших данных об инсайдерских атаках в формате NoSQL (включая модель инсайдера);
- 3) разработка алгоритма обнаружения инсайдеров в КС, основанного на экспертных правилах;
- 4) разработка модели и алгоритмов комбинированного применения экспертных правил и методов машинного обучения в интересах обнаружения инсайдерских атак;
- 5) разработка методики обнаружения инсайдеров в КС с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных;
- 6) построение архитектуры и реализация программного комплекса системы обнаружения инсайдеров в КС на базе предложенной методики, настройка алгоритма

на основе методов машинного обучения с помощью набора данных, характеризующих действия инсайдеров по заданному множеству сценариев атак, и экспериментальная оценка разработанной методики системы обнаружения инсайдеров в КС.

Объектом исследования являются КС, в которых возможно наличие инсайдеров и атаки инсайдеров на КС.

Предметом исследования являются модели, методики и алгоритмы обнаружения инсайдеров в КС с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных.

Разработка модельно-методического аппарата для обнаружения инсайдеров в КС на основе комбинированного использования экспертных правил, методов машинного обучения и обработки больших данных определяет **научную задачу исследования**.

Теоретическая и практическая значимость работы. Теоретическая значимость диссертационной работы определяется ее вкладом в дальнейшее развитие теории и методов информационной безопасности, что проявляется в следующих аспектах: расширены классы атрибутов, необходимых для обнаружения инсайдеров; предложен новый подход к комбинированию двух классов алгоритмов, основанных на экспертных правилах и на методах машинного обучения, для решения задачи обнаружения инсайдеров в КС; методика обнаружения инсайдеров реализует последовательность операций, необходимых для решения задачи обнаружения инсайдеров, основывается на модели в формате NoSQL, алгоритмах, основанных на экспертных правилах, а также алгоритмах, основанных на методах машинного обучения; архитектура программного комплекса системы обнаружения инсайдеров реализует совокупность компонентов, их взаимосвязь, процедуру их выполнения и программную реализацию для решения задачи обнаружения инсайдеров в КС; архитектура основана на модели в формате NoSQL, алгоритмах, основанных на экспертных правилах и методах машинного обучения, предложенных в диссертации.

Практическая значимость диссертационной работы заключается в следующем:

- модель представления больших данных об инсайдерских атаках является основой для формализации данных и знаний о пользователях, устройствах, приложениях и сервисах в КС;

- модель и алгоритмы комбинированного применения экспертных правил и методов машинного обучения позволяют оперировать большими объемами данных и выявлять инсайдеров для достижения наилучших показателей эффективности; произведена настройка алгоритмов на основе методов машинного обучения по типовым сценариям инсайдеров в КС; обосновано комбинированное применение алгоритмов обнаружения инсайдеров;

- методика обнаружения инсайдеров повышает эффективность обнаружения внутренних нарушителей в КС (оперативность повышается за счет использования методов обработки больших данных; результативность – за счет совместного использования алгоритмов на основе экспертных правил и методах машинного обучения, ресурс-экономность – за счет новых высокотехнологичных программно-аппаратных решений);

- архитектура и программная реализация системы способствует эффективному обнаружению инсайдеров в КС с использованием предложенной методики обнаружения инсайдеров, обеспечивающей комбинированное применение технологий обработки больших данных, экспертных правил и методов машинного обучения.

Методология и методы исследований. Для решения поставленных задач использовались как классические, так и современные методы исследования, а именно: системный, причинно-следственный и сравнительный анализ был применен в равной степени для получения практически всех основных научных результатов; теория вероятностей и теория множеств применялись в интересах формирования математической модели представления больших данных для обнаружения инсайдеров; сбор, систематизация и анализ научно-технической информации предметной области, а также функциональный и структурный синтез позволили создать модель и комплекс алгоритмов обнаружения инсайдеров; методы машинного обучения явились центральным звеном одного из алгоритмов комплекса обнаружения инсайдеров; методы обработки больших данных легли в основу методики обнаружения инсайдеров, затрагивая, тем самым, все остальные полученные результаты; для практической оценки методики и программной реализации системы обнаружения инсайдеров был проведен компьютерный эксперимент на базе имитационного моделирования; основой программной реализации системы обнаружения инсайдеров послужила общая методология программирования.

Положения, выносимые на защиту. Основными положениями, выносимыми на защиту, являются:

1. Модель представления больших данных об инсайдерских атаках в формате NoSQL, обеспечивающая хранение и анализ признаков пользователей в компьютерных сетях в различные моменты времени.

2. Модель и алгоритмы комбинированного применения экспертных правил и методов машинного обучения в интересах обнаружения инсайдерских атак.

3. Методика обнаружения инсайдеров в компьютерных сетях с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных.

4. Архитектура и программная реализация системы обнаружения инсайдеров в компьютерных сетях с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных.

Научная новизна результатов диссертационной работы состоит в следующем:

1. Модель представления больших данных об инсайдерских атаках в формате NoSQL отличается от существующих возможностью обеспечения хранения и анализа признаков пользователей, полученных на базе UBA/UEBA-аналитики и характеризующих потенциальную инсайдерскую деятельность в компьютерных сетях, а также возможностью учета динамики изменения этих признаков.

2. Модель и алгоритмы комбинированного применения экспертных правил и методов машинного обучения в интересах обнаружения инсайдерских атак отличаются от существующих применением комплексного подхода к решению задачи обнаружения инсайдеров с учетом признаков и свойств пользователей, устройств, приложений, сервисов, включая параметр времени.

3. Методика обнаружения инсайдеров отличается от существующих использованием предложенной модели представления больших данных об инсайдерских атаках, а также предложенных модели и алгоритмов комбинированного применения экспертных правил, методов машинного обучения и обработки больших данных.

4. Архитектура и программная реализация системы обнаружения инсайдеров в компьютерных сетях отличается от известных архитектур и программных средств использованием предложенной методики обнаружения инсайдеров, обеспечивающей

комбинированное применение технологий обработки больших данных, экспертных правил и методов машинного обучения.

Реализация результатов работы. Отраженные в диссертационной работе исследования проведены в рамках ФЦП 2019-2020 гг. в соответствии с соглашением № 05.607.21.0322 (идентификатор RFMEFI60719X0322) с Минобрнауки России по теме: «Разработка методов, моделей, алгоритмов и программных средств, основанных на выявлении отклонений в эвристиках трафика сверхвысоких объемов, для обнаружения сетевых атак и защиты от них». Полученные результаты внедрены в учебный процесс СПбГУТ (учебные курсы: «Безопасность компьютерных сетей», «Безопасность беспроводных локальных сетей») и СПбГУТПД (учебные курсы: «Комплексная защита информации на предприятии», «Технологии и методы программирования»), применяются в рабочем процессе Роскомнадзора по Северо-Западному федеральному округу, компании ООО «Фаст Лейн». Результаты диссертационного исследования представлены в заявке, победившей на конкурсе субсидий молодым ученым, молодым кандидатам наук вузов, отраслевых и академических институтов, расположенных на территории Санкт-Петербурга, в 2019 г.

Обоснованность и достоверность полученных результатов обеспечивается за счет тщательного анализа состояния исследований предметной области, подтверждается согласованностью результатов с экспериментальными оценками, успешной апробацией основных теоретических положений диссертации на ряде научных конференций всероссийского и международного уровня, а также публикацией основных научных результатов в ведущих рецензируемых научных изданиях.

Апробация результатов работы. Основные положения и результаты работы докладывались на научных конференциях: международной конференции по интеллектуальным распределенным вычислениям IDC-2019 (Санкт-Петербург, 2019), международной конференции IEEE SMARTWORLD ATC-2017 (Сан-Франциско, 2017); Санкт-Петербургской межрегиональной конференции «Информационная безопасность регионов России» (Санкт-Петербург, 2015, 2017, 2019), Международной научно-технической и научно-методической конференции «Актуальные проблемы инфотелекоммуникаций в науке и образовании» в СПбГУТ (Санкт-Петербург, 2015–2019); XV-й Санкт-Петербургской международной конференции «Региональная информатика» (Санкт-Петербург, 2019); Российской мультikonференции по проблемам управления «Информационные технологии в управлении» (Санкт-Петербург, 2016).

Личный вклад. Все результаты, представленные в диссертационной работе, получены лично автором в процессе выполнения научно-исследовательской деятельности.

Публикации. По материалам диссертационной работы опубликовано 40 работ, в том числе 9 – в рецензируемых изданиях из перечня ВАК («Вопросы кибербезопасности», «Защита информации. Инсайд», «Труды СПИИРАН», «Труды учебных заведений связи»), 2 – в изданиях, индексируемых в международных базах Scopus и Web of Science, получено 3 свидетельства о государственной регистрации программ для ЭВМ

Структура и объем диссертационной работы. Диссертационная работа включает введение, три главы, заключение, список литературы (190 наименований) и 2 приложения. Объем работы – 206 страниц машинописного текста; включает 35 рисунков и 13 таблиц.

Содержание работы

Первая глава диссертации посвящена анализу проблемы обнаружения инсайдеров в КС. Проведен анализ атак на КС. Выполнена постановка задачи исследования, которая включает в себя разработку модельно-методического аппарата для обнаружения инсайдеров в КС на основе комбинированного использования экспертных правил, методов машинного обучения и обработки больших данных. Установлены место и роль задачи обнаружения инсайдеров в КС в общем цикле обработки информации в SIEM-системе.

Определено множество функциональных и нефункциональных свойств обнаружения инсайдеров и требований к системе обнаружения инсайдеров. Выделены следующие свойства обнаружения инсайдеров: *своевременность*, как способность обнаруживать инсайдеров в установленный промежуток времени; *обоснованность*, как доля обнаруженных инсайдеров по сравнению с их реальным наличием в сети; *ресурсопотребление* как характеристики программных и аппаратных средств – количество хостов (h); средний сетевой трафик (n); объем занимаемого пространства на SSD/HDD (v); средняя нагрузка на CPU (c); средняя загрузка памяти (m), необходимых системе для обнаружения инсайдеров;

Для вычисления обоснованности введены меры качества: количество пользовательских сессий, определённых как инсайдерские и являющиеся таковыми (TP), количество пользовательских сессий, определённых как инсайдерские, но не являющиеся таковыми (FP), количество пользовательских сессий, не определённых как инсайдерские и являющиеся таковыми (TN), количество пользовательских сессий, не определённых как инсайдерские, но не являющиеся таковыми (FN), а также получаемые из них полнота, точность, аккуратность, ошибка, F-мера.

Сформулирована задача исследования. Она заключается в разработке: (1) модели представления больших данных об инсайдерских атаках в формате NoSQL; (2) модели и алгоритмов комбинированного применения экспертных правил (RB-алгоритм, от англ. Rule-Based – на базе правил) и методов машинного обучения (ML-алгоритм, от англ. Machine Learning – машинное обучение) в интересах обнаружения инсайдерских атак; (3) методики обнаружения инсайдеров в КС с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных; (4) архитектуры и программной реализации системы обнаружения инсайдеров в КС с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных.

Сформулирована цель исследования – повышение защищенности КС от внутренних атак. В диссертации показатель защищенности определяется через показатель обоснованности (F-мера) с учетом ограничений других показателей обоснованности (полноты, точности, аккуратности, ошибки), а также с учетом требований к своевременности и ресурсопотреблению.

Во второй главе представлена разработанная модель представления больших данных об инсайдерских атаках, а также модель и алгоритмы комбинированного применения экспертных правил и методов машинного обучения в интересах обнаружения инсайдерских атак.

Формальный вид модели представления больших данных об инсайдерских атаках имеет следующий вид:

$$M = \langle A, I \rangle, \quad (1)$$

где A – элементы, представляющие собой атрибуты поведения пользователя, I – модель инсайдера и критерии, которые позволяют отнести текущего пользователя к категории инсайдеров.

Выделенные атрибуты поведения пользователей и их взаимосвязи формально имеют следующий вид:

$$A = \langle DataSources, Users, Data, Parser \rangle \quad (2)$$

Элементы, входящие в этот кортеж, включают в себя:

– $DataSources = \langle Netflow, Application, Scanner, Server, Device \rangle$ – источники данных, каждый элемент которых представляет собой соответственно сетевой поток, приложение, файл операционной системы, сканер, сервер, устройство;

– $Users = \bigcup_{i=1}^P User_i$ – пользователи, $User_i = \langle UserID_i, Attr_i, Sessions_i \rangle$ – триплет, представляющий собой соответственно идентификатор пользователя, атрибуты пользователя и соответствующие ему сессии;

– $Data = \{0,1\}^+ = \{0,1,00,01,10,11,000, \dots\}$ – данные, представляющие собой всевозможные битовые цепочки, хранящиеся на источниках данных;

– $Parser: Data \times DataSources \times Time \rightarrow Sessions$ – отображение, формирующее сессию из сырых данных в зависимости от типа источника этих данных и времени.

Модель инсайдера представлена в следующем виде:

$$I = \langle R, L, Q, G \rangle, \quad (3)$$

где R – критерии определения атрибутов инсайдера, L – уровни доступа, Q – квалификация инсайдера, G – цель инсайдера.

Источники данных содержат информацию о пользователях в сыром виде. Для преобразования этих данных в сессии применяется отображение $Parser$. Например, в случае сетевого потока данное отображение позволяет выделить ТСР-соединение или НТТР-сессию, а в случае приложения / файла сформировать характеристики сессии ОС, в рамках которой это приложение было установлено / файл был создан.

Если на одном из источников данных $datasrc$ были сгенерированы данные $data$ в определенный момент времени $time$, тогда идентификатор пользователя uid , в рамках одной из сессий которого были сгенерированы эти данные, определяется следующим образом:

$$uid \in \{userid \mid \langle userid, attr, sessions \rangle \in Users \wedge Parser(data, datasrc, time) \in sessions\} \quad (4)$$

Предложен алгоритм обнаружения инсайдеров, основанный на экспертных правилах (рисунок 1). Исходя из ограниченности применения алгоритма (так как правила, введенные экспертами, могут быть неполными и неточными), было произведено его комбинирование с другим, в качестве которого был выбран алгоритм, основанный на методах машинного обучения – таким образом получен их комплекс.

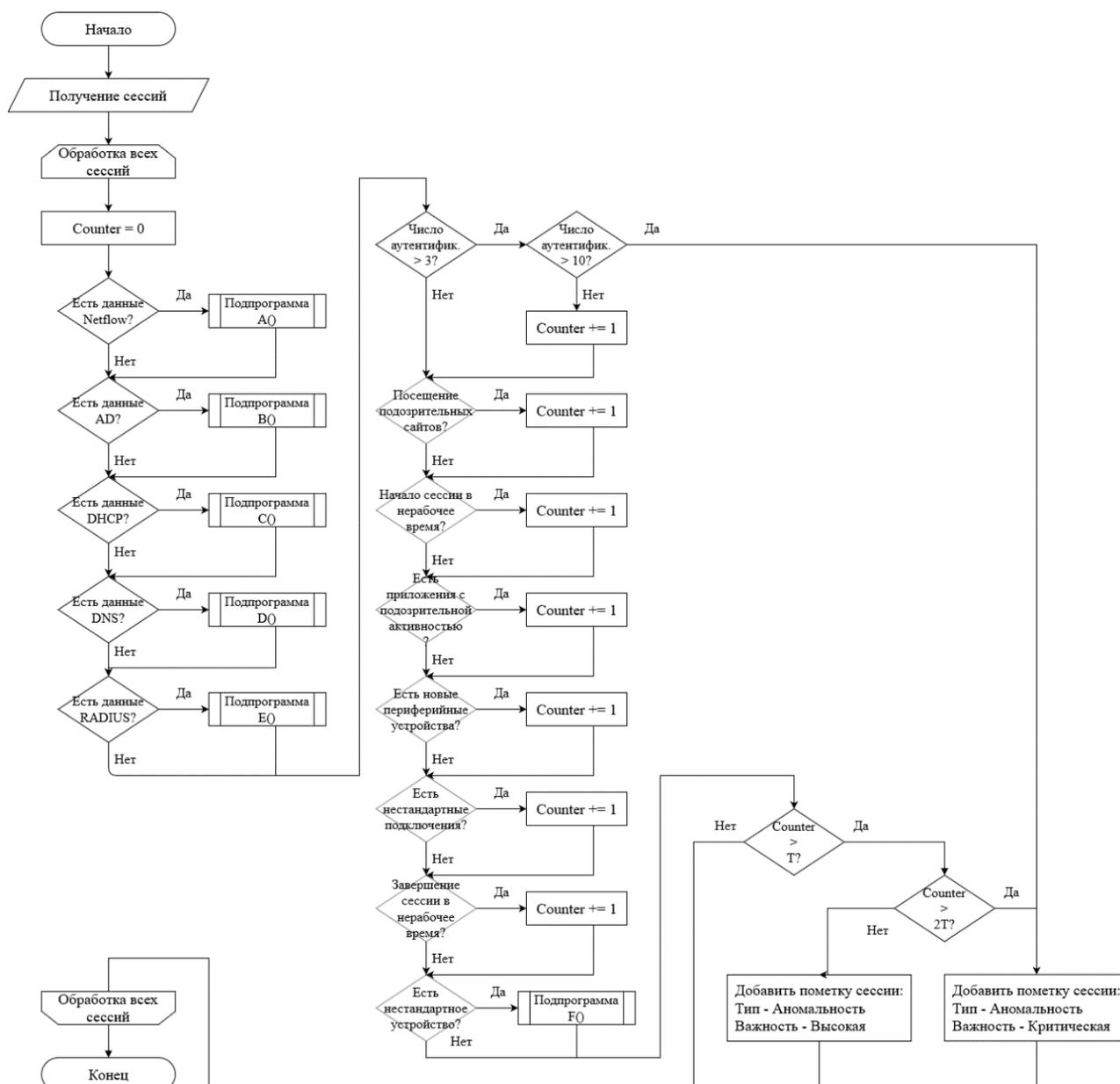


Рисунок 1 – Блок-схема алгоритма определения аномалий на основе экспертных правил

Функциональная структура комплекса алгоритмов (K_A) обнаружения инсайдеров в КС состоит из двух алгоритмов (A_1 и A_2), скомбинированных одним из способов. Формальная запись комплекса имеет следующий вид:

$$\left\{ \begin{array}{l} K_A = \{A_1 \otimes A_2\} \\ \otimes \in \{I, II, V, \Lambda\} \end{array} \right. \quad (5)$$

где \otimes – операции комбинирования, I – результат работы комплекса включает в себя инсайдеров, обнаруженных только первым из алгоритмов, II – результат работы комплекса включает в себя инсайдеров, обнаруженных только вторым из алгоритмов, V – результат работы комплекса включает в себя инсайдеров, обнаруженных любым из алгоритмов, Λ – результат работы комплекса включает в себя инсайдеров, обнаруженных обоими алгоритмами одновременно.

Комбинированное применение алгоритмов может быть представлено в виде модели, представленной в графическом виде на рисунке 2.

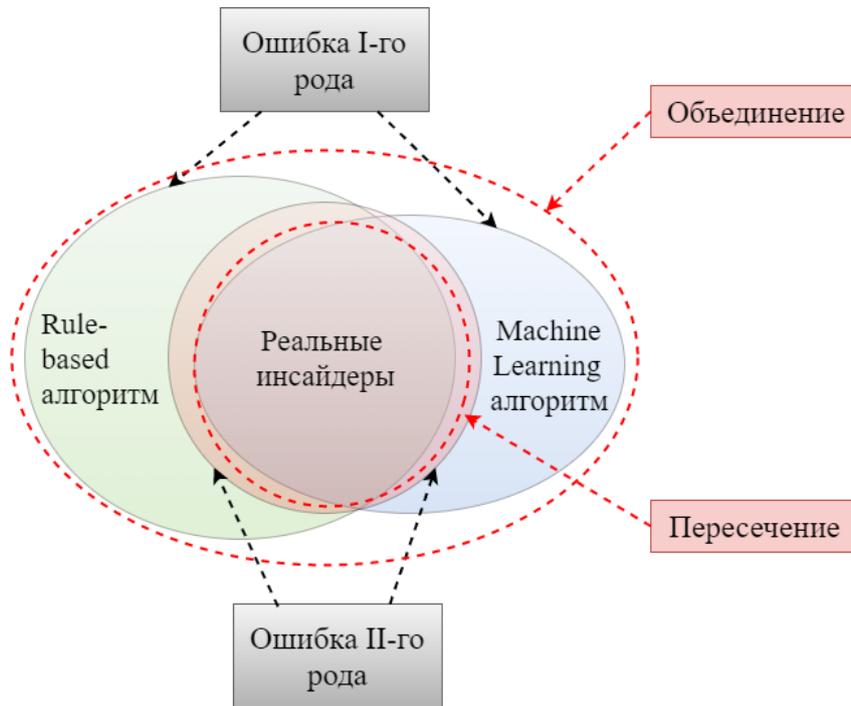


Рисунок 2 – Графическая интерпретация модели комбинирования алгоритмов обнаружения инсайдеров в КС

Модель отражает взаимосвязь следующих сущностей, представленных на рисунке с помощью 3-х эллиптических областей, ассоциированных с инсайдерами: реальные инсайдеры в КС (I_0), инсайдеры, обнаруженные алгоритмом на основе экспертных правил (I_{RB}), инсайдеры, обнаруженные алгоритмом на основе методов машинного обучения – для некоторого классификатора (I_{ML}).

В рамках формальной записи комплекса алгоритмов модель может быть описана следующим образом (для каждой комбинации):

$$\text{а) для } I = I_{RB}: \begin{cases} I_{TP} = I_0 \wedge I_{RB} \\ I_{TN} = \overline{I_0 \vee I_{RB}} \\ I_{FP} = I_{RB} \setminus I_0 \\ I_{FN} = I_0 \setminus I_{RB} \end{cases} \quad (6)$$

$$\text{б) для } I = I_{ML}: \begin{cases} I_{TP} = I_0 \wedge I_{ML} \\ I_{TN} = \overline{I_0 \vee I_{ML}} \\ I_{FP} = I_{ML} \setminus I_0 \\ I_{FN} = I_0 \setminus I_{ML} \end{cases} \quad (7)$$

$$\text{в) для } I = I_{RB} \wedge I_{ML}: \begin{cases} I_{TP} = I_0 \wedge (I_{RB} \wedge I_{ML}) \\ I_{TN} = \overline{I_0 \vee (I_{RB} \wedge I_{ML})} \\ I_{FP} = (I_{RB} \wedge I_{ML}) \setminus I_0 \\ I_{FN} = I_0 \setminus (I_{RB} \wedge I_{ML}) \end{cases} \quad (8)$$

$$\text{г) для } I = I_{RB} \vee I_{ML}: \begin{cases} I_{TP} = I_0 \wedge (I_{RB} \vee I_{ML}) \\ I_{TN} = I_0 \vee (I_{RB} \vee I_{ML}) \\ I_{FP} = (I_{RB} \vee I_{ML}) \setminus I_0 \\ I_{FN} = I_0 \setminus (I_{RB} \vee I_{ML}) \end{cases}, \quad (9)$$

где $I_{TP}, I_{TN}, I_{FP}, I_{FN}$ – множества, соответствующие каждой из мер. Очевидно, что для

идеального случая работы комплекса (то есть для $I = I_0$): $\begin{cases} I_{TP} = I_0 \\ I_{TN} = \bar{I}_0 \\ I_{FP} = \emptyset \\ I_{FN} = \emptyset \end{cases}$, где $\bar{I}_0 = U - I_0$ –

все законные пользователи (для всего множества пользователей КС – U).

Для комплексирования алгоритмов обоснована необходимость выбора способа комбинирования результатов их работы: объединение – результат работы комплекса включает в себя инсайдеров, обнаруженных любым из алгоритмов; пересечение – результат работы комплекса включает в себя инсайдеров, обнаруженных любым из алгоритмов; только первый – результат работы комплекса включает в себя инсайдеров, обнаруженных только первым алгоритмом; только второй – результат работы комплекса включает в себя инсайдеров, обнаруженных только вторым алгоритмом. Также показана необходимость выбора одного из классификаторов, задающих конкретный метод машинного обучения (DT, NB, k-NN, SVM), или их комбинации (голосование большинством (PV), взвешенное голосование (WV), мягкое голосование (SV), Adaboost).

Показано, что всего возможно 25 пар вариантов комбинаций алгоритмов, результат каждой из которых может иметь собственные показатели качеств, определяющие свой выбор при использовании в методике.

Предложенная схема комплекса алгоритмов обнаружения инсайдеров представлена на рисунке 3.

Схема работы комплекса алгоритмов состоит из следующих этапов:

- 1) ввод данных о сетевой активности КС, которые затем передаются на вход алгоритмов на базе экспертных правил и методов машинного обучения;
- 2) параллельный запуск алгоритма на базе экспертных правил и вариаций алгоритма на базе методов машинного обучения – для каждого из классификаторов (1-я вариация) и комбинирование результатов обоих алгоритмов путем различных способов (2-я вариация);
- 3) вывод множества результатов работы алгоритмов по каждой из вариаций.

В качестве типовых сценариев атак инсайдеров были выбраны следующие семь сценариев, поделенные на две группы и представляющие наибольшую распространенность среди типовых атак в КС:

- группа 1 – сканирование внутренних MAC/IP-адресов и TCP/UDP портов, отказ в обслуживании с помощью «SYN-флуд», неудачная попытка входа в сервис через HTTP/FTP, аномальное получение чрезмерно большого количества данных с внутренних ресурсов, аномальная отправка чрезмерно большого количества данных вне организации;

- группа 2 – кража информации сотрудником в нерабочее время и их вынос через собственный ноутбук, скрытый сбор и отправка сотрудником в рабочее время на внешний IP чужих персональных данных.

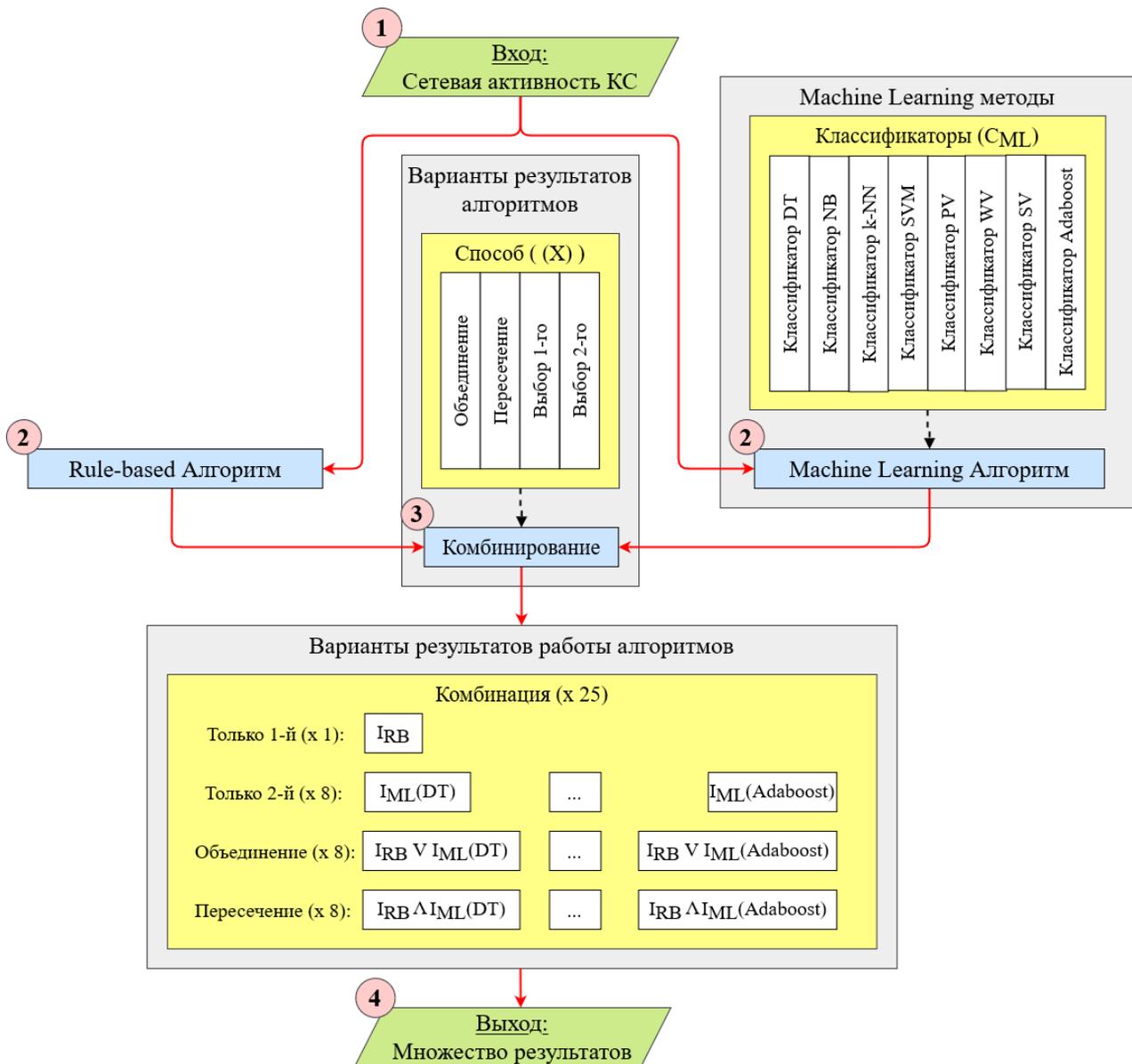


Рисунок 3 – Схема работы комплекса алгоритмов обнаружения инсайдеров

В третьей главе приведены методика, архитектура и программная реализация системы обнаружения инсайдеров в КС. Представлены результаты экспериментов и сравнение предложенной методики с существующими аналогами, а также предложения по применению разработанного модельно-методического аппарата для обнаружения инсайдеров в КС.

Методика обнаружения инсайдеров в КС с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных определяет основные этапы и шаги использования разработанных моделей и алгоритмов. Методика состоит из трех этапов (рисунок 4):

- 1) сбор информации;
- 2) применение алгоритмов обнаружения инсайдеров;
- 3) анализ выходных данных.

Архитектура системы обнаружения инсайдеров включает три уровня обработки информации: (1) уровень сети и данных (данные от источников, специальные агенты сбора и отправки данных); (2) уровень предварительной обработки информации и со-

бытий безопасности (компонент балансировки и разделения нагрузки, компоненты корреляции и индексирования (MapReduce), хранилище (NoSQL база данных), компонент работы с базой данных, компонент предварительного анализа информации и событий безопасности); (3) уровень аналитической обработки информации и событий безопасности (компонент поиска аномалий и инцидентов безопасности, компонент визуализации).

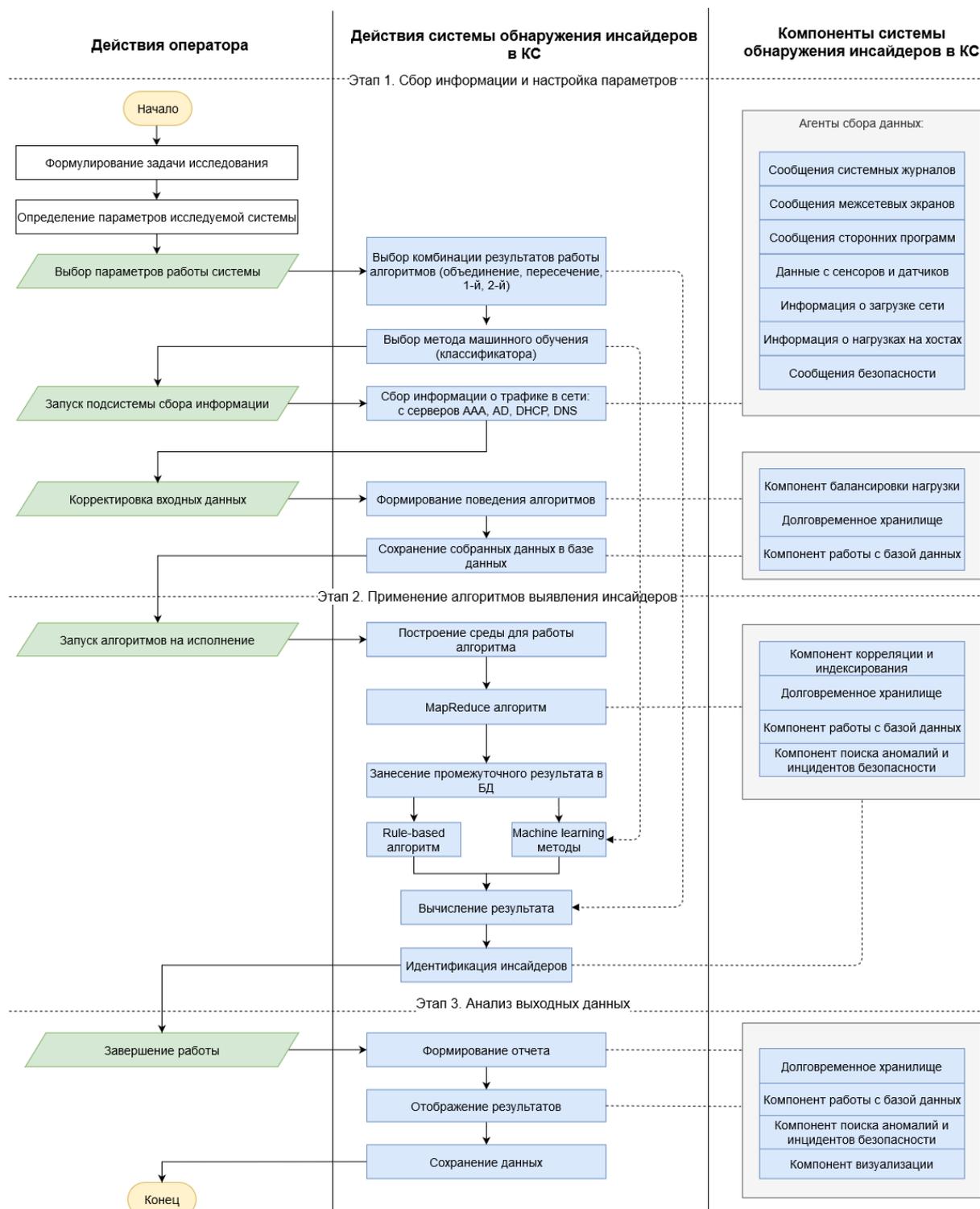


Рисунок 4 – Структура методики обнаружения инсайдеров в КС

Разработана функциональная схема проведения экспериментов по оценке обоснованности и оперативности разработанной системы обнаружения, состоящая из 3-х этапов: (1) подготовка стенда, (2) настройка алгоритмов и (3) проведение измерений.

Настройка методов машинного обучения производилась следующим образом. В качестве обучающей выборки были взяты 100 000 записей, из числа которых 70 000 (т.е. 70%) соответствовали нормальному трафику пользователей, а 30 000 (т.е. 30%) были связаны с инсайдерскими атаками. Эксперимент по обнаружению инсайдеров проводился на выборке данных, состоящей из 20 000 записей (т.е. 20% от обучающей выборки), число относящихся к инсайдерской деятельности составляло 10 000 (т.е. 50%).

С использованием экспериментального стенда были оценены показатели обоснованности каждого из вариантов комбинирования алгоритмов для каждой группы сценариев. В таблице 1 представлены отдельные значения показателей обоснованности для второй группы сценариев (Клс. ML – тип классификатора машинного обучения). Анализ результатов показал, что наилучшей с точки зрения F-меры является вариация, использующая пересечение результатов работы алгоритма на основе экспертных правил и метода машинного обучения при использовании SV (строка 24). Обоснованием результата является факт того, что каждый из алгоритмов имеет достаточное количество как пропусков инсайдеров, так и ошибочных срабатываний. По сравнению с работой каждого из алгоритмов, объединение их результатов улучшает обнаружение инсайдеров, но в то же время увеличивает число ложных срабатываний. Пересечение же наоборот, ухудшает обнаружение инсайдеров, но уменьшает число ложных срабатываний. Однако, усредненный эффект от работы пересечения (то есть совместная оценка с позиции полноты и точности) оказывается выше; при этом метод SV также оказывается с этой позиции лучшим среди остальных методов.

В качестве альтернативных систем, для сравнения с разработанной, были выбраны продукты Cisco StealthWatch и PacketFence.

Таблица 1 – Сравнение значений показателей обоснованности

Вариации			Меры								
№	RB+ML	Клс. ML	TP	TN	FP	FN	r	p	a	e	f
1	I_{RB}	–	8969	7781	2219	1031	0.90	0.80	0.84	0.16	0.85
2	I_{ML}	DT	9201	8111	1889	799	0.92	0.83	0.87	0.13	0.87
...											
9	I_{ML}	Adaboost	9362	8215	1785	638	0.94	0.84	0.88	0.12	0.89
10	$I_{RB} \vee I_{ML}$	DT	9822	7538	2462	178	0.98	0.80	0.87	0.13	0.88
...											
17	$I_{RB} \vee I_{ML}$	Adaboost	9901	7746	2254	99	0.99	0.81	0.88	0.12	0.89
18	$I_{RB} \wedge I_{ML}$	DT	8724	9607	393	1276	0.87	0.96	0.92	0.08	0.91
...											
24	$I_{RB} \wedge I_{ML}$	SV	8968	9989	11	1032	0.90	1.00	0.95	0.05	0.95
25	$I_{RB} \wedge I_{ML}$	Adaboost	8885	9809	191	1115	0.89	0.98	0.93	0.07	0.93

Результаты сравнения показали преимущество системы обнаружения инсайдеров в КС с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных по показателям обоснованности и оперативности, а также равенство по показателю ресурсопотребления.

Таким образом, можно утверждать о достижении более высокой эффективности разработанной системы, что доказывает реализацию итоговой цели исследования – относительного повышения защищенности КС за счет усовершенствования методик, моделей и алгоритмов обнаружения инсайдеров КС с использованием методов машинного обучения и обработки больших данных на $\frac{0.95-0.92}{0.92} = 0.033$ (3.3%), где 0.95 – полученное значение F-меры для метода SV (таблица 1, строка 24); 0.92 – максимальное значение F-меры, полученное в существующих исследованиях.

Произведен сравнительный анализ разработанной методики с существующими методиками по используемым функциональным возможностям: А – использование методов обработки больших данных, Б – распределенная обработка, В – использование методов машинного обучения, Г – использование экспертных правил, Д – учет гетерогенности признаков.

Результаты приведены в таблице 2 (используются следующие обозначения и баллы: «+» – наличие параметра в работе, 1 балл; «+/-» – частичное соответствии параметру, 0.5 балл; «-» – его отсутствии, 0 балл).

Таблица 2 – Сравнительный анализ разработанной методики с существующими аналогами

Методика обнаружения инсайдеров	Учитываемые параметры					Оценка
	А	Б	В	Г	Д	
Anagi Gamachchi, Li Sun and Serdar Boztas, 2017	+/-	+	+	-	+	3.5
Kara Nance, Raffael Marty, 2011	+/-	-	+/-	+	+	4
You Chen, Steve Nyemba, Wen Zhang, Bradley Malin, 2012	+	+/-	+	-	+	3.5
Chen T. и др., 2015	+/-	-	+	-	+/-	2
Ignacio J. Martinez-Moyano, Eliot Rich, Stephen Conrad, David F. Andersen, Thomas R. Stewart, 2008	+	+/-	-	+	+	3.5
Разработанная методика	+	+	+	+	+	5

Анализ результатов сравнения позволяет утверждать о преимуществе разработанной методики (5 баллов) перед аналогами (2, 3.5 и 4 балла).

Предложены возможные пути применения системы обнаружения инсайдеров.

В заключении представлена итоговая оценка проделанной работы, приведены основные результаты исследования и описаны перспективы дальнейшего исследования в рамках темы.

Основные выводы и результаты работы

В диссертационной работе в целях повышения защищенности КС решена задача разработки модельно-методического аппарата для обнаружения инсайдеров в сети на основе комбинированного использования экспертных правил, методов машинного обучения и обработки больших данных, имеющая важное значение для развития технологий в области информационной безопасности. При решении данной задачи были получены следующие результаты:

1. Модель представления больших данных об инсайдерских атаках в формате NoSQL, обеспечивающая хранение и анализ признаков пользователей в компьютерных сетях в различные моменты времени.

2. Модель и алгоритмы комбинированного применения экспертных правил и методов машинного обучения в интересах обнаружения инсайдерских атак.

3. Методика обнаружения инсайдеров в компьютерных сетях с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных.

4. Архитектура и программная реализация системы обнаружения инсайдеров в компьютерных сетях с использованием комбинирования экспертных правил, методов машинного обучения и обработки больших данных.

Рекомендации по применению разработанного модельно-методического аппарата для обнаружения инсайдеров в КС включают использование модели представления больших данных об инсайдерских атаках для наблюдения за законными пользователями в интересах повышения эффективности работы организации. Использование методики и алгоритмов обнаружения инсайдеров после необходимой доработки для анализа внешних сетевых атак и их последующей нейтрализации. Интеграция архитектуры системы обнаружения инсайдеров в КС в крупные SIEM-системы. Перспективы дальнейшей разработки темы заключаются в расширении сценариев инсайдеров в КС, обнаруживаемых предложенными алгоритмами, методикой и системой, интеграция механизмов автоматической нейтрализации инсайдерской деятельности, имеющей на основании работы алгоритмов критический уровень важности, расширение комплекса алгоритмов другими (такими, как отклонение от «сезонного поведения», применение дискретного вейвлет-анализа и др.), способствующими итоговому повышению качества и скорости обнаружения инсайдеров.

Полученные результаты соответствуют п. 3 «Методы, модели и средства выявления, идентификации и классификации угроз нарушения информационной безопасности объектов различного вида и класса» паспорта специальности 05.13.19 – «Методы и системы защиты информации, информационная безопасность».

Список основных публикаций по теме диссертации

В рецензируемых изданиях из Перечня ВАК:

1) Ушаков, И.А. Обнаружение инсайдеров в корпоративной компьютерной сети на основе технологий обработки больших данных. / И.А. Ушаков // Вестник Санкт-Петербургского государственного университета технологии и дизайна. Серия 1: Естественные и технические науки. – 2019. – № 4. – С. 38-43.

2) Ушаков, И.А. Масштабируемое honeypot-решение для обеспечения безопасности в корпоративных сетях / А.В. Красов, Р.Б. Петрив, Д.В. Сахаров, Н.Л. Сторожук, И.А. Ушаков // Труды учебных заведений связи. – 2019. – Т. 5. – №. 3. – С. 86-97.

3) Ушаков, И.А. Исследование модели сети ЦОД на основе политик Cisco ACI / Н.В. Савинов, К.А. Токарева, И.А. Ушаков, А.В. Красов, Д.В. Сахаров // Защита информации. Инсайд. – 2019. – № 4 (88). – С. 32-43.

4) Ушаков, И.А. Комплексный подход к обеспечению безопасности киберфизических систем на основе микроконтроллеров / И.В. Котенко, Д.С. Левшун, А.А. Чечулин, И.А. Ушаков, А.В. Красов // Вопросы кибербезопасности. – 2018. – № 3 (27). – С. 29-38.

5) Ушаков, И.А. Обеспечение безопасности передачи multicast-трафика в ip-сетях / А.В. Красов, Д.В. Сахаров, И.А. Ушаков, Е.П. Лосин // Защита информации. Инсайд. – 2017. – № 3 (75). – С. 34-42.

6) Ушаков, И.А. Гибридная модель базы данных NoSQL для анализа сетевого трафика . И.В. Котенко, И.А. Ушаков, Д.В. Пелёвин, А.Ю. Овраменко // Защита информации. Инсайд. – 2019. – № 1 (85). – С. 46-54.

7) Ушаков, И.А. Система сбора, хранения и обработки информации и событий безопасности на основе средств Elastic Stack / И.В. Котенко, А.А. Кулешов, И.А. Ушаков // Труды СПИИРАН. – 2017. – № 5 (54). – С. 5-34.

8) Ушаков, И.А. Технологии больших данных для мониторинга компьютерной безопасности / И.В. Котенко, И.А. Ушаков // Защита информации. Инсайд. – 2017. – № 3 (75) – С. 23-33.

9) Ушаков, И.А. Выявление инсайдеров в корпоративной сети: подход на базе UBA и UEBA / И.В. Котенко, И.А. Ушаков, Пелевин Д.В., Преображенский А.И., Овраменко А.Ю. // Защита информации. Инсайд. – 2019. – № 5 (89). – С. 2-11.

В зарубежных изданиях из баз данных Web of Science и Scopus:

10) Ushakov, I. Aggregation of Elastic Stack Instruments for Collecting, Storing and Processing of Security Information and Events / I. Kotenko, A. Kuleshov, I. Ushakov // 2017 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computed, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/CBDCOM/IOP/SCI) 2017. – Pp. 1-8.

11) Ushakov, I. Approach to Detection of Denial-of-Sleep Attacks in Wireless Sensor Networks on the base of Machine Learning / A. Balueva, V. Desnitsky, I. Ushakov. – Pp. 350-355.

Свидетельства о государственной регистрации программ для ЭВМ:

12) Ушаков И.А. Компонент предобработки трафика в корпоративной компьютерной сети с использованием алгоритма Map Reduce в Hadoop кластере: свидетельство о государственной регистрации программы для ЭВМ / И.А. Ушаков, И.В. Котенко, А.Ю. Овраменко. – Рег. № 2019666737. – 13.12.2019.

13) Ушаков И.А. Система обнаружения инсайдеров в корпоративной компьютерной сети с использованием технологий машинного обучения: свидетельство о государственной регистрации программы для ЭВМ / И.А. Ушаков, И.В. Котенко, Ю.В. Твердохлебова. – Рег. № 2019666738. – 13.12.2019.

14) Ушаков И.А. Система обнаружения инсайдера в корпоративной компьютерной сети, используя алгоритмы, основанные на экспертных правилах: свидетельство о государственной регистрации программы для ЭВМ / И.А. Ушаков, И.В. Котенко, Д.В. Пелёвин – Рег. № 2019666959. – 17.12.2019.