

**МИНИСТЕРСТВО ОБРАЗОВАНИЯ РЕСПУБЛИКИ БЕЛАРУСЬ**

**БЕЛОРУССКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ**

**Факультет прикладной математики и информатики  
Кафедра математического моделирования и анализа данных**

Дорофеев Герман Сергеевич

**ОБНАРУЖЕНИЕ ВРЕДНОСНЫХ DOCX-ФАЙЛОВ С  
ПОМОЩЬЮ МЕТОДОВ МАШИННОГО ОБУЧЕНИЯ**

Аннотация к дипломной работе

**Научный руководитель:**

кандидат физико-математических наук,

М.С. Абрамович

---

Минск, 2022

## РЕФЕРАТ

**Дипломная работа:** 50 стр., 19 рис., 4 таблицы, 12 источников, 3 приложения.

**Ключевые слова:** docx, malware, macro, machine learning, classification.

**Объект исследования:** docx-файлы.

**Цель работы:** исследовать применение методов машинного обучения в обнаружении вредоносных docx-файлов, разработать приложение.

**Методы исследования:** методы машинного обучения.

**Результат:** изучено применение методов машинного обучения в обнаружении вредоносных docx-файлов, разработано приложение.

**Область применения:** сфера информационной безопасности.

## РЭФЕРАТ

**Дыпломная работа:** 50 стр., 19 рыс., 4 табл, 12 крыніц, 3 прыкл.

**Ключавыя словы:** docx, malware, macro, machinelearning, classification.

**Аб'ект даследавання:** docx-файлы.

**Мэты працы:**

даследаваць ужыванне метадаў машыннага навучання ў выяўленні шкодных docx-файлаў, распрацаваць праграму.

**Метады даследавання:** метады машыннага навучання.

**Вынікі:** вывучана ўжыванне метадаў машыннага навучання ў выяўленні шкодных docx-файлаў, распрацавана праграма.

**Вобласць ужывання:** абарона інфармацыі.

## ABSTRACT

**Diploma thesis:** 50 pages, 19 figures, 4 tables, 12 sources, 3 attachments.

**Keywords:** docx, malware, macro, machine learning, classification.

**Object of research:** docx-files.

**Purpose:** study the application of machine learning methods for detection of malicious docx-files, develop the corresponding software application.

**Methods:** machine learning methods.

**Result:** the application of machine learning methods for detection of malicious docx-files has been studied. The corresponding software application has been developed.

**Scope:** information security field.

## ВВЕДЕНИЕ

Вредоносный макрос (также называемый «макровирус») - это код, который использует функциональность документов Office (обычно это Excel и Word) для выполнения злонамеренных действий против системы при открытии файла. Этот тип вредоносного ПО был очень популярен в конце 90-х - начале 2000-х годов. С 2014-го года макровирусы начали использоваться снова, на этот раз как метод распространения других вредоносных программ. Так или иначе макровирусы продолжают представлять угрозу для пользователя.

В этой работе изучается возможность улучшения обнаружения макровирусов с помощью методов машинного обучения, применяемых к свойствам кода.

Изначально макросы добавляли дополнительную функциональность к документам, предоставляя им динамические свойства, которые позволяют, например, выполнять действия над набором ячеек в документе Excel или встраивать объекты мультимедиа в файлах Word. Но к концу 90-х они стали использоваться злоумышленниками для выполнения кода в операционной системе. Злоумышленники программировали макросы, наделяя их дополнительными свойствами для выполнения вредоносных действий (например, загрузка и запуск исполняемых файлов).

В конце 1990-х одним из самых распространенных и вредных вирусов была «Мелисса» [1]. Традиционно, наиболее популярным способом распространения этого типа вредоносных программ была (и остается) электронная почта. Жертва получает электронное письмо с вложением, и при его открытии макрос, содержащийся в документе, исполняется и заражает операционную систему. Однако в последующие годы Microsoft предоставили механизм, встроенный в пакет Office, предотвращающий автоматическое выполнение макросов. Из-за этого данный вирус потерял свою актуальность. Этому также поспособствовало существование других, более прямых методов, не зависящих от конфигурации системы Office (например, использование уязвимостей). Однако в 2014 году Microsoft заявили, что макровирусы вновь стали использоваться злоумышленниками, на этот раз как метод распространения другого вредоносного ПО.

В 2014 и 2015 годах макровирусы использовались вполне успешно, несмотря на всеческие меры противодействия. Таким образом, даже спустя более чем 15 лет после своего появления макровирусы все еще остаются угрозой, и механизмы их обнаружения все еще необходимы.

Многие исследования были сосредоточены на обнаружении вредоносных PDF-документов [2]. Однако, исследование, фокусирующееся на файле PDF,

не является достаточно общим, чтобы его можно было применять к другим файловым структурам, таким как офисные документы MS. Файлы PDF значительно отличаются от файлов docx в двух важных пунктах. Во-первых, файлы PDF состоят из набора связанных объектов, а файлы docx - это архивные файлы, состоящие из папки файлов XML. Учитывая это, применение структурного подхода для обнаружения вредоносных PDF-файлов [3] на вредоносных файлах docx было бы неудачным. Во-вторых, для файлов PDF idocx используются разные методы атаки, и даже когда используется общая методика (например, встроенные файлы), атака запускается совершенно по-разному, так что одна и та же атака влияет на структуру файла по-разному в каждом случае. Следовательно, для обнаружения вредоносных файлов docx требуется другой подход, основанный на специальной файловой структуре этих файлов.

Цель данной работы - продемонстрировать, действительно ли обнаружение и анализ наиболее распространенных методов программирования вредоносных программ может способствовать правильной и автоматической классификации и возможности отличить документы, содержащих допустимые макросы, от тех, что содержат вредоносные.

Обнаружение макровирусов все же является специфической задачей антивирусных систем, поэтому вовсе не предполагается, что методы машинного обучения могут полностью заменить их. Это лишь иной взгляд на проблему. В этой работе планируется провести обзор наиболее известных и удачных способов и подходов к применению техник машинного обучения для обнаружения вредоносных docx-файлов.