

**«Санкт-Петербургский государственный электротехнический университет
«ЛЭТИ» им. В.И. Ульянова (Ленина)»
(СПбГЭТУ «ЛЭТИ»)**

Направление	09.04.02 – Информационные системы и технологии
Профиль	Распределенные вычислительные комплексы систем реального времени
Факультет	ФКТИ
Кафедра	ИС

К защите допустить

Зав. кафедрой

подпись

Цехановский В.В.

**ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА
МАГИСТРА**

**Тема: Методика анализа политик безопасности на основе
онтологического представления предметной области**

Студент

подпись

Кузнецов М.Д.

Руководитель

к.т.н., доцент
(Уч. степень, уч. звание)

подпись

Новикова Е.С.

Санкт-Петербург

2021

ЗАДАНИЕ НА ВЫПУСКНУЮ КВАЛИФИКАЦИОННУЮ РАБОТУ

Утверждаю
Зав. кафедрой ИС
Цехановский В.В.
_____ подпись
« ____ » _____ 2021 г.

Студент Кузнецов М.Д.

Группа 5374

Тема работы: Методика анализа политик безопасности на основе онтологического представления предметной области.

Место выполнения ВКР: Санкт-Петербургский государственный электротехнический университет «ЛЭТИ» им. В.И.Ульянова (Ленина)

Исходные данные (технические требования): —

Содержание ВКР: В разделе «Анализ предметной области» произведен анализ литературы и работ в данной области, в разделе «Проектирование» проведено проектирование инструментария для сбора датасета, «Реализация» приведены некоторые аспекты реализации.

Перечень отчетных материалов: пояснительная записка, иллюстрационный материал.

Дополнительные разделы: —

Дата выдачи задания

« ____ » _____ 2021 г.

Дата представления ВКР к защите

« ____ » _____ 2021 г.

Студент

_____ подпись

Кузнецов М.Д.

Руководитель

к.т.н., доцент
(Уч. степень, уч. звание)

_____ подпись

Новикова Е.С.

КАЛЕНДАРНЫЙ ПЛАН ВЫПОЛНЕНИЯ ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЫ

Утверждаю

Зав. кафедрой ИС

Цехановский В.В.

подпись

« ____ » _____ 2021 г.

Студент Кузнецов М.Д.

Группа 5374

Тема работы: Методика анализа политик безопасности на основе онтологического представления предметной области.

№ п\п	Наименование работ	Срок выполнения
1	Обзор литературы по теме работы	01.01 – 02.01
2	Анализ предметной области	01.01 – 02.01
3	Проектирование	01.01 – 02.01
4	Реализация	01.01 – 02.01
5	Оформление пояснительной записки	01.01 – 02.01
6	Оформление иллюстративного материала	01.01 – 02.01

Студент

подпись

Кузнецов М.Д.

Руководитель

к.т.н., доцент

(Уч. степень, уч. звание)

подпись

Новикова Е.С.

РЕФЕРАТ

Поясн. зап. 33 стр., 9 рис., 6 табл., 9 ист., 1 прил.

КЛЮЧЕВЫЕ СЛОВА И СЛОВСОЧЕТАНИЯ, НЕ БОЛЕЕ ПЯТНАДЦАТИ, ЧЕРЕЗ ЗАПЯТУЮ

Объектом исследования (разработки) являются указать объект исследования или разработки.

Цель работы – кратко (в 2–3 строки) указать цель работы.

Кратко (в 10–12) строк описать основное содержание работы, методы исследования (разработки), полученные результаты.

ABSTRACT

Briefly (10–15 lines) the content of graduating work is specified

ТЕРМИНЫ И ОПРЕДЕЛЕНИЯ

В настоящей пояснительной записке к выпускной квалификационной работе используются следующие термины с соответствующими определениями.

Датасет — набор данных для обучения моделей анализа естественного языка

Вэб-скрейпинг — это технология извлечения данных из вэб-страниц путем из скачивания и обработки

ПЕРЕЧЕНЬ СОКРАЩЕНИЙ И ОБОЗНАЧЕНИЙ

В настоящей пояснительной записке к выпускной квалификационной работе используются следующие сокращения и обозначения.

LSA — (от англ. Latent Semantic Search) латентно-семантический анализ

POS — (от англ. Part Of Speech) разложение по частям речи

TF-IDF — (от англ. Term Frequency – Inverse Document Frequency) разбор предложения по частям речи

СОДЕРЖАНИЕ

ТЕРМИНЫ И ОПРЕДЕЛЕНИЯ	3
ПЕРЕЧЕНЬ СОКРАЩЕНИЙ И ОБОЗНАЧЕНИЙ	4
ВВЕДЕНИЕ	7
1 Анализ предметной области	9
1.1 Статистические модели текстовых документов	9
1.2 Подход основанный на латентно-семантическом анализе текста	10
1.3 Подход основанный на латентном размещении Дирихле	15
1.4 Подход основанный на применении контекстно-свободных грамматик и синонимическом поиске	15
1.5 Выводы по строгим методам текстового анализа	20
1.6 Подход основанный на глубоком обучении	20
2 Проектирование инструментария	22
2.1 Техническое задание «Инструментарий для сбора датасета»	22
2.1.1 Основные положения технического задания	22
2.1.2 Скрейпер вэб-страниц	22
2.1.3 Очистка скачанных страниц политик	22
2.1.4 Инструмент разметки датасета	22
2.1.5 Фреймворк глубокого обучения	23
2.2 Первичная декомпозиция и планирование	23
2.3 Приложение вэб-скрейпер	24
2.3.1 Структура приложения вэб-скрейпера	24
2.3.2 Потенциальные проблемы	27
2.4 Инструмент разметки датасета	27
2.4.1 Объектное моделирование приложения	27
2.4.2 Реляционная модель приложения	28
2.4.3 Проектирование пользовательского интерфейса	28
2.5 Выбор средств разработки	29

3 Реализация инструментария	30
3.1 Полученные в результате реализации исходные коды	30
3.2 Полученный в результате реализации пользовательский интерфейс инструмента разметки.	30
3.3 Результаты решения поставленной задачи с помощью разработанного инструментария	30
ЗАКЛЮЧЕНИЕ.....	31
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ.....	32
ПРИЛОЖЕНИЕ А	33

ВВЕДЕНИЕ

В настоящее время персональные данные широко используются в предоставлении цифровых услуг, их персонализации и улучшении. Персональные данные – это любые данные, идентифицировать физическое лицо [1]. Таким образом, личные данные – это не только биометрическая информация, данные о состоянии здоровья человека, а также фото абонента услуги, местонахождение, информация о приложении и устройстве, которое можно использовать для отслеживания действий и информации о потребителе. Несколько массовых утечек персональных данных за последнее десятилетие привело к ужесточению законодательных требований во многих страны по всему миру. В настоящее время требуется, чтобы все личные данные обрабатывались надежно, а действия с ними были ясны и прозрачны для субъекта данных в соответствии с его или ее явно указанным согласием. Политики конфиденциальности поставщиков услуг, онлайн-согласие пользователей – единственные законные документы, сообщающие конечным пользователям, как собираются, обрабатываются их личные данные и передается третьим лицам. Однако в большинстве случаев эти документы написаны так, что их довольно сложно понять. И в настоящее время ситуация такова, что законодательные требования соблюдаются производителями продукции и поставщиками услуг, но конечные пользователи дают свое согласие без четкого понимания того, как обрабатываются их личные данные, потому что политика конфиденциальности и онлайн-согласие пользователя читаются редко из-за их сложности и низкой читабельности. Это ведет к ситуациям, когда конечные пользователи не знают о рисках для конфиденциальности связанных с использованием определенной услуги или устройства.

В настоящее время сфера информационных технологий является одной из самых быстрорастущих, в ней решается множество задач прикладного характера. Одной из прогрессивных технологий является технология глубокого

обучения. Полученные с помощью глубокого обучения модели способны решать широкий спектр прикладных задач. Однако, у данного подхода имеется существенный недостаток – необходимость датасета для обучения. Датасет играет критически важную роль в формировании результата в целом. Если качество датасета будет посредственным, либо он окажется недостаточно объемным, то поставленная задача не будет решена с адекватной точностью. В то же время сбор датасета – работа кропотливая и рутинная. Отличным решением является автоматизация данного процесса, возможно не полная, но частичная. Она абсолютно точно увеличит скорость сбора информации, что позволит за то же время собирать более объемные датасеты, и как следствие более точные модели будут получены после обучения.

Целью выпускной квалификационной работы является разработка программного пакета для автоматизированной формализации политик безопасности с последующей оценкой рисков для персональных данных пользователей сервисов. В ходе выполнения предполагается реализация инструментов для сбора датасета, который будет применен для обучения классификатора. Классификатор позволит автоматизированно формализовать политики безопасности. По формализованному описанию политик станет возможной оценка рисков для персональных данных пользователей.

1 Анализ предметной области

1.1 Статистические модели текстовых документов

Были протестированы две модели векторизованного представления текста – «мешок слов» и модель TF-IDF. Модель «мешок слов» представляет документ в виде матрицы, представленной на рисунке 1. Здесь слова каждого абзаца подсчитываются и сопоставляются с абзацами, в которых они встретились.

Amounts of words in paragraphs			
Paragraph	Word	Word 1	Word n
	Par.	Word 1	Word n
	Doc. 1	Count (w1, d1)	Count (wn, d1)
	Doc. n	Count (w1, dn)	Count (wn, dn)

Рисунок 1 – Bag-of-Words матрица

Модель TF-IDF представляет документ в виде матрицы, представленной на рисунке 2. Формула (1) показывает, как можно получить метрику TF-IDF.

$$tfidf(t, d, D) = \frac{n_t}{\sum_k n_k} \times \log \frac{|D|}{|\{d_i \in D : t \in d_i\}|}, \quad (1)$$

где t – термин или слово;

d – конкретный абзац;

D – набор абзацев.

Итак, модель TF-IDF придает больший вес словам которые использованы меньше раз. Это может быть полезно, когда тексты схожи с точки зрения

используемых слов, как в нашем случае, для политик безопасности.

		TF-IDF metrics		
Paragraph	Word	Word 1	...	Word n
	Par.			
	Doc. 1	tfidf (w1, d1, D)	...	tfidf (wn, d1, D)

	Doc. n	tfidf (w1, dn, D)	...	tfidf (wn, dn, D)

Рисунок 2 – Матрица TF-IDF

1.2 Подход основанный на латентно-семантическом анализе текста

Современные методы кластеризации текстов позволяют определять тематику текстов с высокой точностью. Однако большинство из этих методов принимают тексты с самыми разными темами как вход для алгоритмов. Но тексты со схожими тематиками можно проанализировать с помощью латентно-семантического анализа дважды: группировать тексты по темам один раз, и предоставить еще более детальное разделение их по подтемам во второй раз. Такой подход можно использовать для более точной классификации абзацев с точки зрения их характеристик и аспектов использования персональных данных. Следует отметить, что латентно-семантический поиск сильно зависит от глобального текстового контекста с потерями информации о локальных контекстных отношениях между словами. Были выделены девять тем конфиденциальности, которые следует сопоставить с абзацами согласия пользователя сайта – «сбор личных данных», «сбор данных третьими лицами», «управление личными данными», «механизмы защиты персональных данных» и др. Очевидно, что аспекты обращения с данными состоят из нескольких слов, и в некоторых случаях перекрываются. На основании этих фактов была выдвинута гипотеза о том, что латентно-семантический поиск

способен обнаружить даже незначительную разницу в тексте абзацев при пропуске частых слов. Перед применением латентно-семантического анализа требуется предварительная обработка входных данных. Обычно эта процедура включает очистку данных, удаление гиперссылок, пунктуации и т. д. Также текст политик конфиденциальности был разбит на набор абзацев. Каждый абзац был преобразован в массив слов, которые он содержит. Следующим шагом было удаление наиболее частых, но не столь значимых слов, так называемых стоп-слов. Также была применена операция стемминга, чтобы рассматривать только основную часть всех слов полученных от единого корня.

Пусть A – это матрица абзацев и слов, тогда используя формулу (2)

$$A = U \times S \times V^T \quad (2)$$

где A – матрица слов и параграфов;

U – ортонормированная матрица U ;

V – ортонормированная матрица V ;

S – диагональная матрица S , значения которой сингулярны для A .

После того, как матрица была разделена на три компонента, матрица U содержит n -мерные векторы, которые можно интерпретировать как координаты в n -мерном пространстве [6]. Документы могут быть распределены по кластерам по значениям этих координат. Проведенные эксперименты с латентно-семантическим анализом выполнялись с использованием набора данных с открытым исходным кодом, который включает 115 политик безопасности, которые были размечены вручную, и все абзацы присвоены одному или нескольким сценариям использования персональных данных [3]. Результаты экспериментов для модели «мешок слов» представлены в таблице 1, в ней показаны полученные кластеры и соответствующие значения

координат.

Таблица 1 – Кластеры политик безопасности для модели Bag-of-Words

№	Coordinate 1	Coordinate 2	Coordinate 3	Coordinate 4
0	0.634"inform"	0.28"may"	0.276"use"	0.232"servic"
1	0.202"cooki"	0.466"inform"	0.336"site"	0.257"use"
2	0.524"privaci"	0.433"polici"	0.388"cooki"	0.219"site"
3	-0.589"servic"	0.344"site"	0.244"parti"	-0.240"third"
4	-0.504"parti"	0.486 "third"	-0.449"servic"	0.235"advertis"
5	-0.594"site"	0.278"cooki"	0.272"websit"	0.264"privaci"
6	-0.326"may"	0.311"site"	0.307"servic"	-0.293"email"
7	-0.437"may"	-0.369"advertis"	0.345"person"	0.319"cooki"
8	0.501"may"	-0.315"email"	-0.281"use"	-0.264"address"
9	-0.488"user"	-0.384"use"	0.310"provid"	-0.301"websit"

Как видно, результаты противоречивы, поэтому трудно понять, какая из тем каким смыслом обладает. Затем рассчитывалась метрика принадлежности к теме с помощью библиотеки Gensim [7] и результаты снова не были обнадеживающими. Результаты расчета метрики принадлежности кластеру представлены в таблице 2.

Таблица 2 – Принадлежность кластерам

Topic	0	1	2	3	4
Affiliation	2.27	-0.8	0.15	-0.22	-1.2
Topic	5	6	7	8	9
Affiliation	-0.17	-0.15	-0.2	0.22	-0.07

Другие результаты с параграфами, относящимися к другому аспекту

обращения с данными, были почти такими же. Результаты представлены в таблице 3.

Таблица 3 – Принадлежность кластерам

Topic	0	1	2	3	4
Affiliation	2.59	-0.76	0.64	0.74	0.13
Topic	5	6	7	8	9
Affiliation	0.14	-0.12	0.23	0.12	0.41

Все протестированные абзацы были сопоставлены с кластером 0, что не может быть верным так как абзацы относились к заведомо разным аспектам обращения с персональными данными.

Результаты экспериментов для модели TF-IDF представлены далее, в таблице 4. Также показывались десять кластеров и значения атрибутов. И, как в первом случае с «мешком слов», по значениям координат невозможно судить о теме кластера.

Таблица 4 – Кластеры политик безопасности для модели Bag-of-Words

№	Coordinate 1	Coordinate 2	Coordinate 3	Coordinate 4
0	0.202“cooki”	0.2“may”	0.198“inform”	0.198“site”
1	0.573“cooki”	0.262“browser”	0.195“advertis”	0.182“web”
2	-0.406“media”	0.291“cooki”	0.282“health”	0.279“advertis”
3	-0.453“health”	0.258“email”	-0.204“kaleida”	0.191“address”
4	0.423“health”	0.215“media”	0.205“kaleida”	-0.199“secur”
5	-0.299“advertis”	0.262“health”	-0.252“media”	-0.213“privaci”
6	-0.325“media”	0.263“polici”	0.249“privaci”	0.197“chang”
7	0.280“cooki”	-0.216“device”	-0.183“health”	-0.166“social”
8	-0.223“advertis”	-0.206“teenag”	-0.206“inelig”	0.176“child”

Продолжение таблицы 4

№	Coordinate 1	Coordinate 2	Coordinate 3	Coordinate 4
9	-0.263'' child''	-0.26''wireless''	0.245''message''	0.239''parent''

Результаты кластеризации снова противоречивы, поэтому трудно сказать, какая конкретная тема описывает какой аспект политики конфиденциальности. В разных темах встречаются одни и те же слова с изменением веса. Для аспектов политики конфиденциальности, которые мы искали нет тем, которые могли бы их точно описать, поскольку многие из них могут. Затем с помощью библиотеки Gensim был рассчитан показатель принадлежности к теме, и результаты снова не были обнадеживающими. Результаты расчета аффилированности по абзацу одной из политик конфиденциальности представленные в таблице 5.

Таблица 5 – Принадлежность кластерам

Topic	0	1	2	3	4
Affiliation	2.18	-0.97	-0.69	-0.27	0.65
Topic	5	6	7	8	9
Affiliation	0.98	-1.17	0.8	0.27	0.01

Результат для другого абзаца, относящегося к другой политике конфиденциальности, был почти такой же. Результаты представлены в таблице 6.

Таблица 6 – Принадлежность кластерам

Topic	0	1	2	3	4
Affiliation	1.82	0.25	0.49	0.29	-0.04
Topic	5	6	7	8	9

Продолжение таблицы 6

Affiliation	0.74	0.52	-0.04	-0.58	-1.33
-------------	------	------	-------	-------	-------

Как можно заметить, результаты для модели TF-IDF аналогичны результатам модели «мешка слов», за исключением нескольких незначительных изменений. Все абзацы снова были сопоставлены с кластером 0, что неверно, потому что они на самом деле описывают разные сценарии использования персональных данных. Эти эксперименты позволили сделать вывод, что использование латентно-семантического анализа не дает ценной информации о содержании онлайн-согласия пользователя. Проблема может быть связана с тем, что сценарии использования персональных данных очень похожи между собой, и для того, чтобы различать разные сценарии необходимо учитывать локальный контекст.

В результате апробации алгоритма латентно-семантического анализа было выяснено что для кластеризации экстремально схожих между собой текстов он подходит не лучшим образом. В связи с этими обстоятельствами было решено обратить внимание на несколько иной подход анализа текста, основанный на контекстно-свободных грамматиках, тегировании по частям речи и синонимическом поиске.

1.3 Подход основанный на латентном размещении Дирихле

1.4 Подход основанный на применении контекстно-свободных грамматик и синонимическом поиске

Другой предложенный подход – подход, основанный на анализе с помощью контекстно-свободных грамматик и синонимического поиска. Синонимический поиск в данном случае – это подмена ключевых слов и их синонимов метками, например «__FP_A__» означает, что это слово и его синонимы считаются акторами первой стороны. Этот метод можно применить ко многим другим словам. Например, сообщения электронной почты, авата-

ры, местоположение также могут быть объектами и синонимами абстрактной метки «__CN__», которая означает существительное сбора или объект сбора. Так все ключевые слова могут быть преобразованы в их смыслы в контексте предметной области. Маркировка выполняется легко, все слова совпадающие с пулами заменяются метками этих пулов.

Предварительная обработка данных в данном случае состоит из токенизации и лемматизации для более гибкой замены слов на метки их пулов.

При анализе пользовательского согласия сайта недостаточно найти ключевые слова, относящиеся к разным типам персональных данных, например цель и правовую основу распознать гораздо сложнее. Следующий шаг - установить слова отношения в предложениях, чтобы можно было определенно сказать, что ярлыки пулы синонимов связаны друг с другом и формируют логическая цепочку. Один из возможных способов определения отношений слов в тексте на естественном языке – это синтаксический анализ предложения, основанный на частеречной разметке [8]. Имея размеченное по частям речи предложение, парсер грамматики NLTK [9] строит деревья предложений по правилам грамматики. Одно из таких деревьев в обозначениях NLTK можно увидеть на рисунке 3 [9], где «S» – основа предложения, «NP» – именная фраза, «VP» – глагольная фраза, «Adj» – прилагательное, «НОМ» – именное словосочетание, «ПП» – предлог фраза, «Det» – артикль, «V» – глагол, «N» – существительное, «P» – предлог.

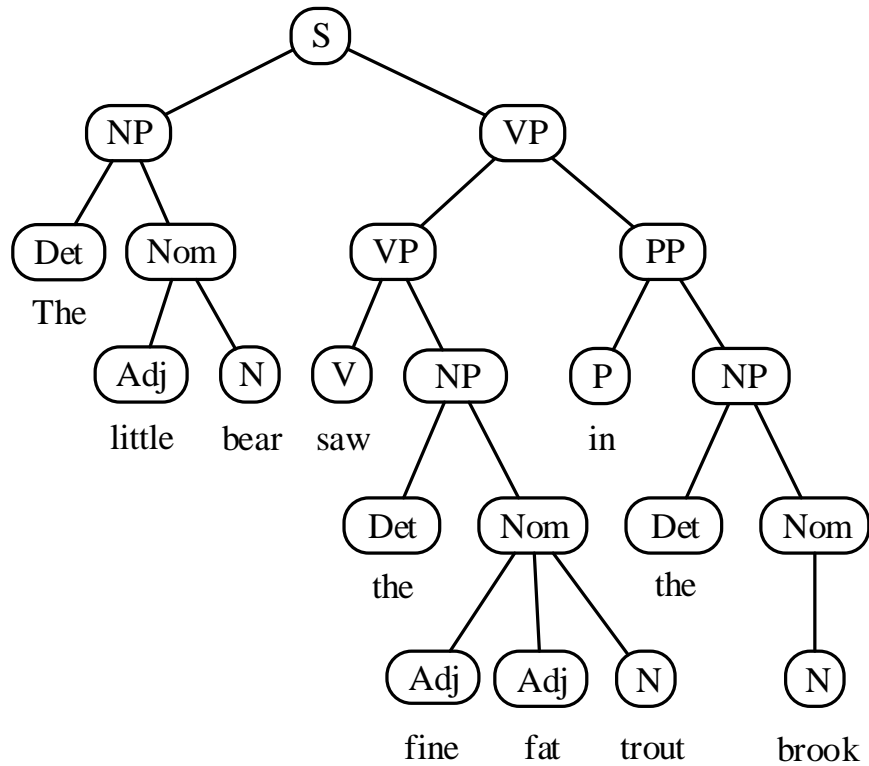


Рисунок 3 – Матрица TF-IDF

В предлагаемом подходе немного другая грамматическая запись. Созданная грамматика представлена в (3).

$$\left\{ \begin{array}{l} D \rightarrow S \mid S D \mid S U D \\ S \rightarrow NPG \ VBG \\ VPG \rightarrow VP \mid VP \ VPG \mid VP \ U \ VPG \\ NPG \rightarrow NP \mid NP \ NPG \mid NP \ U \ NPG \\ AJPG \rightarrow AJ \mid AJ \ APG \mid AJ \ U \ APG \\ AVPG \rightarrow AV \mid AV \ APG \mid AV \ U \ APG \\ VP \rightarrow VAPG \mid V \ PPG \mid V \ PP \ APG \\ NP \rightarrow NOM \mid DET \ NOM \\ NOM \rightarrow N \mid AJPG \ N \\ PP \rightarrow NPG \mid P \ NPG \end{array} \right. , \quad (3)$$

где «D» – документ, «S» – предложение, «NPG» – группа именных фраз, «VPG» – группа глагольных фраз, «AJPG» – группа прилагательных, «AVPG»

– группа наречий, «VP» – глагольная фраза, «NP» – именная фраза, «NOM» – именная фраза, «PP» – предлог, «DET» – артикль. Грамматика из формулы (3) позволяет рекурсивно выделять основу предложения и последовательности глагола, существительного, прилагательного, наречия и т.д. Это все еще не идеальное решение, но попытка найти более сложные предложения в политиках безопасности. Этот подход требует использования пулов синонимов, которые соответствуют различным ключевым словам. Поэтому в грамматику включены метки пулов синонимов, привязанных к части речи. Метки пулов вручную назначены частям речи для преобразования привязок частей речи NLTK, это показано в формуле (4).

$$\left\{ \begin{array}{l} U \rightarrow NLTK_CC \\ DET \rightarrow NLTK_DT \\ AJ \rightarrow NLTK_JJ \\ AV \rightarrow NLTK_RB \\ N \rightarrow _CN_ | _FP_A_ | _TP_A_ | NLTK_N \\ V \rightarrow _CV_ | NLTK_V \end{array} \right. \quad (4)$$

Здесь в формуле (4) «NLTK_CC» – соединение NLTK, «NLTK_N» – все формы существительных NLTK, «NLTK_V» – все формы глаголов NLTK, «NLTK_DET» – определители NLTK, «NLTK_JJ» – все NLTK формы прилагательных, «NLTK_RB» – все формы наречий NLTK. Теги, начинающиеся с подчеркивания, являются метками пулов синонимов. Синтаксический анализ выполняет библиотека NLTK. На основе предложенной грамматики, описанной (3) и (4) и разметки лейблами пулов было построено дерево предложения, результат на рисунке 4.

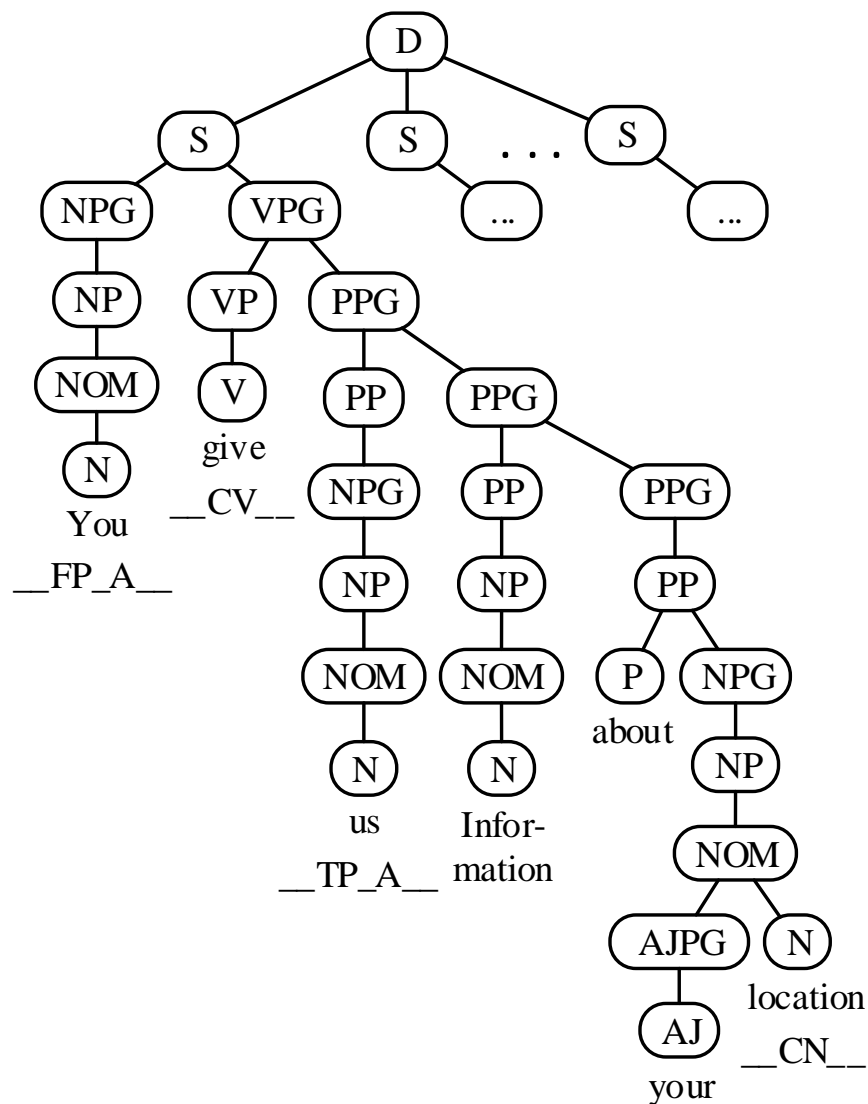


Рисунок 4 – Матрица TF-IDF

Здесь на рисунке 4 «__FP_A__» - метка актора-обладателя персональных данных, «__TP_A__» - третья сторона, «__CV__» - глагол сбора, «__CN__» - существительное сбора. Когда было построено дерево предложений последовательность меток ключевых слов может быть распознана. В этом случае представленная на рисунке 4, последовательность «__FP_A__», «__CV__», «__CN__» хорошо видна. Такие атомарные последовательности, раскрывают значения частей предложения и могут быть объединены в список, после этого весь смысл документов будет описан этим список. Сочетание маркировки ключевых слов и синтаксического анализа дает значения ключевых слов с отношениями между этими словами, определенными в виде древо-

видных структур. Дерево структура данных более гибкая, чем строка предложения, деревья и особенно поддерева показывают важные отношения между словами. Запросы к таким структурам могут дать необходимую информацию для построения логических последовательностей действующих лиц, их действий, субъектов этих действий и, наконец, обстоятельств. Предлагаемый подход определенно имеет такие недостатки, как низкая производительность, вручную определенные пулы синонимов и т.д.

В результате апробации алгоритма латентно-семантического анализа было выяснено что для кластеризации экстремально схожих между собой текстов он подходит не лучшим образом. В связи с этими обстоятельствами было решено обратить внимание на несколько иной подход анализа текста, основанный на контекстно-свободных грамматиках, тегировании по частям речи и синонимическом поиске.

1.5 Выводы по строгим методам текстового анализа

Эксперименты показали, что оба рассмотренных метода имеют как преимущества, так и определенные недостатки. Хотя предложенные подходы, оказались противоречивыми, окончательные результаты заслуживают внимания. Подход с латентно-семантическим поиском оказался не слишком эффективным. Однако, подход основанный на грамматическом анализе предложений и синонимическом поиске дал определенные результаты. Хоть он и не является производительным, с его помощью возможно производить выделение логических цепочек из предложений для получения более формального описания политик безопасности нежели их текстовые варианты.

1.6 Подход основанный на глубоком обучении

Исходя из проведенных исследований стало понятно, что более предпочтительным вариантом решения задачи будет подход с применением моделей с глубоким обучением. Реализация подобного проекта – комплексная задача, ее можно разделить на несколько этапов. Сначала необходимо собрать датасет, потом его разметить для обучения модели, далее обучить модель и

получить результаты. Однако сбор датасета тоже является непростой задачей. Для того чтобы осуществить сбор датасета необходим инструмент для поиска и скачивания веб-страниц из сети интернет. Затем необходимо произвести очистку данных, удалить все теги со страниц, чтобы можно было передать текст аннотаторам. Все этапы сбора датасета полагаются на базу данных. Она лишена сложного объектно-реляционного моделирования, так как в ней по сути необходимо только хранить промежуточные результаты обработки текстовых файлов.

2 Проектирование инструментария

2.1 Техническое задание «Инструментарий для сбора датасета»

2.1.1 Основные положения технического задания

2.1.2 Скрейпер вэб-страниц

Скачивание веб-страниц будет производиться инструментом написанном на языке Python, с помощью библиотек можно скачивать страницы анализировать данные с них, переходить по гиперссылкам и много другое. Такой инструмент позволит просматривать и сохранять содержимое страниц в автоматическом режиме без вмешательства пользователя. Таким образом в автоматическом режиме можно сохранить и проанализировать огромное количество текстовой информации.

2.1.3 Очистка скачанных страниц политик

Для очистки страниц от кода разметки планируется использовать библиотеку «html sanitizer». Очистка кода необходима для того, чтобы аннотаторы могли максимально сфокусироваться на анализе текста, таким образом получая чистый текст они не будут отвлекаться на не имеющие значения в контексте задачи фрагменты.

2.1.4 Инструмент разметки датасета

Инструмент разметки датасета планировалось реализовать с помощью веб-технологий. Серверная часть будет полагаться на приложение, написанное на PHP, которое будет регулировать порядок выдачи текста на аннотирование. Процесс разметки высокодинамичен, поэтому невозможно избежать написания качественной клиентской части приложения на языке javascript. Это позволит сделать работу аннотаторов максимально производительной, в «одну сессию», так как страница не будет перезагружаться, однако все изменения, которые будут вноситься, сохраняться.

2.1.5 Фреймворк глубокого обучения

Для создания и тренировки модели анализа текста планируется использовать фреймворк машинного обучения «Keras». Он позволяет быстро создавать классификаторы с самыми разными конфигурациями и любых типов.

После того как классификатор будет сконфигурирован останется лишь обучить его на датасете, полученном ранее.

Обученный классификатор будет в состоянии определять различные характеристики политики безопасности и аспекты обращения с данными, что позволит в автоматическом режиме формировать краткие отчеты о безопасности предоставляемого соглашения.

2.2 Первичная декомпозиция и планирование

Начальным этапом решения задачи является первичная декомпозиция, в ее результате выделяются подзадачи различной важности, которые должны быть решены для доведения цикла разработки до конца. В данном случае можно выделить следующие подзадачи:

- 1) определение источника информации о различной IoT-продукции,
- 2) отправка поискового запроса,
- 3) получение результатов запроса (список IoT-продуктов),
- 4) определение производителей IoT-продукции,
- 5) поиск официальных сайтов производителей в сети интернет,
- 6) поиск раздела «политика безопасности» на сайтах производителей,
- 7) скачивание политик безопасности,
- 8) очистка скачанных веб-документов от лишних элементов разметки.

Получение списка производителей возможно на электронных торговых площадках, типовая разметка веб-страниц располагает для получения такой информации, так как существует лишь несколько вариантов наполнения страницы продукции.

Получение официальных веб-сайтов производителей задача на первый

взгляд сложная, однако результаты ручной проверки показали, что лучшим вариантом является поисковый запрос с названием производителя «как есть». В таком случае веб-сайт производителя оказывается на первой строчке результата поискового запроса, а если не оказывается, значит у этой компании его с очень большой вероятностью нет.

2.3 Приложение вэб-скрейпер

2.3.1 Структура приложения вэб-скрейпера

Исходя из результатов декомпозиции, эффективным подходом выглядит представление приложения в виде последовательно выполняющихся подпрограмм так, что входом модуля является результат работы предыдущего модуля, то есть в виде конвейера. Схема организации приложения представлена на рисунке 5.

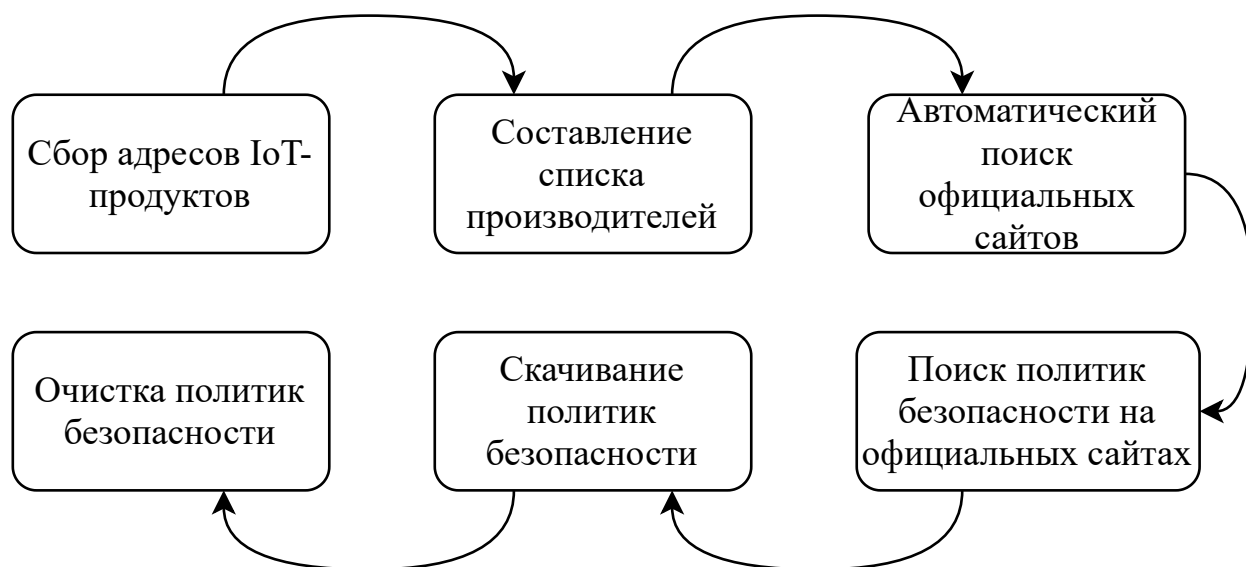


Рисунок 5 – Схема организации приложения

Далее была разработана композиционная модель приложения, на ней присутствуют все необходимые для решения задач модули. Схема представлена на рисунке 6.

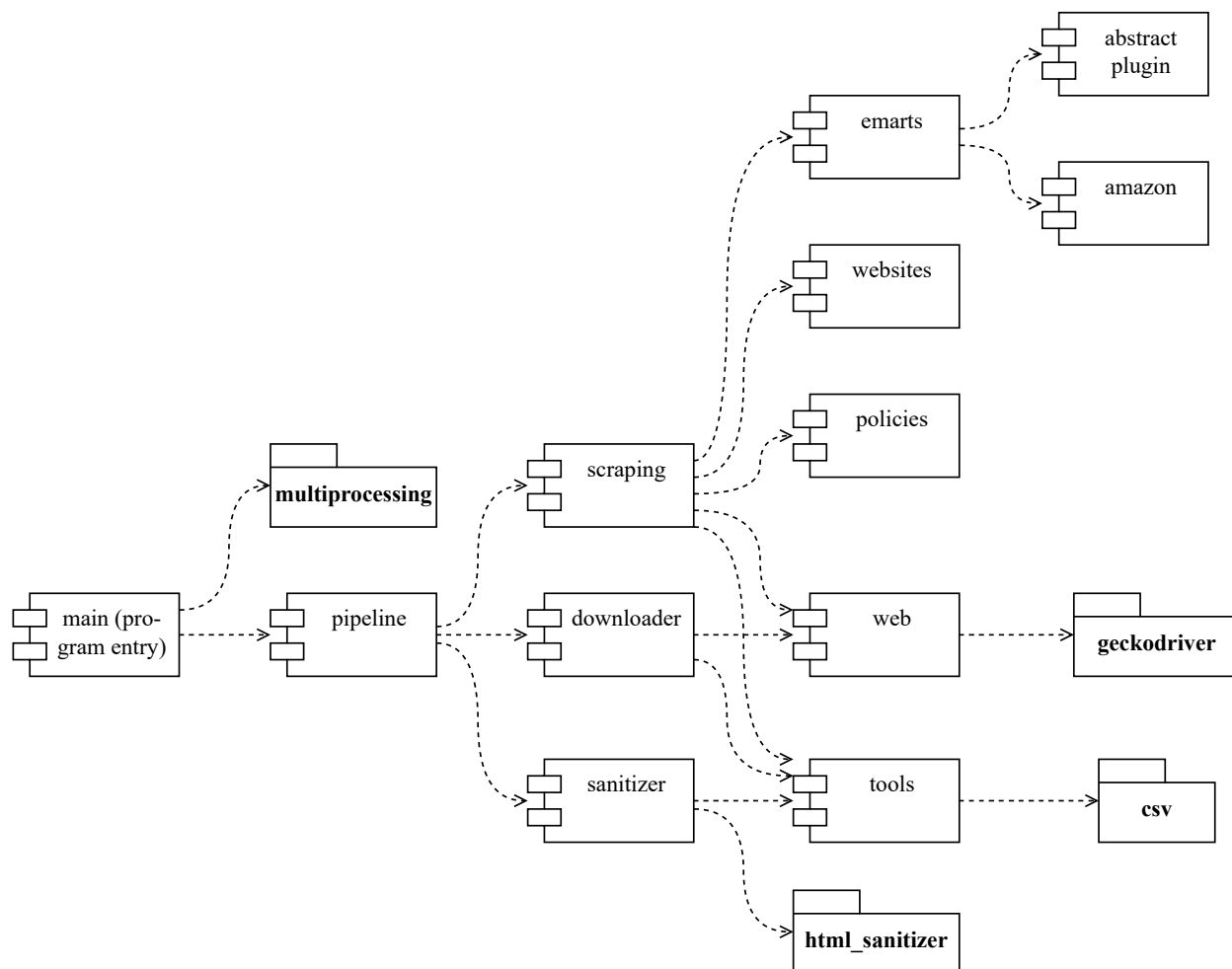


Рисунок 6 – Композиционная модель приложения

Здесь на рисунке 6 модуль «main» отвечает за запуск программы, раз-
вертывание основных ее частей. Там же происходит инициализация пула по-
токов для мультипроцессинга затратных задач таких как, например, взаимо-
действие с «безголовым» браузером.

Модуль «pipeline» отвечает за последовательное исполнение подпро-
грамм элементов конвейера. Он осуществляет прием выходных и передачу
входных данных.

Модуль «scraping» отвечает за получение данных с веб-страниц. Ин-
формация полученная с помощью этого модуля записывается в csv-файл для
большей прозрачности и возможности сохранения результатов между запус-
ками приложения, например, для пропуска данного этапа и использования
его сохраненных результатов работы. Он включает в себя 3 основных дей-

ствия:

- 1) получение производителей IoT-продуктов (соответствует модулю «emarts»),
- 2) получение официальных сайтов производителей (соответствует модулю «websites»),
- 3) получение веб-ссылок на политики безопасности (соответствует модулю «policies»).

Модуль «downloader» отвечает за скачивание страниц и их сохранение в отведенную для этого директорию. Информация полученная с помощью этого модуля записывается в html-файлы для большей прозрачности и возможности сохранения результатов между запусками приложения, например, для пропуска данного этапа и использования его сохраненных результатов работы.

Модуль «sanitizer» отвечает за очистку скачанных веб-страниц от ненужных тегов и ссылок. Информация полученная с помощью этого модуля записывается в html-файлы для большей прозрачности и возможности сохранения результатов между запусками приложения, например, для пропуска данного этапа и использования его сохраненных результатов работы.

Модуль «emarts» включает в себя набор плагинов, каждый из которых адаптирован для получения требуемой информации с определенного шаблона веб-страничной разметки. Некоторое поведение инкапсулировано в абстрактном плагине для увеличения «reusability» кода. Получая адрес на вход, данный плагин производит скачивание страницы и с помощью набора шаблонов пытается извлечь информацию. Данный модуль записывает полученную с помощью плагинов информацию в csv-файл для большей прозрачности и возможности сохранения результатов между запусками приложения, например, для пропуска данного этапа и использования его сохраненных результатов работы.

Модуль «web» отвечает за взаимодействие с веб-сайтами будь то торго-

вые площадки или сайты производителей IoT-продуктов. В нем используется geckodriver для управления «безголовым» браузером.

Модуль «tools» содержит вспомогательные функции, в частности для ввода и вывода данных в формате csv.

2.3.2 Потенциальные проблемы

Потенциально возможные проблемы при реализации приложений подобного типа:

- 1) блокировка из-за подозрительных заголовков браузера,
- 2) блокировка из-за слишком частого обращения с запросами,
- 3) как следствие 2-х предыдущих пунктов требование подтвердить, что это не попытка автоматического доступа (ввод капчи).

Упомянутые проблемы решаются использованием разных заголовков браузера попеременно. Также отправка запросов ограничена по частоте от 15 до 30 секунд, ограничение выбирается случайным образом. Такие решения позволяют крайне редко попадать под подозрения, соответственно процент успеха при попытке получить данные с веб-страницы значительно повышается.

2.4 Инструмент разметки датасета

Инструмент разметки датасета планировалось реализовать с помощью веб-технологий. Серверная часть будет полагаться на приложение, написанное на PHP, которое будет регулировать порядок выдачи текста на аннотирование. Процесс разметки высокодинамичен, поэтому невозможно избежать написания качественной клиентской части приложения на языке javascript. Это позволит сделать работу аннотаторов максимально производительной, в «одну сессию», так как страница не будет перезагружаться, однако все изменения, которые будут вноситься, сохранятся.

2.4.1 Объектное моделирование приложения

Объектная модель инструмента представлена на рисунке 7.

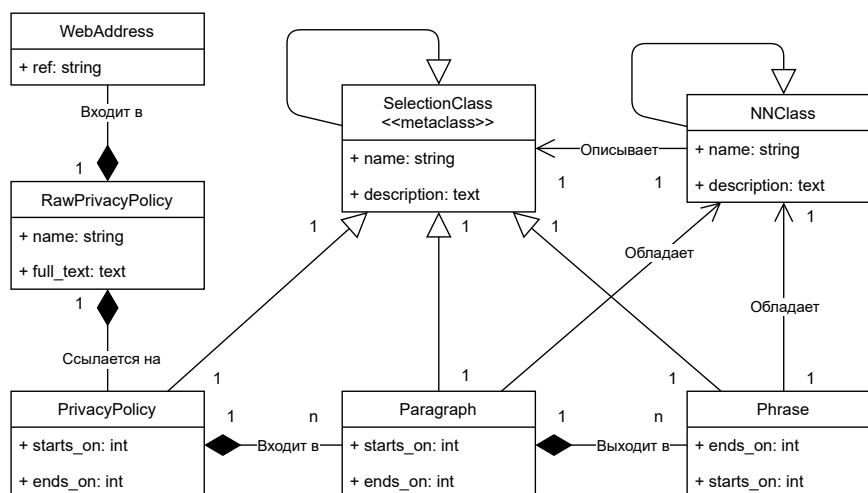


Рисунок 7 – Объектная модель

2.4.2 Реляционная модель приложения

Реляционная модель инструмента представлена на рисунке 8.

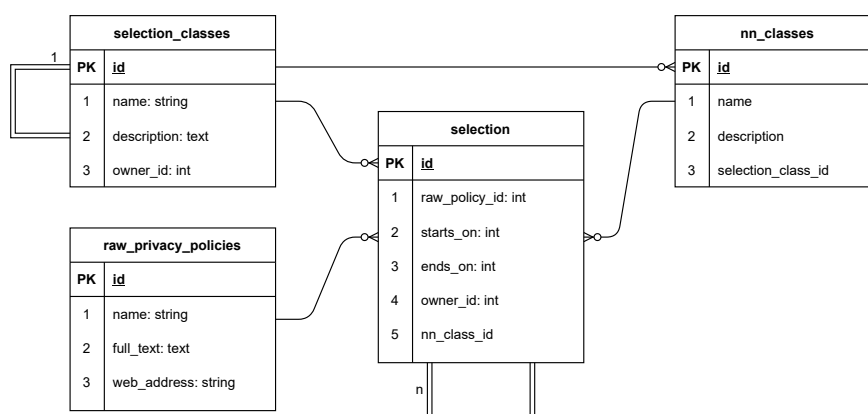


Рисунок 8 – Реляционная модель

2.4.3 Проектирование пользовательского интерфейса

Презентационный прототип интерфейса инструмента представлен на рисунке 9.

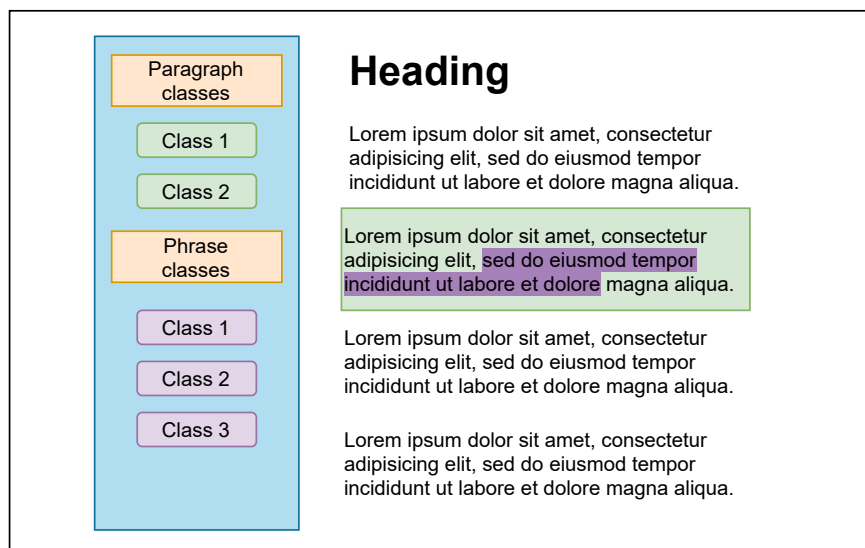


Рисунок 9 – Презентационный прототип интерфейса

2.5 Выбор средств разработки

Для реализации приложения были выбраны следующие средства:

- 1) python 3.8.5,
- 2) «безголовый» браузер Firefox,
- 3) драйвер для управления «безголовым» браузером geckodriver,
- 4) библиотека html sanitizer для очистки скачанных веб-документов.

Выбор «безголового» браузера обусловлен потребностью в отрисовке страниц, так как на некоторых веб-страницах разметка генерируется с помощью javascript. Это делает невозможным использование простого скачивания, необходима страница именно с исполненными скриптами, в противном случае будет невозможно получить требуемую информацию. В то же время браузер лишен графического интерфейса, чем снижается потребление вычислительных ресурсов.

3 Реализация инструментария

3.1 Полученные в результате реализации исходные коды

В соответствии с результатами декомпозиции, выбора средств и проектирования приложение было реализовано. Характеристики полученных классов и функций приведены далее, в таблицах ?-?. Исходные коды представлены в приложении А.

3.2 Полученный в результате реализации пользовательский интерфейс инструмента разметки

3.3 Результаты решения поставленной задачи с помощью разработанного инструментария

В результате работы программы были найдены 65 политик безопасности, разумеется среди них имеется определенный процент промахов, если производитель имеет сходство с каким либо другим более крупным. Поиск осуществлялся по торговой площадке amazon, брались результаты поискового запроса по первым 10-ти страницам, по категориям «smart scales», «smart watches», «smart locks» и «smart bulbs». Всего производителей было найдено приблизительно 160. Стоит отметить, что результат является приемлемым, так как многие производители на данной торговой площадке не имеют выделенного веб-сайта, а пользуются услугами amazon, то есть на таких страницах действует политика безопасности amazon, а не производителя. Также стоит отметить что у некоторых продуктов явно не указан производитель, что сократило количественно результат поиска.

ЗАКЛЮЧЕНИЕ

Исходя из анализа методов формализации политик безопасности, было принято решение продолжать движение в сторону создания инструментов разметки датасетов, и моделей глубокого обучения. Таким образом было проведено первичное планирование процесса выполнения выпускной квалификационной работы магистра. В результате выполнения работы было спроектировано и реализовано требуемое программное средство для сбора датасета, ориентированного на политики безопасности.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- 1 General Data Protection Regulation, GDPR homepage. URL: <https://gdpr.eu> (дата обращения 14.02.2021).
- 2 Electronic Passes in Moscow during lockdown. URL: https://www.cnews.ru/news/top/2020-05-25_moskovskij_sajt_s_propuskami (дата обращения 14.02.2021).
- 3 Harshvardhan J. Pandit, Declan O’Sullivan, and Dave Lewis. Personalised Privacy Policies. 2018.
- 4 Hamza Harkous, Kassem Fawaz, Remi Lebre, Florian Schaub³, Kang G. Shin, and Karl Aberer. Polisis: Automated Analysis and Presentation of Privacy Policies Using Deep Learning. 2018. arXiv:1802.02561v2.
- 5 Evgenia Novikova, Elena Doynikova, and Igor Kotenko. P2Onto: Making Privacy Policies Transparent. Springer, 2020.
- 6 Landauer, T. K., Foltz, P. W., and Laham, D. An Introduction to Latent Semantic Analysis. Discourse Processes, 25, 1998, pp. 259-284. DOI: <https://doi.org/10.1080/01638539809545028>.
- 7 Gensim topic modeling library, Gensim homepage. URL: <https://radimrehurek.com/gensim> (дата обращения 14.02.2021).
- 8 Sachini Weerawardhana, Subhojeet Mukherjee, Indrajit Ray, and Adele Howe. Automated Extraction of Vulnerability Information for Home Computer Security, pages 356–366. Springer, 2015. DOI: https://doi.org/10.1007/978-3-319-17040-4_24.
- 9 Natural Language ToolKit, Analyzing Sentence Structure, NLTK homepage. URL: <https://www.nltk.org/book/ch08.html> (дата обращения 14.02.2021).

ПРИЛОЖЕНИЕ А

Архив с исходными кодами вэб-скрейпера и инструмента для разметки датасета.