

## **Применение алгоритмов текстовой аналитики и практик OSINT в информационно-аналитических системах**

### **Александр Юрьевич Выжигин**

Кандидат технических наук, доцент, заведующий кафедрой  
Российская академия народного хозяйства и государственной службы при Президенте Российской Федерации  
Москва, Россия  
Кандидат технических наук, доцент  
Российский технологический университет  
Москва, Россия  
vijigin\_new2000@mail.ru

### **Илья Сергеевич Москалев**

Студент  
Российский технологический университет  
Москва, Россия  
moskalevilya1@gmail.com

### **Олег Владимирович Трубиенко**

Кандидат технических наук, доцент, заведующий кафедрой  
Российский технологический университет  
Москва, Россия  
trubienko@mail.ru

Поступила в редакцию 27.03.2023

Принята 01.04.2023

### **Аннотация**

Информация – важный ресурс нашей жизни, который принято делить на два вида в зависимости от значимости: 1) точные данные – данные, которые неискаженно отражают события, факты и помогают провести грамотный анализ произошедшего; 2) избыточные (мусорные) данные – данные, которые могут ввести аналитика в заблуждение, таким образом, создав риск возникновения ошибки в ходе анализа. Проблема современного информационного общества заключается в том, что люди не всегда могут извлечь точные данные для последующего анализа информации. С данной проблемой способны справиться информационно-аналитические системы (особый класс информационных систем, обладающих свойством аналитической обработки данных при автоматизации процесса получения ответа на тот или иной запрос). Работу информационно-аналитических систем (ИАС) можно сравнить с информационными системами (ИС). Рассмотрим одну из известных всем поисковых систем – «Яндекс». Поисковые системы также являются информационно-аналитическими. К примеру, мы хотим получить информацию о чемпионате мира по футболу 2014 года. В Яндексе мы формируем запрос – «чемпионат мира по футболу 2014 года», что на языке SQL выглядело бы так (SELECT football FROM world\_championship WHERE year = 2014). В ответе от Яндекса содержится 5000 результатов, и чтобы выяснить необходимую информацию, необходимо или проанализировать эти ответы, или переформулировать поисковые запросы. Поисковые системы похожи на ИАС, но для оптимального решения стоит использовать те, которые содержат в себе удобный интерфейс для пользователя, где по кнопкам из меню можно попасть к нужному функционалу и получить желаемый результат без дополнительного анализа. Научная новизна исследования заключается в следующем: 1) расширены возможности текстовой аналитики за счет использования практик OSINT при выполнении запросов в ИАС; 2) разработан метод извлечения информации из российского сегмента сети Интернет; 3) разработан метод проверки запросов к ИАС на их логичность с точки зрения лексикологии; 4) разработан защищенный доступ к функционалу ИАС. Цель работы – повышение эффективности методов текстовой аналитики за счет применения практик OSINT. В статье использовались такие методы как: OSINT (поиск информации из открытых источников, метод деловой разведки, суть которого заключается в получении конкретного ответа на конкретный вопрос); текстовая аналитика (процесс обработки неструктурированного текста для выявления идей, закономерностей и т. д.). В результате был разработан функционал информационно-аналитической системы, использующий алгоритмы текстовой аналитики и практики OSINT. Полученные результаты могут использоваться работниками различных служб безопасности нашей страны, а также обычными пользователями данной ИАС для оперативного получения нужной информации.

### **Ключевые слова**

информационно-аналитическая система, текстовая аналитика, NLP, Open source intelligence, OSINT, язык запросов Google dorks, язык программирования Python, морфологический анализатор, data science, поисковые системы.

### Введение

В настоящее время из-за больших объемов данных в интернете возникают проблемы с поиском необходимой информации и её анализа. Сегодня существует различные поисковые системы, такие как Google, Bing, и в том числе сделанные в России, – Yandex, Mail, Rambler. Пользователи в бытовом использовании выбирают более удобную для них поисковую систему, но при поиске специфической узконаправленной информации прибегают к результатам сразу нескольких поисковых систем, потому что каждая выдает различные результаты. Основная проблема организации поиска информации состоит в невозможности создания идеальной поисковой системы. И даже при наличии таких гигантов как Google ведутся научные работы по поиску и анализу информации в интернете (Брумштейн, 2017; Расел, 2021; Jason H. S., 2018).

Для решения поставленной проблемы было проведено исследование, целью которого являлось повышение эффективности методов текстовой аналитики за счет применения практик OSINT (Batrinca, Philip, 2015).

Для достижения поставленной цели были решены следующие задачи:

- проведен обзор и анализ существующих информационно-аналитических систем в предметной области;
- описана суть предлагаемого решения;
- описаны практики OSINT;
- описаны особенности предобработки запросов;
- описана интеграция техник OSINT;
- описан метод решения задачи;
- проведено исследование работоспособности предложенного метода на основе разработанной ИАС.

### Материалы и методы исследования

Для разработки улучшенной ИАС требуется провести обзор и анализ уже существующих систем. Рассмотрим некоторые из них (табл. 1).

Таблица 1. Сравнение ИАС

	Медialogия	Интегрум	Public.ru	Park.ru
Год появления на рынке	2003	1996	1990	1995
Количество источников	2910	5000	2000	1300
Источники информации	русскоязычные + иностранные	русскоязычные + иностранные	русскоязычные	русскоязычные
Программная платформа	собственная разработка	собственная разработка	на базе системы RetrievalWare	на базе системы Yandex.Server
Обработка данных	автоматизированная + ручная	автоматизированная	автоматизированная	автоматизированная
Оценка интерфейса	5/5	4/5	3/5	3/5
Поиск информации	объектно-ориентированный + контекстный	контекстный	контекстный	контекстный
Оценка поиска информации	5/5	4/5	4/5	3/5
Оценка формата представления данных	5/5	4/5	3/5	3/5
Оценка способа доставки данных	5/5	4/5	4/5	3/5
Оценка ведения архива	5/5	4/5	отсутствует	отсутствует
Доступность использования	корпоративная	корпоративная + персональная	корпоративная + персональная	корпоративная + персональная
Оценка дополнительных услуг	5/5	5/5	5/5	5/5
Тестовый доступ	есть	есть	есть	есть
Оценка информативности сайта ИАС	5/5	5/5	3/5	4/5

Система «Медialogия» изначально разрабатывалась для внутренних нужд холдинга IBS, специализирующегося на управленческом консалтинге; с 2003 г. система «Медialogия» была запущена в коммерческую эксплуатацию. Уникальной особенностью системы «Медialogия», с точки зрения пользователя,

является удобный и разнообразный инструментарий для анализа информации и наглядность представления результатов.

Система «Интегрум» на сегодняшний день является самым полным электронным архивом русскоязычных СМИ, созданным на основании материалов из открытых источников информации. Уникальность информационно-аналитической системы «Интегрум» для пользователя заключается в возможности получения информации из адресно-справочных и правовых баз данных – Роспатента, Госкомстата, а также из специализированной литературы.

Public.ru. позиционируется как публичная библиотека, которая специализируется на оказании услуг по предоставлению доступа к отечественным периодическим изданиям. Основное отличие системы от рассмотренных выше заключается в простоте поиска информации и существенно более низкой стоимости услуг, что позволяет компании работать не только с организациями, но и частными лицами.

Park.ru. позиционируется как полнотекстовая библиотека российских СМИ. Особенность системы – большое количество готовых тематических мониторингов новостей.

По результатам сравнения лучшей является система Медиалогия, так как при среднем количестве источников, по сравнению с другими, выдает наилучший результат.

### Результаты и обсуждения

Для устранения недостатков в задачах поиска и анализа информации была создана собственная информационно-аналитическая система (рис. 1).

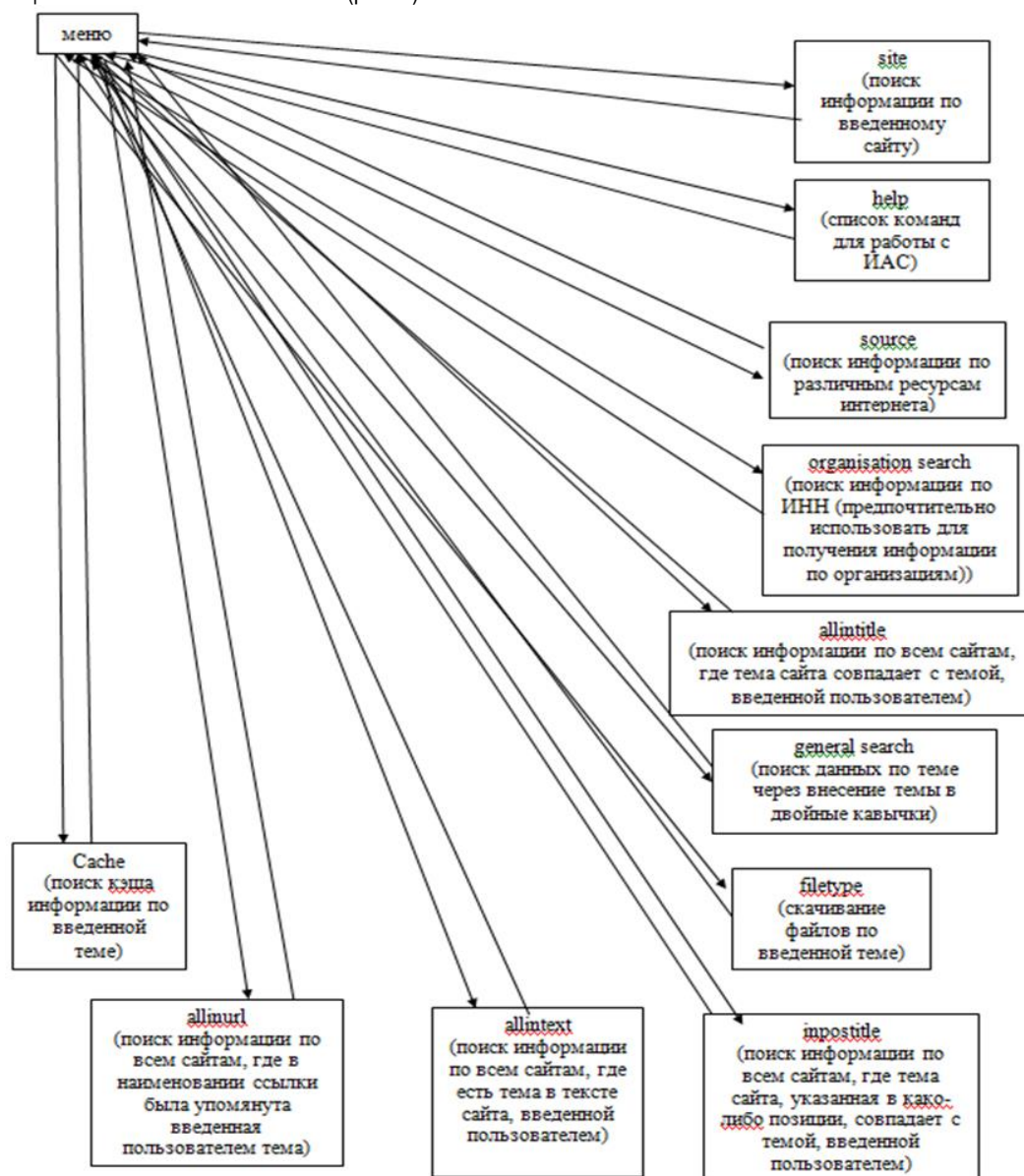


Рис. 1. Схема взаимодействия компонентов в ИАС

Алгоритм работы ИАС:

При запуске системы появится меню, где будет отражен весь функционал ИАС.

Запуск одного из видов поиска информации – потребует ввести одну из комбинаций кодовых значений, «вшитых» в систему под каждый из функционалов технологии OSINT. Команды были закодированы в соответствии с присвоенными им кодами (код строился с помощью 1 линии клавиатуры). Смысл состоит в том, что пользователь, не зная этих кодов, не сможет воспользоваться функционалом данной ИАС. Если пользователь введет неверную комбинацию, выведется сообщение об ошибке и будет предложено попробовать еще раз (Брумштейн, Васьяковский, Куаншалиев, 2017). Такая же возможность будет у него и после выполнения запроса. Если же была введена верная комбинация, то останется ввести данные, запрашиваемые системой, и ждать результат работы.

Результат при удачной обработке информации будет выглядеть следующим образом: сначала выведется таблица со столбцами: res (ссылки на ресурсы), net (домены российского сегмента интернета), work (показатели работоспособности сайтов). Если речь идет не о получении файла, то алгоритм получения данных таков: одна ссылка – извлечение данных из одной ссылки; больше одной ссылки – выбор ссылки будет осуществляться информационно-аналитической системой.

При желании пользователь может сохранить информацию в файл, указав его имя. В случае со скачиванием файлов (xlsx, pptx, csv и т.д.) алгоритм таков же. Его единственное отличие состоит в скачивании уже имеющихся файлов с расширением, указанным выше (Рассел, 2021).

При разработке данной системы использовалась технология получения и анализа информации из открытых источников под названием Open source intelligence (OSINT).

В процессе работы информационно-аналитической системы существуют особенности предобработки запросов.

Проверка запроса к ИАС на логичность с точки зрения лексикологии.

Эта проверка осуществляется с помощью морфологического анализатора, входящего в состав алгоритмов текстовой аналитики. Смысл его работы таков: тема запроса разбивается на отдельные слова, далее каждое из слов обрабатывается. Из обработки получается вероятность существования каждого слова. Потом все вероятности суммируются (1), в конце результат проверяется с условием (2). Если значение проходит, то осуществляется переход к следующему этапу; в противном случае будет предложено отказаться от поиска информации по введенному запросу (Бельдеубаева, 2020).

$$p = \sum_{i=0}^n p_i; \quad (1)$$

формула общей вероятности введенного запроса

$p_i$  – вероятность существования  $i$ -ого слова

$p$  – суммарная вероятность существования слов, указанных в запросе к ИАС

$$p \geq k * \text{threshold}; \quad (2)$$

условие проверки корректности запроса с точки зрения лексикологии

$k$  – количество слов в запросе

threshold – порог вхождения;

Порог вхождения - минимальная вероятность, которую должно иметь каждое слово

threshold=0.66;

Вторая особенность предобработки заключается в методе извлечения ссылок, находящихся в российском сегменте (Таршис, 2020). Для того, чтобы объяснить, как это происходит, прибегнем к анализу одного из исходных кодов, а точнее, к одной из его частей (рис. 2)

```
for i in res:
    t = i
    y = list(t)
    z = y[8:]
    h = z.index("/")
    k = z[:h]
    elem = "."
    idx = np.argwhere(np.array(k) == elem).flatten()
    idx = list(idx)
    d = idx[len(idx) - 1]
    s = k[d] + "+"
    res = s + z[h:]
net.append(res)
```

Рис. 2. Получение доменов из ссылок интернета

Просмотрев данный код, мы видим, что в процессе цикла перебора элементов берется каждая из ссылок, далее берутся срезы массива с учетом полученных индексов и в конце результат обработки будет конвертирован в строку и добавлен в столбец доменов (Голушко, Дрянных, 2019).

Интеграция техник OSINT. В разработанной информационно-аналитической системе присутствует методология поиска информации по открытым источникам OSINT (Москалев, 2022). Для интеграции техник данной методологии был выбран язык Googledorks. Googledorks – набор запросов для выявления грубейших дыр в безопасности. Всего, что должным образом не скрыто от поисковых роботов. Далее был создан в каждом из функционалов ИАС массив, относящийся к тому или иному запросу Googledorks. Пример можно увидеть на рис. 3.

mass = ["allinurl:", query]

Рис. 3. Пример интеграции одного из запросов Googledorks, относящийся к техникам OSINT

В дальнейшем для поиска информации ИАС переведет в массив и выполнит запрос через поисковую систему. Метод решения задач в разработанной информационно-аналитической системе работает следующим образом: пользователь запускает ИАС, далее вводит комбинацию для перехода к нужному для него виду поиска, после вводит необходимые данные и ждет ответа системы. Как появится ответ, пользователь решает, хочет ли он сохранить информацию для дальнейшей работы. Далее пользователь может обратиться к выведенной таблице, чтобы изучить дополнительную информацию по введенному запросу. Выведенные в таблицу ссылки он может спокойно изучить, так как они рабочие и относятся к российскому сегменту сети Интернет (Янгаева, Павленко, 2022).

Исследование работоспособности метода в предлагаемой ИАС. Для исследования работоспособности предложенного метода выполним один из запросов в разработанной ИАС (рис. 4,5).

Суть работы ИАС:

- При запуске системы ввести комбинацию кодовых значений, соответствующую нужному запросу.
- Заполнить поля и дождаться результатов обработки данных в виде таблицы с текстом, извлеченным из сайта, выбранного случайным образом.
- Решить, нужно ли сохранить извлеченную информацию для дальнейшего анализа и рассмотреть остальные ссылки для дополнения к полученной информации.

```
Windows PowerShell
(C) Корпорация Майкрософт (Microsoft Corporation). Все права защищены.

Попробуйте новую кроссплатформенную оболочку PowerShell (https://aka.ms/pscore6)

PS C:\Users\Илья\PcharmProjects\pythonProject18> python start.py
Welcome to contrrus
You can get:
general search
cache
allintext
allintitle
allinurl
site
inposttitle
source
 filetype
organisation search
help
```

Рисунок 4. Ввод данных для поиска информации по общему запросу

```
Enter code:576444567879961
Enter topic:кто такое ринофиз
Enter number of exclusive words:0
Enter counts:22
```

a)

```

| com | + |
5 | https://www.smdoctor.ru/disease/zabolevaniya-gipofiza/
| ru | + |
6 | https://www.center-endo.ru/gipofiz-stroenie-i-funkczii-gipofiza/
| ru | + |
7 | https://ru.wiktionary.org/wiki/%D0%B3%D0%B8%D0%BF%D0%BE%D1%84%D0%B8%D0%B7
| org | + |
8 | https://mrt-don.ru/mrt/info/chto-pokazyvaet-mrt-gipofiza/
| ru | + |
9 | https://medportal.ru/enc/endocrinology/gipofiz/organ-gipofiz/
| ru | + |
10 | https://fnkc-fmba.ru/zabolevaniya/adenoma-gipofiza/
| ru | + |
11 | https://sunmedexpert.ru/napravleniya/endokrinologiya/zabolevaniya-gipofiza/
| ru | + |

```

[illegible]

670



```
12 | https://endokrinet.ru/zabolevaniya/adenoma-gipofiza/
13 | https://spb.medi.ru/articles/art-gipofiza/
14 | https://8378.ru/termin/gipofiz
15 | http://www.imun.ru/main/endokrinologiya/gipofiz/
16 | https://kry-mtexpert.ru/articles/235
17 | https://www.smeclinic-spb.ru/doctor/endokrinolog/zabolevaniya/4534-adenoma-gipofiza
18 | https://www.booksite.ru/fulltext/1/003/008/038/038.htm
19 | https://lapino2.ru/napravleniya/opukholi-golovy-i-shei/opukholi-gipofiza/
20 | https://lapino2.ru/napravleniya/opukholi-golovy-i-shei/opukholi-gipofiza/
```

Рис. 5 а), б), с) и d). Результат по введенному запросу общего поиска

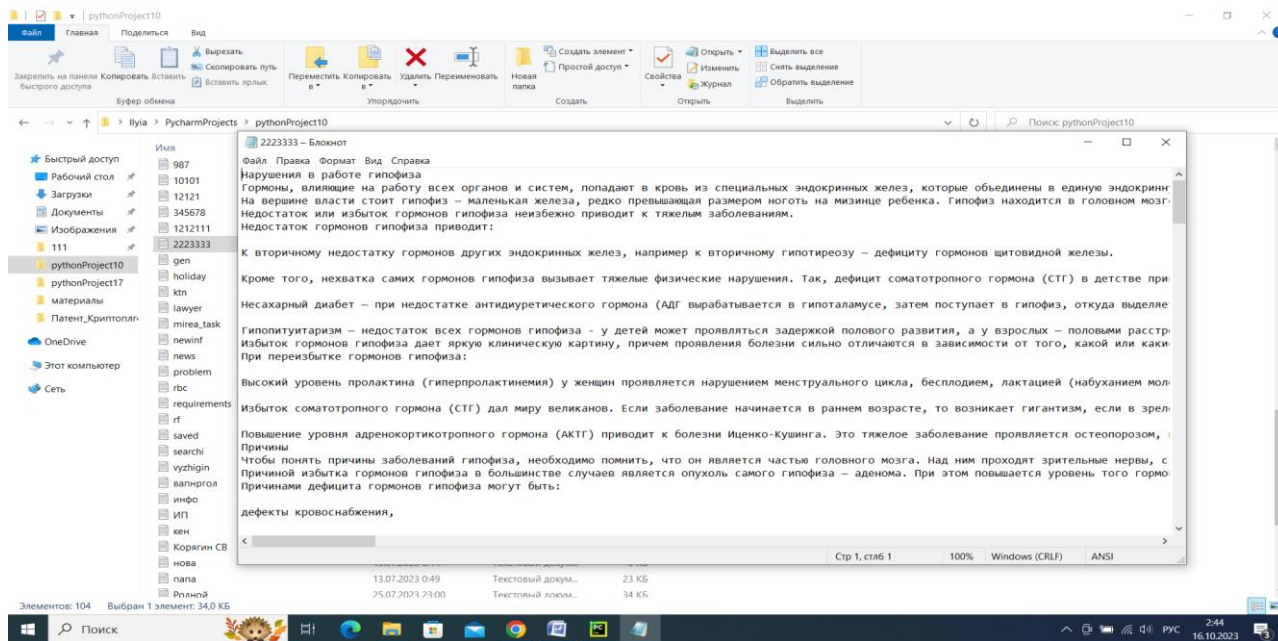
Как ранее было сказано, в случае успешного выполнения запроса появится таблица и номер ссылки, откуда будет извлечена информация. Процесс извлечения информации можно изучить на рис. 6.

**Нарушения в работе гипофиза**  
Гормоны, влияющие на работу всех органов и систем, попадает в кровь из специальных эндокринных желез, которые объединены в единую эндокринную систему. Это надпочечники, щитовидная и паращитовидные железы, яичники (у женщин), семенники и яички – (у мужчин), поджелудочная железа, гипоталамус и гипофиз. Пожалуй, в организме нет более иерархичной и дисциплинированной системы, чем эндокринная.  
На вершине власти стоит гипофиз – маленькая железа, редко превышающая размером ноготь на мизинце ребенка. Гипофиз находится в головном мозге (в самом его центре) и жестко контролирует работу большинства эндокринных желез, выделяя специальные гормоны, которые управляют производством других гормонов. Например, гипофиз выбрасывает в кровь тиреотропный гормон (ТТГ), который заставляет щитовидную железу создавать тироксин и трийодтиронин. Некоторые гормоны гипофиза оказывают непосредственный эффект, например, соматотропный гормон (СТГ), отвечающий за процессы роста и физического развития ребенка.  
Недостаток или избыток гормонов гипофиза неизбежно приводит к тяжелым заболеваниям.  
Недостаток гормонов гипофиза приводит:  
  
К вторичному недостатку гормонов других эндокринных желез, например к вторичному гипотиреозу – дефициту гормонов щитовидной железы.  
Кроме того, нехватка самих гормонов гипофиза вызывает тяжелые физические нарушения. Так, дефицит соматотропного гормона (СТГ) в детстве приводит к карликовости.  
Сахарный диабет – при недостатке антидиуретического гормона (АДГ) вырабатывается в гипоталамусе, затем поступает в гипофиз, откуда выделяется в кровь)  
Гипопитуитаризм – недостаток всех гормонов гипофиза - у детей может проявляться задержкой полового развития, а у взрослых – половыми расстройствами. В целом, гипопитуитаризм ведет к тяжелым нарушениям обмена веществ, которые затрагивают все системы организма.  
Избыток гормонов гипофиза дает яркую клиническую картину, причем проявления болезни сильно отличаются в зависимости от того, какой или какие гормоны превышает норму.  
При переизбытке гормонов гипофиза:  
  
Высокий уровень пролактина (гиперпролактинемия) у женщин проявляется нарушением менструального цикла, бесплодием, лактацией (набуханием молочных желез и секрецией молока). У мужчин гиперпролактинемия ведет к снижению полового влечения, импотенции.  
Избыток соматотропного гормона (СТГ) дал миру великанов. Если заболевание начинается в раннем возрасте, то возникает гигантизм, если в зрелом - акромегалия. Согласно Книге рекордов Гиннесса, самым высоким мужчиной был Роберт Першинг Уодлоу, родившийся в 1918 году в США. Его рост составлял 272 сантиметра (размах рук 288 сантиметров). Однако, по данным отечественной книги рекордов Диво, самым высоким в мировой истории был российский гражданин Федор Махов. Его рост составлял 2 метра 85 сантиметров при весе 182 килограмма. При акромегалии у больного утолщаются кисти рук и ступней, черты лица становятся крупными, увеличивается внутренние органы. Это сопровождается нарушениями работы сердца, неврологическими расстройствами.  
Повышение уровня адренокортикотропного гормона (АКТГ) приводит к болезни Иценко-Кушинга. Это тяжелое заболевание проявляется остеопорозом, повышением артериального давления, развитием сахарного диабета, психическими нарушениями. Болезнь сопровождается характерными изменениями внешности: похуданием ног и рук, ожирением в области живота, плеч, а также лица.  
**Причины**  
Чтобы понять причины заболеваний гипофиза, необходимо помнить, что он является частью головного мозга. Над ним проходят зрительные нервы, с боковых сторон – крупные мозговые сосуды и глазодвигательные нервы.

a)

```
некоторые лекарства,
облучение,
хирургическое вмешательство.
Диагностика и лечение
Диагностикой и лечением заболеваний гипофиза занимается врач-эндокринолог. При первом обращении врач соберет анамнез (жалобы, информацию о перенесенных заболеваниях и наследственной предрасположенности) и на основании этого назначит необходимое исследование гормонального профиля (анализ крови на гормоны), тест с тиролиберинном, тест с синактеном и т.д. При необходимости может быть назначена компьютерная томография головного мозга, магнитно-резонансная томография головного мозга и др.
Лечение заболеваний гипофиза направлено на нормализацию уровня гормонов в крови, а в случае аденомы – уменьшение давления опухоли на окружающие структуры мозга. При недостатке гормонов гипофиза применяется заместительная гормональная терапия: человеку дают лекарства-аналоги нужных гормонов. Такое лечение часто длится пожизненно. К счастью, опухоли гипофиза крайне редко бывают злокачественными. Тем не менее, их лечение – трудная задача для врача.
В лечении опухолей гипофиза используют следующие методы и их сочетание:
лекарственная терапия;
Лечение заболеваний гипофиза направлено на нормализацию уровня гормонов в крови, а в случае аденомы – уменьшение давления опухоли на окружающие структуры мозга. При недостатке гормонов гипофиза применяется заместительная гормональная терапия: человеку дают лекарства-аналоги нужных гормонов. Такое лечение часто длится пожизненно. К счастью, опухоли гипофиза крайне редко бывают злокачественными. Тем не менее, их лечение – трудная задача для врача.
В лечении опухолей гипофиза используют следующие методы и их сочетание:
Заявка на прием
Заполните все обязательные поля, пожалуйста
Save to file?{y/n}
name:111
This file exists
name:2223333
File was written
.
Go ones more yes/no
no
Thanks!
PS C:\Users\Ilna\PycharmProjects\pythonProject18>
```

b)



с)

Рис. 6 а), б) и с). Извлечение и запись информации в файл

В разработке любого программного обеспечения одним из важных этапов является его тестирование. В данном случае тестирование выполнялось в виде подбора темы запроса в ИАС, а также определения максимального количества ссылок, которое является очень важным моментом в тестировании ИАС, поскольку количество ссылок – непосредственная нагрузка на сеть. ИАС тестировалась на разные нагрузки, каждая из которых давалась по несколько раз, чтобы выяснить, какие исключения от интернета может получить пользователь, чтобы сделать работу ИАС более качественной. В качестве подтверждения проведенных экспериментов предоставим пример исключений, которые были получены на сегодняшний день (рис. 7) (Jason, Jeffry, 2018).

```
except urllib.error.HTTPError:
    print("too many requests!")
except requests.exceptions.ContentDecodingError:
    print("decoding error!")
except requests.exceptions.InvalidSchema:
    print("no connection adapters were found")
except UnicodeEncodeError:
    print("cannot encode!!!")
except http.client.RemoteDisconnected:
    print("http mistake")
except http.client.IncompleteRead:
    print("wrong size of bytes")
except requests.exceptions.ChunkedEncodingError:
    print("wrong size of bytes")
except urllib3.exceptions.ProtocolError:
    print("wrong size of bytes")
```

Рис. 7. Исключения от протоколов интернета, полученные в ходе работы ИАС

### Заключение

Проведенные эксперименты доказали эффективность использования методов текстовой аналитики за счет применения практик OSINT. Цель работы была достигнута. Разработан метод извлечения информации из российского сегмента сети Интернет. Разработан метод проверки запросов к ИАС на их логичность с точки зрения лексикологии. Разработан защищенный доступ к функционалам ИАС путем применения комбинаций кодовых значений, что повышает безопасность доступа со стороны пользователей. Функционал информационно-аналитической системы позволяет решать профессиональные задачи.

### Список литературы



1. Брумштейн Ю.М., Васильковский Е.Ю., Куаншкалиев Т.Х. Поиск информации в Интернете: анализ влияющих факторов и моделей поведения пользователей // Известия Волгоградского государственного технического университета. 2017. № 1. С. 50-55.
2. Бельдеубаева Д.Р. Применение OSINT технологий в качестве повышения эффективности деятельности органов внутренних дел // Актуальные вопросы эксплуатации систем охраны и защищенных телекоммуникационных систем. Воронеж. 2020. С. 160-161.
3. Голушко А.П., Дрянных Ю.Ю. Цель и задача поиска информации в открытых источниках (open source intelligence) // Внедрение результатов инновационных разработок: Проблемы и перспективы. Общество с ограниченной ответственностью "Агентство международных исследований". Уфа. 2019. С. 158-161.
4. Москалев И.С. Применение технологии OSINT для получения информации по IP-адресу // IT OPEN 2022. С. 187-191.
5. Рассел Д. Контент-анализ. М.: VSD, 2021. 931 с.
6. Таршис Е.Я. Контент-анализ. Принципы методологии. Построение теоретической базы. Онтология, аналитика и феноменология текста. Программы исследования. М.: ИЛ, 2020. 786 с.
7. Янгаева М.О., Павленко Н.О. OSINT. Получение криминалистически значимой информации из сети интернет // Алтайский юридический вестник. 2022. № 2. С. 131-135.
8. Batrinca B., Philip C.T. Social media analytics: a survey of techniques, tools and platforms // AI & SOCIETY. 2015. № 30. С. 89-116.
9. GoogleDorking или используем Гугл на максимум // Habr. 2020. URL: <https://habr.com/ru/companies/postuf/articles/510766/>
10. Jason H.S., Jeffrey S.B. Is Information Systems Late to the Party? The Current State of DevOps Research in the Association for Information Systems eLibrary // DevOps Research in the AISel. 2018. №. С. 1-8.

#### **Application of text analytics algorithms and OSINT practices in information and analytical systems**

##### **Alexander Yu. Vyzhigin**

Candidate of Technical Sciences, Associate Professor, Head of Department  
Russian Academy of National Economy and Public Administration under the President of the Russian Federation  
Moscow, Russia  
Candidate of Technical Sciences, Associate Professor  
Russian Technological University  
Moscow, Russia  
[vijigin\\_new2000@mail.ru](mailto:vijigin_new2000@mail.ru)

##### **Ilya S. Moskalev**

Student  
Russian Technological University  
Moscow, Russia  
[moskalevilya1@gmail.com](mailto:moskalevilya1@gmail.com)

##### **Oleg V. Trubienko**

Candidate of Technical Sciences, Associate Professor, Head of Department  
Russian Technological University  
Moscow, Russia  
[trubienko@mail.ru](mailto:trubienko@mail.ru)

Received 13.10.2023

Accepted 10.11.2023

#### **Annotation**

Information is an important resource of our life, which is usually divided into two types depending on its significance: 1) accurate data – data that does not distort events, facts and helps to conduct a competent analysis of what happened; 2) redundant (garbage) data – data that can mislead the analyst, thus creating the risk of errors during the analysis. The problem of the modern information society is that people cannot always extract accurate data for subsequent analysis of information. Information and analytical systems are able to cope with this problem (a special class of information systems that have the property of analytical data processing when automating the process of receiving a response to a particular request). The work of information and analytical systems (IAS) can be compared with information systems (IS). Consider one of the well-known search engines – Yandex. Search engines are also information and analytical. For example, we want to get information about the 2014 FIFA World Cup. In Yandex, we generate a query – "2014 FIFA World Cup",

which in SQL would look like this (SELECT football FROM world\_championship WHERE year = 2014). The response from Yandex contains 5,000 results, and in order to find out the necessary information, it is necessary either to analyze these answers or to reformulate search queries. Search engines are similar to IAS, but for an optimal solution it is worth using those that contain a user-friendly interface where you can click on the buttons from the menu to get to the desired functionality and get the desired result without additional analysis. The scientific novelty of the research is as follows: 1) the possibilities of text analytics have been expanded by using OSINT practices when making requests to the IAS; 2) a method for extracting information from the Russian segment of the Internet has been developed; 3) a method for checking requests to the IAS for their consistency from the point of view of lexicology has been developed; 4) Secure access to the IAS functionality has been developed. The purpose of the work is to increase the effectiveness of text analytics methods through the use of OSINT practices. The article used such methods as: OSINT (search for information from open sources, a business intelligence method, the essence of which is to get a specific answer to a specific question); text analytics (the process of processing unstructured text to identify ideas, patterns, etc.). As a result, the functionality of an information and analytical system using text analytics algorithms and OSINT practices was developed. The results obtained can be used by employees of various security services of our country, as well as by ordinary users of this IAS to quickly obtain the necessary information.

### Keywords

information analytical system, text analytics, NLP, Open source intelligence, OSINT, Googledorks query language, Python programming language, morphological analyzer, data science, search engines.

### References

1. Brumshtejn YU.M., Vas'kovskij E.YU., Kuanshkaliev T.H. Poisk informacii v Internetе: analiz vliyayushchih faktorov i modelej povedeniya pol'zovatelej // Izvestiya Volgogradskogo gosudarstvennogo tekhnicheskogo universiteta. 2017. № 1. S. 50-55.
2. Bel'deubaeva D.R. Primenenie OSINT tekhnologij v kachestve povysheniya effektivnosti deyatel'nosti organov vnutrennih del // Aktual'nye voprosy ekspluatatsii sistem ohrany i zashchishchennyh telekommunikacionnyh sistem. Voronezh. 2020. S. 160-161.
3. Golushko A.P., Dryannyh YU.YU. Cel' i zadacha poiska informacii v otkrytyh istochnikah (open source intelligence) // Vnedrenie rezul'tatov innovacionnyh razrabotok: Problemy i perspektivy. Obshchestvo s ogranichennoj otvetstvennost'yu "Agentstvo mezhdunarodnyh issledovaniy". Ufa. 2019. S. 158-161.
4. Moskalev I.S. Primenenie tekhnologii OSINT dlya polucheniya informacii po IP-adresu // IT OPEN 2022. S. 187-191.
5. Rassel D. Kontent-analiz. M.: VSD, 2021. 931 c.
6. Tarshis E.YA. Kontent-analiz. Principy metodologii. Postroenie teoreticheskoy bazy. Ontologiya, analitika i fenomenologiya teksta. Programmy issledovaniya. M.: IL, 2020. 786 c.
7. YAngaeva M.O., Pavlenko N.O. OSINT. Poluchenie kriminalisticheski znachimoy informacii iz seti internet // Altajskij yuridicheskij vestnik. 2022. № 2. S. 131-135.
8. Batrinca B., Philip C.T. Social media analytics: a survey of techniques, tools and platforms // AI & SOCIETY. 2015. № 30. S. 89-116.
9. GoogleDorking ili ispol'zuem Gugl na maksimum // Habr. 2020. URL: <https://habr.com/ru/companies/postuf/articles/510766/>
10. Jason H.S., Jeffrey S.B. Is Information Systems Late to the Party? The Current State of DevOps Research in the Association for Information Systems eLibrary // DevOps Research in the AISel. 2018. №. S. 1-8.