

УДК 004.89

## Современные методы обработки естественного языка

Б. О. Близнюк, Л. В. Васильева, И. Д. Стрельников, Д. С. Ткачук  
*Харьковский национальный университет  
имени В. Н. Каразина, пл. Свободы, 4, 61022, Харьков, Украина*

В данной статье рассмотрены основные проблемы обработки естественного языка. Проанализированы основные направления обработки, методы, инструменты и библиотеки, доступные на текущий момент времени. Проведено два эксперимента, где данные методики были применены для решения реальных задач – анализ тональности новостного фона некоторых криптовалют с целью определения зависимостей между новостным фоном и их обменным курсом и анализ фактов о сотрудничестве компаний, используя их упоминания в различных пресс-релизах. Показано, что анализ текстовых данных имеет большое практическое значение в современном мире.

**Ключевые слова:** обработка естественного языка, анализ текста, обработка текста, анализ тональности текста, классификация, нейронные сети, сбор данных.

У даній статті розглянуті основні проблеми обробки природної мови. Проаналізовано основні напрямки обробки, методи, інструменти та бібліотеки, що наразі доступні. Проведено два експерименти, де дані методики були застосовані для вирішення реальних завдань - аналіз тональності новин щодо деяких криптовалют з метою визначення залежностей між новинами і обмінним курсом і аналіз фактів про співпрацю компаній, використовуючи їх згадки в різних прес-релізах. Показано, що аналіз текстових даних має велике практичне значення в сучасному світі.

**Ключові слова:** обробка природної мови, аналіз тексту, обробка тексту, аналіз тональності тексту, класифікація, нейронна мережа, збір даних.

The main challenges of natural language processing have been covered in the article. The main processing tasks, methods and libraries presently available have been analyzed. Two experiments have been carried out, where these techniques have been used to solve real-life problems, namely, the analysis of the internet news concerning some cryptocurrencies to see if the sentiment of those news correlated with the prices of the cryptocurrencies and the extraction of facts from the various press-releases to find out the companies with established partnerships. It has been shown that natural language processing is a very important and powerful tool in the modern age.

**Key words:** Natural Language Processing, text analysis, text processing, sentiment analysis, classification, neural network, data mining.

### 1 Введение

Мы живем во время, когда объемы производимой человечеством информации больше, чем когда либо и количество этих данных растет с каждым днем. Однако значительную пользу из этой информации можно извлечь лишь при правильной обработке и анализе этих данных.

Сейчас ежесекундно по всему миру создаются гигабайты новых данных различного вида: делаются новые снимки, видеозаписи, пишутся сотни отзывов к товарам в интернет-магазинах, тысячи комментариев под записями на Facebook, десятки рецензий к фильмам в онлайн-кинотеатрах, цены на акции то

взлетают, то падают. И большая часть этих данных в “сыром” виде практически бесполезна. Чтобы извлечь из них какую-то пользу, их нужно отфильтровать и обработать. Во времена, когда технологии еще не были так развиты, все это приходилось делать вручную. На это уходили часы, дни, недели, а то и месяцы. А если учесть, что раньше и самой информации для обработки было в разы меньше, то несложно понять, что сейчас обрабатывать такие объемы вручную просто невозможно. Поэтому было разработано множество алгоритмов, которые позволяют делать это при помощи компьютерной техники. Именно о таких методах, касающихся обработки естественного языка и пойдет речь в данной работе.

## **2 Основные направления**

Обработка естественного языка (Natural Language Processing, NLP) — общее направление искусственного интеллекта и математической лингвистики. Оно изучает проблемы компьютерного анализа и синтеза естественных языков. Применительно к искусственному интеллекту, анализ означает понимание языка, а синтез – генерацию грамотного текста. Решение этих проблем будет означать создание более удобной формы взаимодействия компьютера и человека.[1] К основным направлениям обработки естественного языка относят такие, как извлечение фактов, анализ тональности текста, ответы на вопросы, информационный поиск, генерация текста, перевод и т.д. Подробнее о некоторых из них.

### **2.1 Извлечение информации или фактов**

Под извлечением информации подразумевается поиск в неструктурированном или слабо структурированном документе отдельных интересующих вас фактов. Например, у вас имеется огромное количество статей, в которых фигурирует большое количество разных личностей, и вы хотите составить базу данных, которая будет хранить данные о том, кто из фигурирующих в данных статьях людей, являются мужем и женой. Данный пример был использован для демонстрации[2] возможностей программы под названием DeepDive, созданной группой студентов и работников университета Stanford.

### **2.2 Анализ тональности текста**

Анализ тональности текста подразумевает под собой автоматическое определение эмоциональной окраски текста и выявление отношения человека, написавшего текст, к объекту обсуждения. Данный тип анализа может быть использован, например, продавцами для того, чтобы лучше понять, какой из продаваемых им товаров пользуется большим успехом среди покупателей, анализируя отзывы. Также его могут использовать власти для выявления отношения к ним и их решениям граждан страны и т.д. В наше время наиболее часто используемыми в исследованиях методами являются методы на основе машинного обучения с учителем. Сутью таких методов является то, что на первом этапе обучается машинный классификатор (например, байесовский) на

заранее размеченных текстах, а затем используют полученную модель при анализе новых документов.[3]

### **2.3 Ответы на вопросы**

Под данное определение в целом могут подходить и так называемые чат-боты, которые имитируют реальное общение с людьми посредством передачи текстовых сообщений, и специальные программы, которые сперва анализируют некий текст, а после – отвечают на вопросы, связанные с его содержанием. Результаты одного из последних исследований на эту тему, описаны в статье [4], под авторством John Ball.

### **2.4 Перевод текста**

Также одним из наиболее известных и часто используемых направлений обработки естественного текста является его перевод с одного естественного языка на другой. Одной из наиболее продвинутых методик, используемых в данный момент для достижения правильного перевода, является использование нейронных сетей типа “Seq2Seq”[5] с “Вниманием”[16], что расшифровывается, как “sequence to sequence”, или “последовательность в последовательность”. Позже, мы более подробно рассмотрим данный тип нейронных сетей.

## **3 Способы анализа**

Для решения вышеописанных задач, исследователи пользуются огромным набором инструментов и техник анализа естественного языка. Некоторые из них узкоспециализированные, как seq2seq, другие же могут применяться в различных ситуациях, как Word2Vec, о котором дальше и пойдёт речь.

### **3.1 Word2Vec**

В основе данной технологии лежит представление слов в виде векторов заданной размерности, располагая похожие слова близко друг к другу. То есть, расстояние между векторами слов, обозначающих похожие вещи, например, “кот” и “собака”, будет значительно меньше, чем между словами, значения которых имеют мало общего, например, “кот” и “самолёт”. Данная особенность позволяет более гибко представлять данные, которые в дальнейшем могут быть использованы в обучении нейронных сетей, различных классификаторов и т.д.

Для создания базы соответствий “слово - вектор”, алгоритм сначала просматривает весь выданный ему текст, составляя “словарь”, который в последующих итерациях работы алгоритма, будет использован для определения соответствующих векторов. Существует два основных подхода: CBOW (Continuous Bag of Words) и Skip-gram. CBOW – «непрерывный мешок со словами» модельная архитектура, которая предсказывает текущее слово, исходя из окружающего его контекста. Архитектура типа Skip-gram действует иначе: она использует текущее слово, чтобы предугадывать окружающие его слова.[6]

### **3.2 Определение структуры текста**

Все тексты на естественном языке имеют большое количество слов, которые не несут информации о данном тексте. К примеру, в английском языке такими

словами являются артикли, в русском к ним можно отнести предлоги, союзы, частицы. Данные слова называют шумовыми или стоп-словами. Для достижения лучшего качества классификации на первом этапе предобработки текстов обычно необходимо удалять такие слова. Второй этап предобработки текстов – приведение каждого слова к основе, одинаковой для всех его грамматических форм. Это необходимо, так как слова несущие один и тот же смысл могут быть записаны в разной форме. Например, одно и то же слово может встретиться в разных склонениях, иметь различные приставки и окончания.

### 3.3 Нейронные сети

Искусственные нейронные сети представляют собой систему соединённых и взаимодействующих между собой простых процессоров – искусственных нейронов. Алгоритм работы таких процессоров зачастую крайне прост. Например, процессор может просто преобразовать полученный на входе сигнал, используя некую математическую функцию, в выходной. И, тем не менее, будучи соединёнными в достаточно большую сеть с управляемым взаимодействием, такие по отдельности простые процессоры вместе способны выполнять довольно сложные задачи.

#### *RNN/LSTM*

Рекуррентные нейронные сети [7] отличаются от другого типа сетей тем, что кроме связей, переходящих от одного нейрона к другому напрямую, как в сетях прямого распространения, а также связи, проходящие “во времени”. То есть, сигнал от одного нейрона на этапе  $t$  перейдёт к другому (или этому же) нейрону на этапе  $t+1$ . Таким образом рекуррентные нейронные сети могут сохранять информацию во времени, тем самым “запоминая” некоторые данные. Данная их особенность как раз очень сильно помогает в переводе, классификации и обработке природного текста в целом, так как наш язык устроен таким образом, что некоторые данные в начале блока текста, могут повлиять на понимание и/или перевод в его конце. В данной [17] статье отлично показаны реальные возможности рекуррентных нейронных сетей и их применения.

#### *CNN*

СНС - сверточные нейронные сети лучше всего показали себя в распознавании объектов и образов на картинках, классификации изображений, выделении особенностей и сжатии данных. Однако, им нашлось применение и в обработке текста.

#### *Посимвольный подход.*

Посимвольный подход для классификации текста с помощью сверточных нейронных сетей был предложен в статье [16]. Подробнее опишем данный метод. Алфавит - это упорядоченный набор символов. Пусть данный алфавит состоит из  $m$  символов. Каждый символ алфавита в тексте закодирован с помощью  $1 - m$  - кодировки. (т. е. каждому символу будет сопоставлен вектор длины  $m$  элемент которого равен единице, в позиции равной порядковому номеру символа в алфавите, а нулю во всех остальных позициях.) В случае, когда в тексте встретится символ, который не вошел в алфавит, то его необходимо закодировать вектором длины  $m$  состоящим из одних нулей. Далее, из текста следует выбрать первые  $l$  символов. Параметр  $l$  должен быть большим,

чтобы в первых  $l$  символах содержалось достаточно информации для определения класса всего текста.

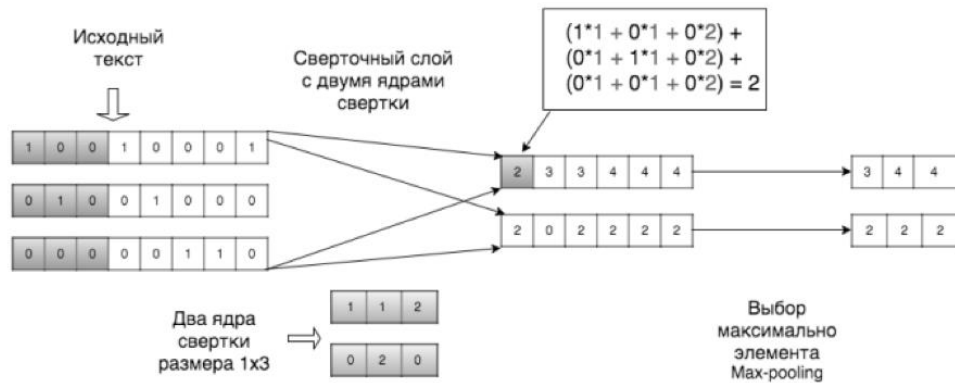


Рис. 1 Посимвольный подход

Затем исходные векторы объединяются в матрицу размера  $m \times l$ , в которой в каждый столбец будет иметь не более одной единицы. Каждая строка полученной матрицы используется как отдельная карта признаков. На вход сверточной нейронной сети подается  $m$  карт признаков размера  $1 \times l$  аналогично изображению. Архитектуру сети необходимо выбирать исходя из задачи. На Рис. 1 приведен пример посимвольного подхода для  $l = 6$ ,  $m = 3$ . В примере показан один сверточный и один субдискретизирующий слой.

*Подход с использованием кодирования слов.*

В данном подходе каждому слову в тексте сопоставляется вектор фиксированной длины, затем из полученных векторов для каждого объекта выборки составляется матрица, которая аналогично изображениям подается на вход сверточной нейронной сети. На Рис. 2 приведен пример сверточной нейронной сети с использованием кодирования слов.

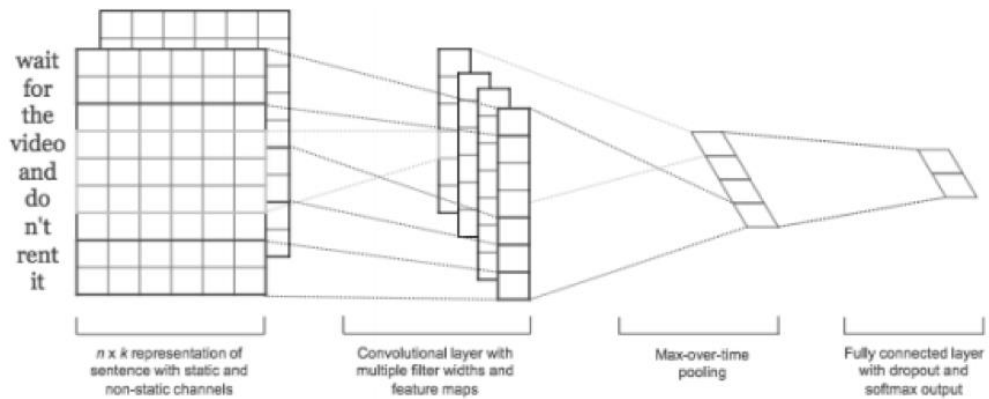


Рис. 2 Подход с использованием кодирования слов

### Seq2Seq

Універсальна бібліотека кодировщик-декодер для Tensorflow[5], которая может использоваться для машинного перевода, определения содержания текста, моделирования диалогов, описания содержания изображений и т. д. Seq2Seq позволяет создавать и обучать модели нейронных сетей вида 'sequence to sequence'.

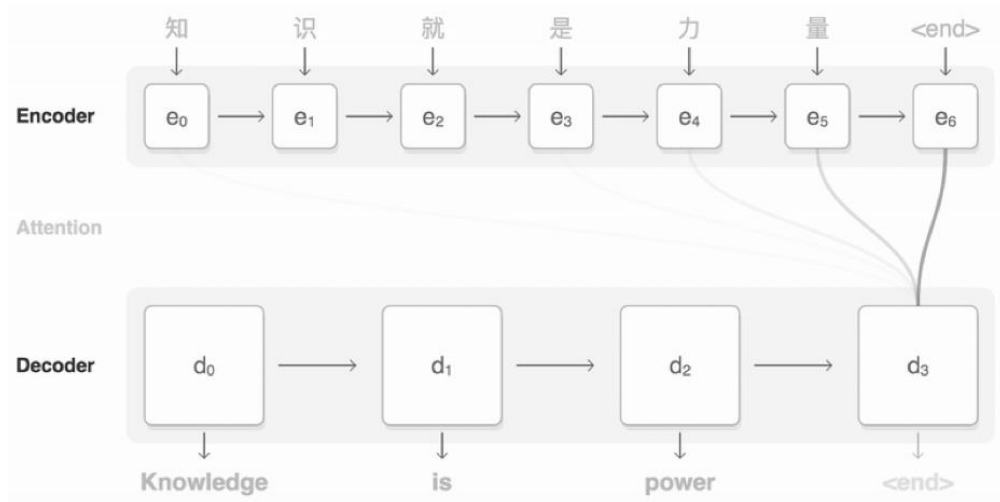


Рис. 3 - Визуализация принципа работы библиотеки Seq2Seq, источник: <https://github.com/google/seq2seq>

## 4 Другие классификаторы

### 4.1 Байесовский классификатор

Наивный байесовский классификатор – простой вероятностный классификатор, основанный на применении теоремы Байеса со строгими (наивными) предположениями о независимости.

Достоинством наивного байесовского классификатора является малое количество данных для обучения, необходимых для оценки параметров, требуемых для классификации.

Пусть  $X$  – множество описаний объектов,  $Y$  – множество наименований классов. На множестве пар «объект, класс» определена вероятностная мера  $P$ . Имеется конечная обучающая выборка независимых наблюдений  $X^m = \{(x_1, y_1), \dots, (x_m, y_m)\}$ , полученных согласно вероятностной мере  $P$ .

При работе с непрерывными данными основное предположение состоит в том, что непрерывные значения, связанные с каждым классом, распределяются в соответствии с распределением Гаусса.

### 4.2 SVM/SVC

Метод опорных векторов (англ. SVM, support vector machine) – набор схожих алгоритмов обучения с учителем, использующихся для задач классификации и

регрессионного анализа. Особым свойством метода опорных векторов является непрерывное уменьшение эмпирической ошибки классификации и увеличение зазора, поэтому метод также известен как метод классификатора с максимальным зазором.

Основная идея метода – перевод исходных векторов в пространство более высокой размерности и поиск разделяющей гиперплоскости с максимальным зазором в этом пространстве. Две параллельных гиперплоскости строятся по обеим сторонам гиперплоскости, разделяющей классы. Разделяющей гиперплоскостью будет гиперплоскость, максимизирующая расстояние до двух параллельных гиперплоскостей. Алгоритм работает в предположении, что чем больше разница или расстояние между этими параллельными гиперплоскостями, тем меньше будет средняя ошибка классификатора.

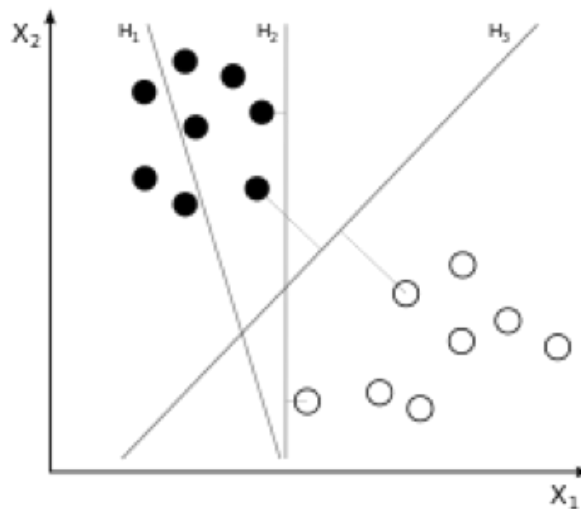


Рис. 4  $H_1, H_2, H_3$  – гиперплоскости.  $H_3$  – гиперплоскость максимальной разности

## 5 Существующие библиотеки для NLP

### 5.1 Stanford CoreNLP

Данный бесплатный программный продукт был создан общими усилиями студентов и научных работников университета Stanford. Основной задачей, которая была поставлена в начале его разработки, было создание набора современных инструментов, позволяющих обрабатывать неструктурированный текст. По сей день этот продукт является одним из лучших в своей нише. С его помощью можно провести полный анализ частей речи в тексте, структуры текста, провести распознавание именованных объектов, определить, где в тексте разные существительные обозначают один и тот же объект, провести анализ тональности текста и многое другое. Больше информации о продукте и примеры использования можно найти на официальном сайте [12].

### 5.2 Natural Language Toolkit, NLTK

Данная бесплатная библиотека для языка программирования Python является одной из лучших для создания различных программных продуктов на этом языке. Она предоставляет большой набор инструментов, корпусов текста, имеет предусмотренные обертки для использования других библиотек внутри себя. Например, для анализа тональности текста и разметки предложений, есть возможность подключения вышеописанного продукта Stanford CoreNLP. Также для различных классификаций в NLTK был предусмотрен интерфейс для подключений классификаторов из другой библиотеки – Scikit learn, о которой пойдёт речь дальше. А больше информации об использовании, устройстве данной можно найти на официальном сайте[13].

### 5.3 Scikit learn

Хотя эта библиотека не имеет никаких специфических инструментов для обработки природного языка, в ней имеется огромное количество классификаторов, основанных на различных алгоритмах; моделей нейронных сетей и прочих общих инструментов для машинного обучения. Используя её вместе с другими, узконаправленными инструментами, можно создавать очень сложные и качественные системы для обработки, анализа, классификации и даже генерации природного текста. Ярким примером использования данной библиотеки является следующий[14] обучающий материал, написанный человеком под ником sentdex.

### 5.4 Tensorflow

Также как и вышеописанная библиотека, Tensorflow[15] абсолютно не задумывалась, как библиотека для обработки природного текста, а как инновационная библиотека для машинного обучения и, в особенности, искусственных нейронных сетей. Но она также может быть использована для анализа и генерации природного текста. Вышеописанный seq2seq, как и Word2Vec, являются частью данной библиотеки, постоянно дорабатываются и уже успешно используется для различных исследований и создания программных продуктов. В данной [5] обучающей статье можно прочитать о том, как можно использовать модель seq2seq для перевода текста.

## 6 Практическое применение анализа текста

В течение всего 2017 года технология блокчейн стремительно набирала обороты. Цены на криптовалюты очень сильно отличаются друг от друга и завязаны на огромном количестве факторов. Таких, как общее количество “монет” или “токенов”, наличие верхней границы их количества, способ “добычи” данной валюты и многие другие. Так как в целом технология всё ещё очень нова, движения цен на данный тип валют в основном зависят от одного ключевого фактора – мнения людей о данной валюте, которое формируется исходя из новостей ключевых крипто-изданий, анонсов от разработчиков, информации о партнерстве с другими организациями и остальном новостном фоне. Хороший новостной фон вокруг валюты приводит к повышенному вниманию со стороны общественности, что привлекает больше



заинтересованных людей, что повышает спрос, и, как правило, позитивно сказывается на стоимости данной валюты.

Обратную ситуацию мы можем наблюдать при негативном новостном фоне вокруг какой-то криптовалюты – плохие новости зачастую негативно влияют на динамику роста данной валюты.

Однако, данные выводы были сделаны в ходе наблюдений за рынком и не были подкреплены никакими статическими данными. Для подтверждения гипотезы было решено спроектировать и реализовать систему, которая позволила бы проанализировать новостной фон некоторых криптовалют и сравнить результаты анализа с графиком изменения стоимости.

В ходе предварительного исследования было найдено несколько сервисов, предоставляющих анализ тональности текста, среди которых Sentiment140, sentiment\_viz, AlchemyAPI и другие. Функционал данных сервисов довольно ограничен и позволяет проанализировать лишь небольшое количество записей.

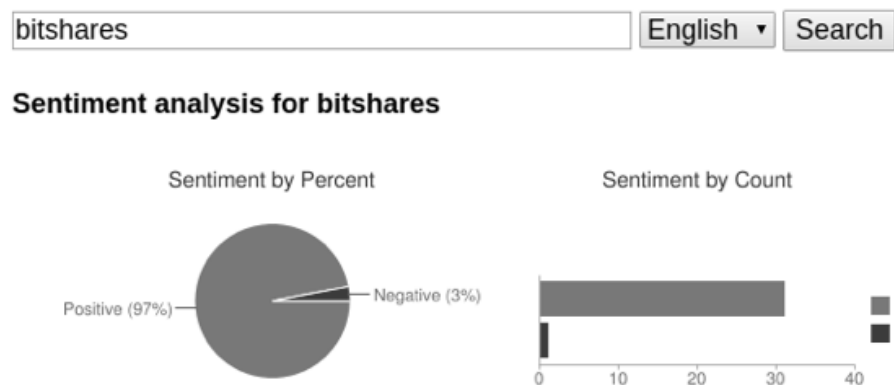


Рис. 5 – Пользовательский интерфейс сервиса Sentiment140

Нами была спроектирована система для обработки новостей, статей, комментариев и прочих текстовых данных, которые непосредственно касались криптовалют. На момент написания данной статьи система всё еще находится в ранних этапах разработки, однако уже имеется рабочий прототип для анализа тональности текста новостей и анонсов, которые люди размещают на площадке “Twitter”.

Для начальных этапов разработки системы было решено использовать более простые в создании и обслуживании методы, а именно – создать классификатор для определения тональности новостей, который будет включать в себя разные подходы к классификации. Было принято решение для каждой поступающей новости вычислять sentiment score – индикатор, показывающий, насколько позитивный или негативный окрас имеет данная новость.

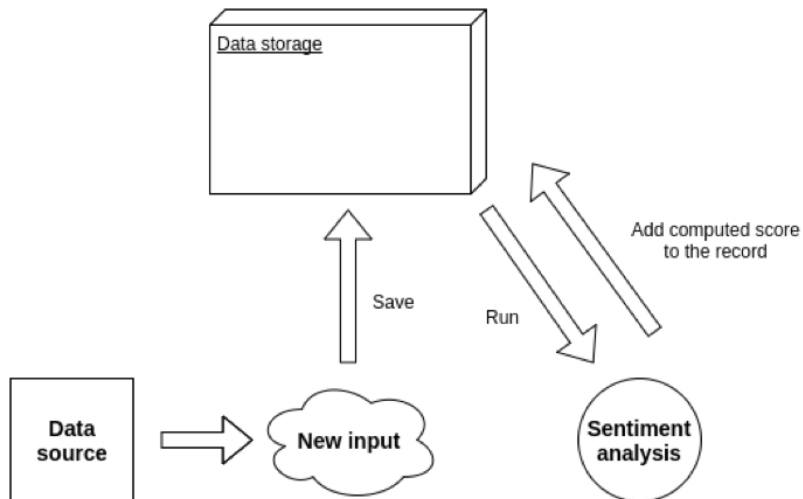


Рис. 6 - Высокоуровневый обзор архитектуры системы для анализа новостей о криптовалютах

Ранее нами было проведено еще одно исследование, в ходе которого была поставлена задача – выяснить, можно ли, используя очень простые правила, извлечь из большого количества неструктурированного текста нужную информацию.

Было решено проверить данное предположение, извлекая факты о сотрудничестве компаний, используя их упоминания в различных пресс-релизах. В исследовании использовались пресс-релизы таких компаний, как Cloudera, Microsoft, Amazon и др.

Для изначальной сборки данных была использована модифицированная версия Apache Nutch. Была скачана информация из общего архива данных, а также более актуальные статьи были собраны непосредственно с официальных сайтов компаний. Ярким примером использованных данных является пресс-релиз компании Cloudera [8].

Система работала следующим образом - на первом этапе были загружены все нужные нам данные, далее данные были помещены Elasticsearch. После, используя написанный нами скрипт на языке Ruby, итеративно, к каждой статье, был применен следующий алгоритм:

1. Проверить правильность регистра слов в заголовке пресс-релиза.
2. Устранить неправильно написанные слова, используя составленные словари, анализ текста самой статьи
3. Проверить, что все оставшиеся слова, написанные с большой буквы, являются названиями компаний, используя поиск по базе данных компаний на Crunchbase [9], Wikipedia [10] и графе знаний Google [11].
4. Внести информацию о данных компаниях в базу данных, если таковых еще нет
5. Проверить заголовок на наличие отобранных нами слов, обозначающих сотрудничество между двумя сторонами

б. Если таковые имеются, то проверить расположение названий компаний и найденных слов относительно друг друга и, если они соответствуют правилам, то внести в базу данных запись о сотрудничестве этих двух компаний

Из-за своей простоты, алгоритм потребовал много времени отладки, заполнения словарей и отбора исключений. Однако, он всё равно дал очень неплохой результат. Из нескольких десятков тысяч статей мы смогли собрать информацию о сотрудничестве более 150 компаний. На Рис. 7 показан скриншот собранных данных в PostgreSQL на ранних этапах работы алгоритма.

	id [PK] serial	names character varying[]		company1_id integer	company2_id integer
1	1	{Cloudera,"Cloudera Inc.",	1	4	5
2	2	{Google,"Google Inc","Goog	2	6	1
3	3	{Dell}	3	7	1
4	4	{Cloudera39s}	4	1	8
5	5	{Hadoop}	5	1	9
6	6	{Talend}	6	1	10
7	7	{Informatica}	7	11	1
8	8	{"Digital Reasoning"}	8	1	12
9	9	{Cseries}	9	13	1
10	10	{"Hadoop MetaScale"}	10	14	1
11	11	{Oracle}	11	15	1
12	12	{Converse}	12	1	16
13	13	{Caggemini}	13	17	1
14	14	{Splunk}	14	1	5
15	15	{FICO}	15	18	1
16	16	{"Cloudera Enterprise"}	16	1	19
17	17	{Persado}	17	1	20
18	18	{Accenture}	18	1	21
19	19	{Udacity}	19	22	1
20	20	{"Red Hat"}	20	1	3
21	21	{"Microsoft Azure"}	21	23	1
22	22	{Telkomsel}	22	24	25
23	23	{NEC}	23	26	1
24	24	{"Cloudera Science Week"}			
25	25	{AsiaPacific}			
26	26	{Experian}			
*					

Рис. 7 – Данные о компаниях

Полученную информацию мы позже преобразовали и внесли в базу данных, основанную на графах – Neo4j. Ниже, на Рис. 8 и Рис. 9, приведены примеры собранных данных.

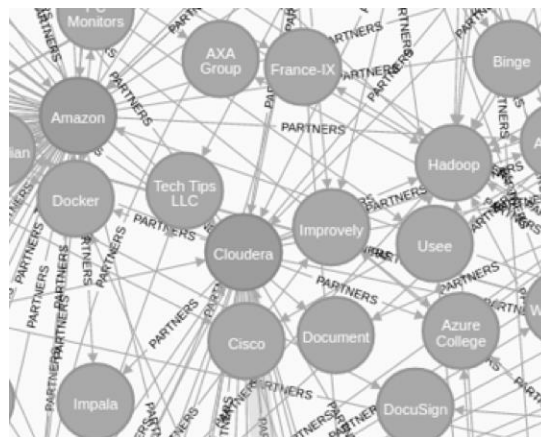


Рис. 8 – Связи между собранными данными

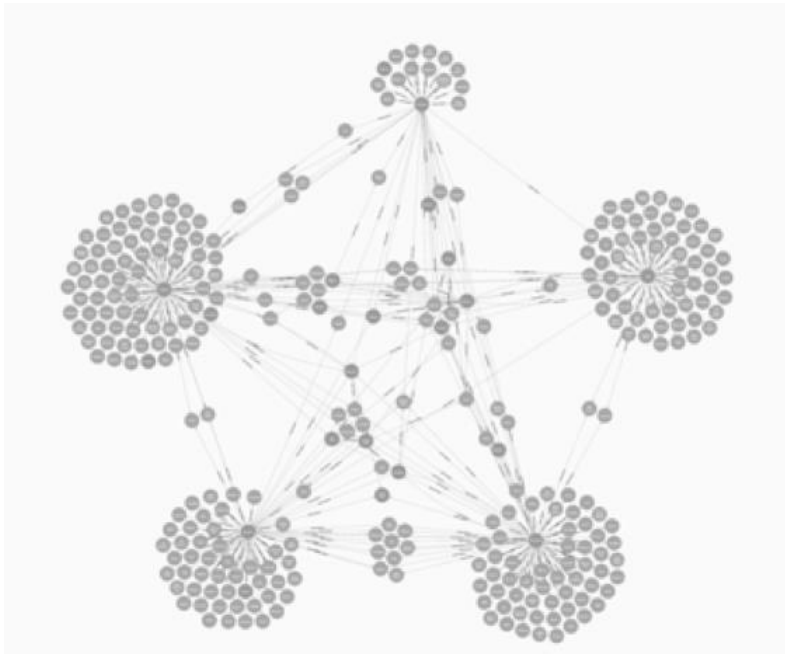


Рис.9 – Связи между собранными данными

## 7 Выводы

В данной работе были рассмотрены основные направления и способы анализа естественного языка, а также методы обработки текстовых данных. Приведено описание существующих инструментов и библиотек для Natural Language Processing. В приведенных примерах используются некоторые из упомянутых в статье методы обработки и классификации текста.

В эксперименте с оценкой котировок криптовалют можно полагать, что есть некоторая зависимость между значением тональности текста и обменным курсом.

В эксперименте с анализом фактов о сотрудничестве компаний с помощью их упоминаний в пресс-релизах был создан граф взаимодействия более 150 компаний.

Подводя итог можно сделать вывод что в современном мире, в среде постоянно растущих объемов информации, анализ текстовых данных имеет большой потенциал и широкое применение.

## ЛИТЕРАТУРА

1. Обработка естественного языка – Режим доступа: [https://ru.wikipedia.org/wiki/Обработка\\_естественного\\_языка](https://ru.wikipedia.org/wiki/Обработка_естественного_языка)
2. DeepDive Tutorial. Extracting mentions of spouses from the news – Режим доступа: <http://deepdive.stanford.edu/example-spouse>
3. Анализ тональности текста – Режим доступа: [https://ru.wikipedia.org/wiki/Анализ\\_тональности\\_текста](https://ru.wikipedia.org/wiki/Анализ_тональности_текста)

4. John S. Ball Using NLU in Context for Question Answering: Improving on Facebook's bAbI Tasks –ARXIV, Электронная версия печ. публикации arXiv:1709.04558, 09/2017 – PDF формат, версия 2 – Режим доступа: <https://arxiv.org/ftp/arxiv/papers/1709/1709.04558.pdf>
5. Neural Machine Translation (seq2seq) Tutorial – Режим доступа: <https://www.tensorflow.org/tutorials/seq2seq>
6. Word2vec – Режим доступа: <https://ru.wikipedia.org/wiki/Word2vec>
7. LSTM – сети долгой краткосрочной памяти – Режим доступа: <https://habrahabr.ru/company/wunderfund/blog/331310/>
8. Cloudera Broadens its Collaboration with Thorn to Include Software and Services to Fight Child Sexual Exploitation – Режим доступа: <https://www.cloudera.com/more/news-and-blogs/press-releases/2016-09-28-cloudera-broadens-its-donation-to-thorn-to-include-software-services-fight-child-sexual-exploitation.html>
9. Crunchbase – Режим доступа: <https://www.crunchbase.com/>
10. Wikipedia – Режим доступа: <https://www.wikipedia.org/>
11. Knowledge – Inside Search – Google - Режим доступа: <https://www.google.com/intl/bn/insidesearch/features/search/knowledge.html>
12. Stanford CoreNLP – Режим доступа: <https://stanfordnlp.github.io/CoreNLP/>
13. Natural Language Toolkit – Режим доступа: <http://www.nltk.org/>
14. Creating a module for Sentiment Analysis with NLTK – Режим доступа: <https://pythonprogramming.net/sentiment-analysis-module-nltk-tutorial/>
15. TensorFlow – Режим доступа: <https://www.tensorflow.org/>
16. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin Attention Is All You Need – ARXIV, Электронная версия печ. публикации arXiv:1706.03762, 06/2017 – PDF формат, версия 5 – Режим доступа: <https://arxiv.org/pdf/1706.03762.pdf>
17. Zhang, X. Character-level convolutional networks for text classification / Xiang Zhang, Junbo Zhao, Yann LeCun // In Advances in Neural Information Processing Systems. — 2015. — Feb. — 649 - 657 p.
18. Andrej Karpathy The Unreasonable Effectiveness of Recurrent Neural Networks – 04/2015– Режим доступа: <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>