

## АНАЛИЗ ТЕХНОЛОГИИ КЛАССИФИКАЦИИ ТЕКСТА

ANALYSIS OF TEXT  
CLASSIFICATION TECHNOLOGYZheng Jing  
Wei Xiaoyu

*Summary.* The article provides an overview of the main approaches to the analysis of textual information. Particular attention is paid to Text Mining technologies and a bag (or cloud) of words. An algorithm for the complex classification of texts is considered. The software solution for constructing the semantic core of the text in Python based on the collections library is described.

*Keywords:* classification of texts, semantic core, signs of documents, key phrases, distributive semantics, text mining, bag (or cloud) of words.

Чжэн Цзини

МГТУ им. Н.Э. Баумана  
sofazjy@gmail.com

Вэй Сяюй

МГТУ им. Н.Э. Баумана  
569006420@mail.ru

*Аннотация.* В статье дается обзор основных подходов к анализу текстовой информации. Особое внимание уделяется технологиям TextMining и мешка (или облака) слов. Рассмотрен алгоритм комплексной классификации текстов. Описывается программное решение построения семантического ядра текста на языке Python на базе библиотеки collections.

*Ключевые слова:* классификация текстов, семантическое ядро, признаки документов, ключевые фразы, дистрибутивная семантика, text mining, мешок (или облако) слов.

**С**ырые неструктурированные данные составляют не менее 90% информации, с которой имеют дело пользователи. Однако все эти огромные запасы данных бесполезны при отсутствии удобных инструментов поиска и/или быстрой классификации информации.

Вопросы классификации текстовой информации изучаются давно и плодотворно [см., например, 1–7]. Существующие сегодня системы классификация применяется, например, в таких задачах как группировка документов в intranet-сетях и на Web-сайтах, размещение документов в определенных папках, сортировка сообщений электронной почты, персонализированная доставка новостей подписчикам.

Задача классификации (или рубрикации, или кластеризации, или распознавания) имеет две основные постановки — бинарную и мультиклассовую.

Бинарная классификация даёт ответ на какой-то один вопрос исследования в стиле «да/нет», например, является ли данное электронное письмо спамом или нет, или представляет ли литературный источник научный интерес или нет. В этом случае говорят о байесовской или наивной классификации информации.

Мультиклассовая классификация ставит перед исследователем задачу отнесения тематики документа к одному из тематических классов предметной области. Количество классов может в некоторых случаях достигать нескольких десятков, количество объектов и их

атрибутов может быть очень большим; поэтому должны быть предусмотрены интеллектуальные механизмы оптимизации процесса классификации.

Мультиклассовая классификация может быть статической, в которой выходные категории заранее предопределены, или динамической, когда рубрики классификации формируются динамически в процессе обработки информации с использованием математических, лингвистических или онтологических моделей [10].

Технология эффективного анализа текста Text Mining способна выступить в роли концентратора, который делает экстракт из наиболее ключевой и значащей информации предметной области. К основным элементам Text Mining относятся суммаризация (summarization), выделение феноменов, понятий (feature extraction), кластеризация (clustering), классификация (classification), ответ на запросы (question answering), тематическое индексирование (thematic indexing) и поиск по ключевым словам (keyword searching). Также в некоторых случаях набор дополняют средства поддержки и создания таксономии (oftaxonomies) и тезаурусов (thesauri).

Результат — таксономия или визуальная карта, которая обеспечивает эффективный охват больших объемов данных. Семантические сети или анализ связей, которые определяют появление дескрипторов (ключевых фраз) в документе для обеспечения и навигации. Извлечение фактов предназначено для получения некоторых фактов из текста с целью улучшения классификации, поиска и кластеризации.

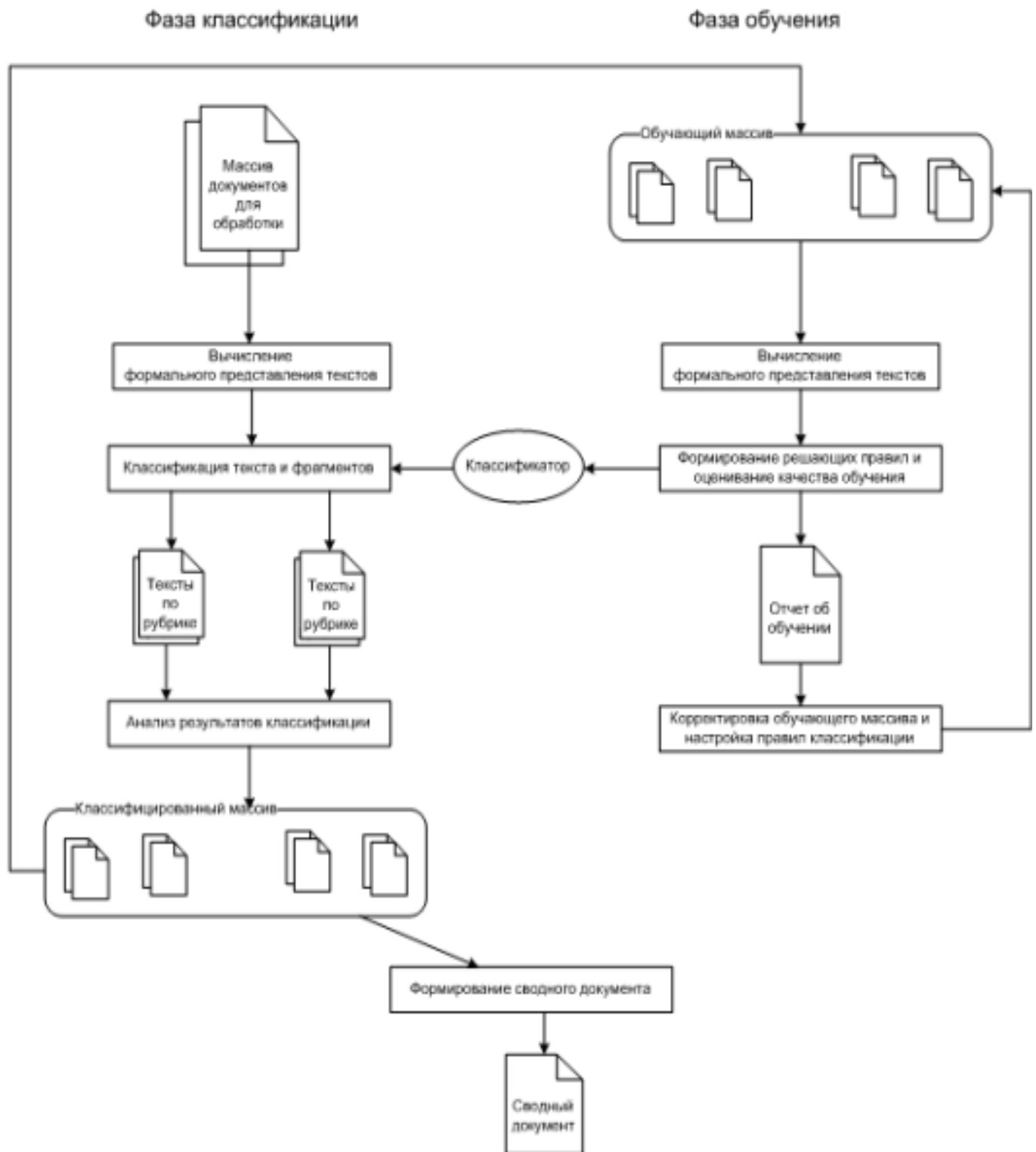


Рис. 1. Алгоритм комплексной классификации тестов [6, с.83]

```

In [1]: runfile('C:/Users/ПК/Bag/properties.py', wdir='C:/Users/ПК/Bag')
-----Исходный текст-----
До поезда одна минута! Одна минута и прощай! В кармане дребезжит валюта. Отдай
проводнику за чай! О чем же о таком ты бредишь, Отчаянно топча перрон? Куда ж
ты от себя уедешь, Когда оно со всех сторон? И запах дыма и мазута Отчаянно
объял весь свет, И волны полыхают жутко, И от судеб защиты нет!
-----Частотный анализ-----
Общее количество словоформ в тексте 57
Из них уникальных 46
-----Словарь уникальных словоформ-----
Counter({'и': 5, 'одна': 2, 'минута': 2, 'о': 2, 'ты': 2, 'отчаянно': 2, '': 2,
'от': 2, 'до': 1, 'поезда': 1, 'прощай': 1, 'в': 1, 'кармане': 1, 'дребезжит':
1, 'валюта': 1, 'отдай': 1, 'проводнику': 1, 'за': 1, 'чай': 1, 'чем': 1, 'же':
1, 'таком': 1, 'бредишь': 1, 'топча': 1, 'перрон': 1, 'куда': 1, 'ж': 1, 'себя':
1, 'уедешь': 1, 'когда': 1, 'оно': 1, 'со': 1, 'всех': 1, 'сторон': 1, 'запах':
1, 'дыма': 1, 'мазута': 1, 'объял': 1, 'весь': 1, 'свет': 1, 'волны': 1,
'полыхают': 1, 'жутко': 1, 'судеб': 1, 'защиты': 1, 'нет': 1})

```

Рис. 2. Пример реализации частотного анализа

Вторая задача — кластеризация — выделение компактных подгрупп объектов с близкими свойствами. Система должна самостоятельно найти признаки и разделить объекты по подгруппам. Она, как правило, предшествует задаче классификации, поскольку позволяет определить группы объектов.

Можно назвать еще несколько задач технологии Text Mining, например, прогнозирование, которое состоит в том, чтобы предсказать по значениям одних признаков объекта значения остальных.

Еще одна задача — нахождение исключений, то есть поиск объектов, которые своими характеристиками сильно выделяются из общей массы. Для этого сначала выясняются средние параметры объектов, а потом исследуются те объекты, параметры которых наиболее сильно отличаются от средних значений. Как известно, поиск исключений широко применяется, например, в работе спецслужб. Подобный анализ часто проводится после классификации, для того чтобы выяснить, насколько последняя была точна.

Несколько отдельно от задачи кластеризации стоит задача поиска связанных признаков (полей, понятий) отдельных документов. От предсказания эта задача отличается тем, что заранее не известно, по каким именно признакам реализуется взаимосвязь; цель именно в том и состоит, чтобы найти связи признаков. Эта задача сходна с кластеризацией, но не по множеству документов, а по множеству присущих им признаков.

Общий алгоритм комплексной классификации текстов представлен на рисунке 1.

Технология Text Miner позволяет определять, насколько правдив тот или иной текстовый документ. Обнаружение лжи в документах производится путем анализа текста и выявления изменений стиля письма, которые могут возникать при попытке исказить или скрыть информацию [8].

Дистрибутивная семантика — модели построения метрик подобия текстов — существует несколько программных реализаций: Word2Vec, GloVe, AdaGram, Text2Vec, Seq2Vec и другие [9]. Анализ взаимозаменяемости слов (модель Skipgram), анализ ассоциированных слов (модель BagOfWords).

Технология мешка (или облака) слов заключается в следующем. Из текста исключаются знаки препинания и стоп-символы (предлоги, союзы, междометия). Текст переводится в какой-то один регистр, как правило, нижний. Все слова исходного текста приводятся к своей основной форме — для существительных — именительный падеж единственного числа в мужском роде [10].

Научная новизна заключается в том, что предлагается актуальное программное решение на языке Python на базе библиотеки collections, которое строит семантическое ядро текста в виде словаря уникальных слов и проводит его частотный анализ (рисунок 2).

Язык Python обладает большим спектром инструментов, делающим его наиболее перспективной платформой для анализа текстов. К ним, в частности, относятся библиотека NLTK для анализа текста, библиотека wordcloud

для построения облака слов и модуль PyEnchant, который не только позволяет проверять правописание слов, но и предлагает варианты исправления ошибок, библиотека matplotlib для визуализации результатов.

#### ЛИТЕРАТУРА

1. Агеев М.С. и др. Автоматическая рубрикация текстов: методы и проблемы // Учен. зап. Казан. ун-та. Сер. Физ.-матем. науки. — 2008. — № 4 [Электронный ресурс] Режим доступа: <https://cyberleninka.ru/article/n/avtomaticheskaya-rubrikatsiya-tekstov-metody-i-problemy>
2. Амиева А.М. и др. Основные методики исследования структуры текста [Электронный ресурс] Режим доступа: [https://elar.urfu.ru/bitstream/10995/31691/1/conf\\_rtf\\_2015\\_28.pdf](https://elar.urfu.ru/bitstream/10995/31691/1/conf_rtf_2015_28.pdf)
3. Андреев А.М. и др. Автоматическая классификация текстовых документов с использованием нейросетевых алгоритмов и семантического анализа [Электронный ресурс] Режим доступа: [http://www.ixlab.ru/pub/docs/RCDL\\_2003.pdf](http://www.ixlab.ru/pub/docs/RCDL_2003.pdf)
4. Батура Т.В. Методы автоматической классификации текстов. // Программные продукты и системы. — 2017. — № 1 [Электронный ресурс] Режим доступа: <https://cyberleninka.ru/article/n/metody-avtomaticheskoy-klassifikatsii-tekstov>.
5. Васильев В.Г. Автоматическое выделение значимых фрагментов в текстах. // Обозрение прикладной и промышленной математики. Выпуск 3, том 14. — М., 2007. — 518 с.
6. Васильев В.Г. Комплексная технология автоматической классификации текстов. // Труды международной конференции «Диалог 2008» [Электронный ресурс] Режим доступа: <https://docplayer.com/55296280-Kompleksnaya-tehnologiya-avtomaticheskoy-klassifikatsii-tekstov-complex-technology-of-automatic-text-classification.html>
7. Добров А.В. Автоматическая рубрикация текстов средствами комплексного лингвистического анализа [Электронный ресурс] Режим доступа: <http://aiire.org/pubs/2012...>
8. Кузнецов И.В. Введение в анализ текстовой информации с помощью Python и методов машинного обучения [Электронный ресурс] Режим доступа: <https://habr.com/ru/post/205360/>
9. Ландэ Д. Глубинный анализ текстов. Технология эффективного анализа текстовых данных [Электронный ресурс] Режим доступа: <http://www.visti.net/~dwl/art/dz/>
10. Сидорова Е.А. Подход к моделированию процесса извлечения информации из текста на основе онтологии. // Онтология проектирования. — 2018. — № 1 (27) [Электронный ресурс] Режим доступа: <https://cyberleninka.ru/article/n/podhod-k-modelirovaniyu-protsessa-izvlecheniya-informatsii-iz-tekstana-osnove-ontologii>.
11. Частотный анализ русского текста и облако слов на Python [Электронный ресурс] Режим доступа: <https://habr.com/ru/post/517410/?amp&amp>