

УДК: 004.912

DOI: <https://doi.org/10.47813/2782-5280-2023-2-2-0301-0313>

EDN: [JKPDAC](https://oajiem.com/)



## Об одном подходе к извлечению именованных сущностей из неструктурированных текстов

А. А. Ворошилова<sup>1,2</sup>, С. Ю. Пискорская<sup>3</sup>

<sup>1</sup>Красноярский краевой Дом науки и техники РосСНИО, Красноярск, Россия

<sup>2</sup>Сибирский федеральный университет, Красноярск, Россия

<sup>3</sup>СибГУ им. М.Ф. Решетнева, Красноярск, Россия

**Аннотация.** В статье рассматривается один из возможных подходов к извлечению именованных сущностей из неструктурированных текстов. Отмечается сложность и трудоемкость наиболее распространенных методов решения данной задачи, базирующихся на использовании создаваемых вручную конечных автоматов. Возникает ряд сложностей при реализации данного подхода при обработке мультилингвистических текстов, так как для каждого нового языка и для каждого нового класса сущностей требуется вмешательство человека для создания вручную нового набора шаблонов для работы с новыми языками и новыми классами. Предлагаемый подход предполагает использование принципов машинного обучения. Дана постановка задачи и описана используемая модель марковской цепи при распознавании именованных сущностей. На основе данной модели для выделения именованных объектов ставится задача нахождения наиболее вероятной последовательности состояний, генерирующих последовательность лексем. В статье описан лексический материал, включающий состав признаков и их описания, представлена методика декодирования и оценка параметров модели. В данной работе для решения задачи используется алгоритм Витерби, который предназначен для нахождения последовательности состояний, для которых вероятность порождения наблюдаемой цепочки символов максимальна. В качестве экспериментальных результатов представлены характеристики точности распознавания типов лексем при различных размерах обучающей выборки и диаграмма количества ошибок по классам лексем.

**Ключевые слова:** обработка информации, неструктурированный текст, именованная сущность, лексема, скрытая марковская цепь.

**Для цитирования:** Ворошилова, А. А., & Пискорская, С. Ю. (2023). Об одном подходе к извлечению именованных сущностей из неструктурированных текстов. Информатика. Экономика. Управление - Informatics. Economics. Management, 2(2), 0301–0313. <https://doi.org/10.47813/2782-5280-2023-2-2-0301-0313>

## To one approach to extracting named entities from unstructured texts

A. A. Voroshilova<sup>1,2</sup>, S. Yu. Piskorskaya<sup>3</sup>

<sup>1</sup>*Krasnoyarsk Science and Technology City Hall, Krasnoyarsk, Russia*

<sup>2</sup>*Siberian Federal University, Krasnoyarsk, Russia*

<sup>3</sup>*Reshetnev Siberian State University of Science and Technologies, Krasnoyarsk, Russia*

**Abstract.** The article considers one of the possible approaches to the extraction of named entities from unstructured texts. The complexity and laboriousness of the most common methods for solving this problem, based on the use of manually created finite automata, are noted. There are a number of difficulties in implementing this approach when processing multilingual texts, since for each new language and for each new class of entities, human intervention is required to manually create a new set of templates for working with new languages and new classes. The proposed approach involves the use of machine learning principles. The statement of the problem is given and the model of the Markov chain used in the recognition of named entities is described. On the basis of this model for the selection of named objects, the task is to find the most probable sequence of states that generate a sequence of tokens. The article describes the lexical material, including the composition of features and their descriptions, presents the decoding technique and estimation of the model parameters. In this paper, to solve the problem, the Viterbi algorithm is used, which is designed to find a sequence of states for which the probability of generating the observed chain of symbols is maximum. As experimental results, the characteristics of the accuracy of recognition of types of lexemes for different sizes of the training sample and a diagram of the number of errors by classes of lexemes are presented.

**Keywords:** information processing, unstructured text, named entity, lexeme, hidden Markov chain.

**For citation:** Voroshilova, A. A., & Piskorskaya, S. Y. (2023). To one approach to extracting named entities from unstructured texts. *Informatics. Economics. Management*, 2(2), 0301–0313. <https://doi.org/10.47813/2782-5280-2023-2-2-0301-0313>

---

### ВВЕДЕНИЕ

Извлечение информации – это задача из области обработки естественно-языковых текстовых массивов, которая включает автоматическое извлечение predetermined типов информации из текста [1-4]. Примером задачи извлечения информации может служить задача получения сведений об организации из массивов информации, представленной в текстовом виде [5]. Входными данными системы извлечения информации является неструктурированный или слабоструктурированный текст на естественном языке; на выходе - заполненные структуры данных, позволяющие проводить дальнейшую автоматическую или ручную обработку информации.

В качестве частного случая данной задачи можно рассмотреть задачу извлечения именованных сущностей (примером может служить выявление в неструктурированном

или слабоструктурированном тексте всех вхождений упоминаний о различных организациях, персонах, географических названий и т.д.) [6, 7]. Ряд авторов, например, в [8-11] используют внутриязыковые ассоциативные поля в мультилингвистической адаптивно-обучающей технологии, а также исследуют системы поиска, анализа и обработки мультилингвистических текстов, интегрированные с информационно-поисковыми системами [12].

Распространенные подходы к решению задачи извлечения именованных сущностей из неструктурированных текстов основываются на использовании создаваемых вручную конечных автоматов (patterns) [13]. Однако для каждого нового языка (в рамках мультилингвистической технологии [11]) и для каждого нового класса сущностей требовалось вмешательство человека для создания вручную нового набора шаблонов для работы с новыми языками и новыми классами. Предлагаемый подход предполагает использование принципов машинного обучения.

Процесс извлечения информации состоит из двух этапов: первый – поиск возможных *кандидатур* лексем, представляющих интерес, второй – определение типа каждой из *кандидатур*. Алгоритм распознавания типа лексемы на выходе должен выдавать единственный и однозначный тип для каждой лексемы в тексте.

Задача автоматического распознавания некоторых типов лексем достаточно тривиальна, в то время как для ряда лексем могут возникнуть неоднозначные толкования. Например, автоматическое распознавание лексем адресов электронной почты и дат может осуществляться при помощи механизмов стандартных регулярных выражений. Однако использование регулярных выражений для некоторых типов лексем, например имен, довольно затруднительно. Например, для лексемы “Владимир” может возникнуть неоднозначность - к какому типу отнести эту лексему: имя человека или название города? В естественных языках, как правило, не существует каких-либо конкретных ограничений на правила формирования названий именованных объектов.

## ПОСТАНОВКА ЗАДАЧИ

Будем рассматривать входной текст как выход некоторой порождающей системы, которая порождает текст, состоящий из лексем, причем у каждой лексемы имеется набор характеристик, включающих семантический тип лексемы. Типам лексем в рассматриваемой модели соответствуют такие типы именованных сущностей как: географические названия, названия организаций, персоны и т. д. Можно ввести

следующую аналогию: при прохождении текстовой информации через некоторый зашумленный канал информация о семантических типах утратилась. Таким образом, задача состоит в восстановлении этой информации.

Соотнесем набор состояний  $S = \{S_1, \dots, S_N\}$  с введенным набором семантических типов, то есть каждое состояние  $S_i$  будет соответствовать некоторому семантическому типу лексемы. Так как входной текст представляет собой последовательность лексем, которую обозначим как  $O = O_1 O_2, \dots, O_T$ , причем у каждой лексемы имеется свой семантический тип, мы можем представить систему, которая в каждый момент времени  $t=1, \dots, T$  находится в одном из состояний  $S_1 \dots S_N$ .

В данной работе в качестве модели системы используются скрытые марковские цепи: переход системы из одного состояния в другое происходит в моменты времени  $t = 1, 2, \dots$  в соответствии с вероятностями перехода, соотнесенными с состояниями. Вероятность перехода из состояния  $S_i$  в  $S_j$  определяется матрицей  $A$ , состоящей из следующих элементов:

$$a_{ij} = P[q_t = S_j | q_{t-1} = S_i], 1 \leq i, j \leq N.$$

Для элементов матрицы  $A$  выполняются стандартные ограничения:  $a_{ij} \geq 0$ , и условие нормировки вероятностей  $\sum_{i=1}^N a_{ij} = 1$ .

Полностью модель скрытой марковской цепи определяется следующими элементами:

1. Набор из  $N$  состояний  $S = \{S_1, \dots, S_N\}$ .
2. Алфавит, состоящий из  $M$  символов  $V = \{v_1, \dots, v_M\}$ .
3. Вероятности перехода из состояния  $S_i$  в  $S_j$  определяемые элементами матрицы  $A$ , где  $a_{ij} = P[q_t = S_j | q_{t-1} = S_i], 1 \leq i, j \leq N$ .
4. Вероятности порождения символа алфавита  $v_k$ , если система находится в состоянии  $S_j$ , определяются элементами матрицы  $B$ , где  $b_i(k) = P[v_k | q_t = S_i], 1 \leq j \leq N, 1 \leq k \leq M$ .
5. Вероятности начальных состояний  $\pi_i = P[q_1 = S_i], 1 \leq i \leq N$ .

Данный вид цепей называют скрытыми, так как наблюдается только последовательность порожденных символов  $O = O_1 O_2, \dots, O_T$ , где  $O_i \in V$ , при этом, последовательный набор состояний, породивший данную последовательность, остается скрытым.

Для краткости, совокупность параметров модели обозначим  $\lambda = (A, B, \pi)$ .

На основе этой модели для скрытых марковских цепей возможно оценить:

- вероятность, с которой модель  $\lambda$  порождает последовательность  $O = O_1O_2, \dots, O_T$ , где  $O_i \in V$ ;
- наиболее вероятную последовательность состояний, генерирующую последовательность  $O$ ;
- параметры  $\lambda = (A, B, \pi)$ , которые максимизируют вероятность порождения последовательности  $O$ :  $\max P(O|\lambda)$ .

Используемая модель марковской цепи при распознавании именованных сущностей представлена на рисунке 1.

Здесь для каждого типа лексемы предусмотрено соответствующее состояние в марковской модели.

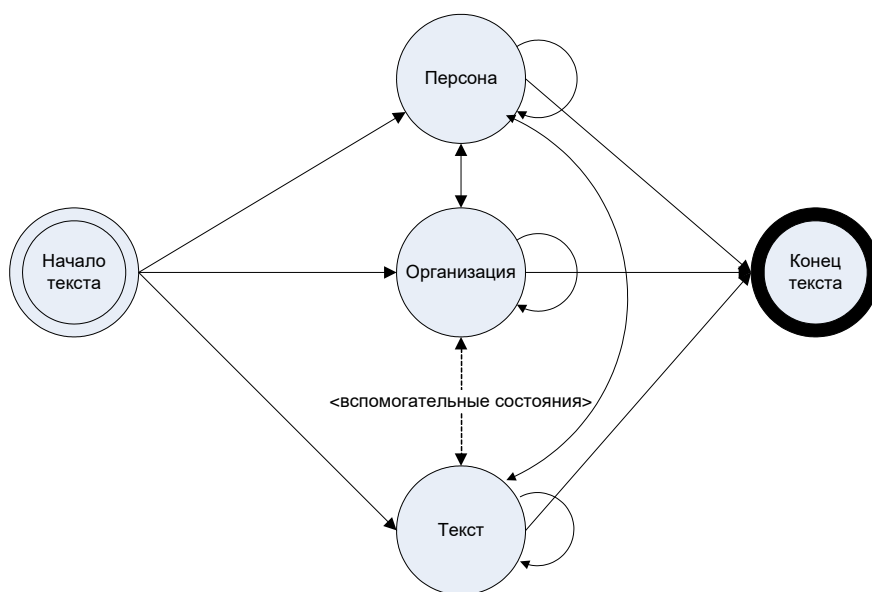


Рисунок 1. Модель марковской цепи при распознавании именованных сущностей.

Figure 1. Markov chain model for named entity recognition.

В данном случае, для выделения именованных объектов, ставится задача нахождения наиболее вероятной последовательности состояний, генерирующих последовательность  $O$ .

## ЛЕКСИЧЕСКИЙ МАТЕРИАЛ

При распознавании лексем, алфавит  $V$  включает в себя лексемы, полученные на этапе обучения, а также ряд признаков лексем, составленных из набора базовых признаков. Состав признаков и их описание представлено в таблице 1.

Таблица 1. Состав признаков и их описание.

Table 1. Composition of features and their description.

Признак	Описание
FIRST_CAP	Первая буква заглавная.
ALL_CAP	Все буквы заглавные.
IN_QUOTES	В кавычках.
NON_VOCAB	Неизвестное слово русского языка.
LETTER_WITH_DOT	Заглавная буква с точкой.
LETTER_THEN_DOT_THEN_CAP	Две заглавные буквы, разделенные точкой.
PREFIX1	Словарный префикс.
SUFFIX1	Словарный суффикс.
OTHER	Другое.

Исходя из данного набора элементарных признаков видно, что одной лексеме может соответствовать несколько элементарных признаков. В таком случае происходит их объединение в новый признак, который добавляется в словарь.

Набору состояний соответствуют семантические типы лексем вида:

- ORG – название организации (АО “ОКБ”);
- PERSON – упоминание о человеке (Владимир Баранов);
- GEOGRAPHIC – географическое название (Усть-Илимский район);
- PLAIN\_TEXT – обычный текст.

При обучении для разметки лексем используется расширяемый язык разметки XML.

Пример разметки текста выглядит следующим образом.

*Как отмечает начальник отдела продаж* <ORGPREFIX>ИФ</ORGPREFIX>  
<ORG>ОЛМА</ORG><FIRSTNAME>Владимир</FIRSTNAME><PERSON>Баранов</

*PERSON*>, подъем курса акций <ORG>ОКБ</ORG> проходит на фоне известий, что компания выиграла конкурс на оказание услуг.

Элементы матрицы  $A$  определяет вероятность того, что текущая лексема принадлежит классу  $j$ , учитывая, что предыдущая лексема принадлежала классу  $i$ .

Элементы матрицы  $B$  определяют вероятность того, что лексема с определенным набором элементарных признаков принадлежит определенному классу.

## МЕТОДИКА ДЕКОДИРОВАНИЯ

Предположим, что имеется некоторая последовательность  $O = O_1 O_2, \dots, O_T$ , где  $O_i \in V$ , также известны элементы матриц  $A, B, \pi$ .

Задача нахождения наиболее вероятного набора состояний  $Q = Q_1 Q_2, \dots, Q_T, Q_i \in S$  называется декодированием.

В данном случае, для решения этой задачи используется алгоритм Витерби [15], который предназначен для нахождения последовательности состояний, для которых вероятность порождения наблюдаемой цепочки символов максимальна.

Максимальную вероятность того, что на  $t$ -м моменте времени символ  $o_t$  был порожден состоянием  $S_j$ , обозначим как:

$$\delta_t(j) = \max_{q_1, q_2, \dots, q_{t-1}} P[q_1 q_2 \dots q_{t-1}, o_1 o_2 \dots o_{t-1}, q_t = S_j, o_t]. \quad (1)$$

В рекуррентном виде данное соотношение может быть записано как:

$$\delta_{t+1}(j) = [\max_{1 \leq i \leq N} \delta_t(i) \cdot a_{ij}] \cdot b_j(o_{t+1}), \quad (2)$$

$$1 \leq j \leq N, 1 \leq t \leq T - 1.$$

В начальный момент времени

$$\delta_1(j) = \pi_j \cdot b_j(o_1). \quad (3)$$

Таким образом, алгоритм начинается с вычисления  $\delta_1(j)$  для  $1 \leq j \leq N$ , затем, используя рекуррентную формулу (2) вычисляются значения последующих  $\delta_t(j)$  до  $t=T$ , для получения оптимальной последовательности состояний. Последнее состояние  $j$  определяется как:

$$j = \arg \max_{1 \leq j \leq N} \delta_T(j). \quad (4)$$

Критерий оптимальности для алгоритма Витерби может быть записан следующим образом:

$$Q^* = \max_{q_1, q_2, \dots, q_T} \prod_{t=1}^T P[q_t = S_j | q_{t-1} = S_i] \cdot P[o_t = v_k | q_t = S_j], \quad (5)$$

где  $Q^*$  - оптимальный набор состояний, который максимизирует вероятность порождения наблюдаемой последовательности символов алфавита.

### ОЦЕНКА ПАРАМЕТРОВ МОДЕЛИ

Элементы матриц  $A$  и  $B$  – вероятности переходов состояний и вероятности порождения символов алфавита  $V$  возможно оценить при наличии обучающей выборки. Обучающая выборка представляет собою размеченный текст, где при помощи тэгов XML разметки указан тип лексемы.

На основании обучающей выборки элементы матрицы  $A$  можно оценить следующим образом [14]:

$$a_{ij} = \frac{c(S_i \rightarrow S_j)}{c(S_i)}. \quad (6)$$

Для элементов матрицы  $B$  имеем:

$$b_j(k) = \frac{c(V_k \uparrow S_j)}{c(S_j)}. \quad (7)$$

Элементы вектора  $\pi$  имеют следующий вид:

$$\pi_j = \frac{c(Start \rightarrow S_j)}{c(Start)}. \quad (8)$$

В данных формулах введены следующие обозначения:

- $c(X)$  – число появления события  $X$ ;
- $S_i \rightarrow S_j$  - переход системы из состояния  $i$  в состояние  $j$ ;
- $V_k \uparrow S_j$  - порождение символа  $V_k$ , когда система находится в состоянии  $S_j$ ;
- $Start \rightarrow S_j$  - появления в качестве первого состояния системы состояния  $j$ .

### ЭКСПЕРИМЕНТАЛЬНЫЕ РЕЗУЛЬТАТЫ

В качестве источников данных для обучения и тестирования использовались новостные ленты из агентства РосБизнесКонсалтинг <http://www.rbcdaily.ru> и <http://www.quote.ru/>.



Характеристики точности распознавания типов лексем при размере обучающей выборки 740 и 1460 лексем соответственно представлены на рисунке 2.

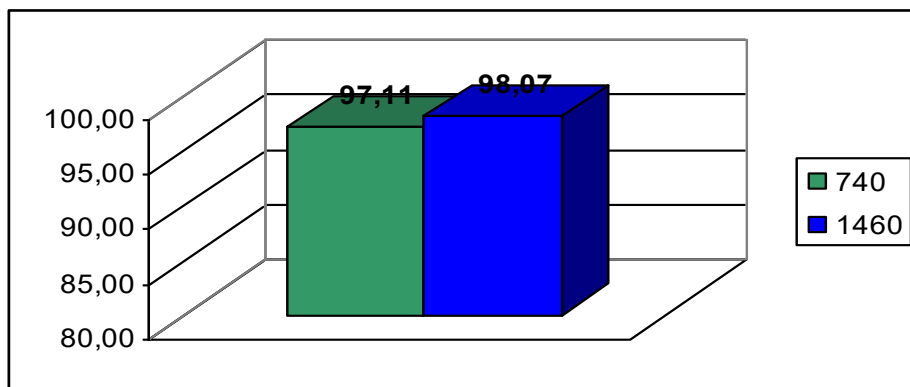


Рисунок 2. Характеристики точности распознавания типов лексем при различных размерах обучающей выборки.

Figure 2. Characteristics of the accuracy of recognition of types of lexemes for different sizes of the training sample.

Диаграмма количества ошибок для различных классов лексем представлена на рисунке 3.

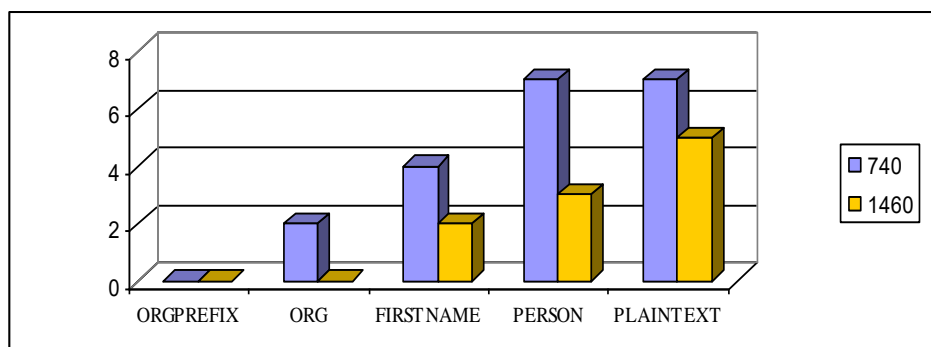


Рисунок 3. Диаграмма количества ошибок по классам лексем.

Figure 3. Diagram of the number of errors by token class.

## ЗАКЛЮЧЕНИЕ

На основе лексического материала, представленного в статье и включающего состав признаков и их описание, разработана методика декодирования и оценка параметров модели для извлечения именованных сущностей из неструктурированных

текстов. Для решения задачи использован алгоритм Витерби, который предназначен для нахождения последовательности состояний, для которых вероятность порождения наблюдаемой цепочки символов максимальна. Критерий оптимальности для алгоритма Витерби представлен в работе в виде формулы (5). Экспериментальные результаты продемонстрировали высокий уровень точности распознавания типов лексем при различных размерах обучающей выборки: 97,11% для 740 лексем и 98,07% для 1460 лексем. Также продемонстрирован достаточно низкий уровень количества ошибок по всем классам лексем (см. рисунок 3). Предложенный подход может быть расширен для мультилингвистического базиса лексем, принадлежащих разным языкам, в рамках мультилингвистической адаптивно-обучающей технологии [16].

### СПИСОК ЛИТЕРАТУРЫ

- [1] Распопин Н.А., Карасева М.В., Зеленков П.В., Каюков Е.В., Ковалев И.В. Модели и методы оптимизации сбора и обработки информации. Сибирский аэрокосмический журнал. 2012; 2(42): 69-72.
- [2] Коровиков Н.А., Гончаров М.А., Кадров М.С. Анализ методов выделения именованных сущностей из неструктурированных документов. Международный журнал прикладных наук и технологий «Integral». 2019; 3: 328-332.
- [3] Абрамов П.С. Извлечение ключевой информации из текста. Новые информационные технологии в автоматизированных системах. 2018; 21: 217-219.
- [4] Киселев С.Л., Ермаков А.Е., Плешко В.В. Поиск фактов в тексте естественного языка на основе сетевых описаний. Компьютерная лингвистика и интеллектуальные технологии: труды международной конференции Диалог'2004. М.: Наука; 2004: 180-185.
- [5] Nadeau D., Sekine S. A survey of named entity recognition and classification. *Linguisticae Investigationes*. 2007; 1(30): 3-26. <https://doi.org/10.1075/li.30.1.03nad>
- [6] Gentile A. L. et al. Cultural Knowledge for Named Entity Disambiguation: A Graph-Based Semantic Relatedness Approach. *Serdica Journal of Computing*. 2010; 4(2): 217-242. <https://doi.org/10.55630/sjc.2010.4.217-242>
- [7] Bikel D. M., Miller S., Schwartz R., Weischedel R. Nymble: a high performance learning namefinder. *Proceedings of the Fifth Conference on Applied Natural Language Processing (ANLP-97)*; 1997: 194-201. <https://doi.org/10.3115/974557.974586>
- [8] Brester C., Semekin E., Kovalev I., Zelenkov P., Sidorov M. Evolutionary feature selection for emotion recognition in multilingual speech analysis. *IEEE Congress on*

Evolutionary Computation (CEC 2015); 2015: 2406-2411.  
<https://doi.org/10.1109/CEC.2015.7257183>

[9] Ковалев И.В., Лесков О.В., Карасева М.В. Внутриязыковые ассоциативные поля в мультилингвистической адаптивно-обучающей технологии. Системы управления и информационные технологии. 2008; 3-1(33): 157-160.

[10] Зеленков П.В., Ковалев И.В., Карасева М.В., Рогов С.В. Мультилингвистическая модель распределенной системы на основе тезауруса. Сибирский аэрокосмический журнал. 2008; 1(18): 26-28.

[11] Ковалев И.В. Системная архитектура мультилингвистической адаптивно-обучающей технологии и современная структурная методология. Телекоммуникации и информатизация образования. 2002; 3: 6.

[12] Ковалев И.В., Полянский К.В., Зеленков П.В., Брезицкая В.В., Сидорова Г.А. Система поиска, анализа и обработки мультилингвистических текстов, интегрированная с информационно-поисковыми системами. Сибирский аэрокосмический журнал. 2013; 1(47): 48-52.

[13] Appelt D., Hobbs J., Bear J., Israel D., Tyson M. FASTUS: A finitestate processor for information extraction from real-world text. Proceedings of the 13th International Joint Conference on Artificial Intelligence (IJCAI-93). Chambéry, France; 1993: 1172–1178.  
<https://doi.org/10.3115/1075671.1075701>

[14] Rabiner L.R. A tutorial on hidden Markov models and selected applications in speech recognition. Proceedings of the IEEE. 1989; 77(2): 257-286. <https://doi.org/10.1109/5.18626>

[15] Wen Y. Text Mining Using HMM and PPM. Master's thesis. Department of Computer Science, University of Waikato. 2001.

[16] Ковалев И.В., Карасева М.В., Суздалева Е.А. Системные аспекты организации и применения мультилингвистической адаптивно-обучающей технологии. Образовательные технологии и общество. 2002; 5(2): 198-212.

## REFERENCES

[1] Raspopin N.A., Karaseva M.V., Zelenkov P.V., Kayukov E.V., Kovalev I.V. Modeli i metody optimizacii sbora i obrabotki informacii. Sibirskij aerokosmicheskij zhurnal. 2012; 2(42): 69-72. (in Russian)

[2] Korovikov N.A., Goncharov M.A., Kadrov M.S. Analiz metodov vydeleniya imenovannyh sushchnostej iz nestrukturirovannyh dokumentov. Mezhdunarodnyj zhurnal

- prikladnyh nauk i tekhnologij «Integral». 2019; 3: 328-332. (in Russian)
- [3] Abramov P.S. Izvlechenie klyuchevoj informacii iz teksta. Novye informacionnye tekhnologii v avtomatizirovannyh sistemah. 2018; 21: 217-219. (in Russian)
- [4] Kiselev S.L., Ermakov A.E., Pleshko V.V. Poisk faktov v tekste estestvennogo yazyka na osnove setevykh opisaniy. Komp'yuternaya lingvistika i intellektual'nye tekhnologii: trudy mezhdunarodnoj konferencii Dialog'2004. M.: Nauka; 2004: 180-185. (in Russian)
- [5] Nadeau D., Sekine S. A survey of named entity recognition and classification. *Linguisticae Investigationes*. 2007; 1(30): 3-26. <https://doi.org/10.1075/li.30.1.03nad>
- [6] Gentile A. L. et al. Cultural Knowledge for Named Entity Disambiguation: A Graph-Based Semantic Relatedness Approach. *Serdica Journal of Computing*. 2010; 4(2): 217-242. <https://doi.org/10.55630/sjc.2010.4.217-242>
- [7] Bikel D. M., Miller S., Schwartz R., Weischedel R. Nymble: a high performance learning namefinder. *Proceedings of the Fifth Conference on Applied Natural Language Processing (ANLP-97)*; 1997: 194-201. <https://doi.org/10.3115/974557.974586>
- [8] Brester C., Semenkin E., Kovalev I., Zelenkov P., Sidorov M. Evolutionary feature selection for emotion recognition in multilingual speech analysis. *IEEE Congress on Evolutionary Computation (CEC 2015)*; 2015: 2406-2411. <https://doi.org/10.1109/CEC.2015.7257183>
- [9] Kovalev I.V., Leskov O.V., Karaseva M.V. Vnutriyazykovye asociativnye polya v mul'tilingvisticheskoj adaptivno-obuchayushchej tekhnologii. *Sistemy upravleniya i informacionnye tekhnologii*. 2008; 3-1(33): 157-160. (in Russian)
- [10] Zelenkov P.V., Kovalev I.V., Karaseva M.V., Rogov S.V. Mul'tilingvisticheskaya model' raspredelennoj sistemy na osnove tezaurusa. *Sibirskij aerokosmicheskij zhurnal*. 2008; 1(18): 26-28. (in Russian)
- [11] Kovalev I.V. Sistemnaya arhitektura mul'tilingvisticheskoj adaptivno-obuchayushchej tekhnologii i sovremennaya strukturnaya metodologiya. *Telekommunikacii i informatizaciya obrazovaniya*. 2002; 3: 6. (in Russian)
- [12] Kovalev I.V., Polyanskij K.V., Zelenkov P.V., Brezickaya V.V., Sidorova G.A. Sistema poiska, analiza i obrabotki mul'tilingvisticheskikh tekstov, integrirovannaya s informacionno-poiskovymi sistemami. *Sibirskij aerokosmicheskij zhurnal*. 2013; 1(47): 48-52. (in Russian)
- [13] Appelt D., Hobbs J., Bear J., Israel D., Tyson M. FASTUS: A finitestate processor for information extraction from real-world text. *Proceedings of the 13th International Joint Conference on Artificial Intelligence (IJCAI-93)*. Chambéry, France; 1993: 1172–1178.

<https://doi.org/10.3115/1075671.1075701>

[14] Rabiner L.R. A tutorial on hidden Markov models and selected applications in speech recognition. Proceedings of the IEEE. 1989; 77(2): 257-286. <https://doi.org/10.1109/5.18626>

[15] Wen Y. Text Mining Using HMM and PPM. Master's thesis. Department of Computer Science, University of Waikato. 2001.

[16] Kovalev I.V., Karaseva M.V., Suzdaleva E.A. Sistemnye aspekty organizatsii i primeneniya mul'tilingvisticheskoy adaptivno-obuchayushchej tekhnologii. Obrazovatel'nye tekhnologii i obshchestvo. 2002; 5(2): 198-212. (in Russian)

#### ИНФОРМАЦИЯ ОБ АВТОРАХ / INFORMATION ABOUT THE AUTHORS

**Ворошилова Анна Анатольевна**, кандидат философских наук, доцент, кафедра Информатики, Сибирский федеральный университет, Красноярск, Россия  
e-mail: [krasnio@bk.ru](mailto:krasnio@bk.ru)  
ORCID: <https://orcid.org/0000-0002-4556-813X>

**Anna Voroshilova**, Candidate of Philosophical Sciences, Associate Professor, Department of Informatics, Siberian Federal University, Krasnoyarsk, Russia

**Пискорская Светлана Юрьевна**, доктор философских наук, профессор, директор Института социального инжиниринга, СибГУ имени академика М.Ф. Решетнева, Красноярск, Россия  
e-mail: [piskorskaya1@rambler.ru](mailto:piskorskaya1@rambler.ru)  
ORCID: <https://orcid.org/0000-0002-5589-801X>

**Svetlana Piskorskaya**, Doctor of Philosophy, Professor, Director of the Institute of Social Engineering, Reshetnev Siberian State University of Science and Technologies, Krasnoyarsk, Russia

*Статья поступила в редакцию 25.06.2023; одобрена после рецензирования 17.07.2023; принята к публикации 18.07.2023.*

*The article was submitted 25.06.2023; approved after reviewing 17.07.2023; accepted for publication 18.07.2023.*