

Министерство науки и высшего образования
Российской Федерации

УНИВЕРСИТЕТ ИТМО

Некоммерческое партнерство ПРИОР Северо-Запад

КОМПЬЮТЕРНАЯ ЛИНГВИСТИКА И ВЫЧИСЛИТЕЛЬНЫЕ ОНТОЛОГИИ

Выпуск 4

**Труды XIII Международной
объединенной научной конференции
«Интернет и современное общество»,
IMS-2020, Санкт-Петербург,
17–20 июня 2020 г.**

Сборник научных трудов

 **УНИВЕРСИТЕТ ИТМО**

Санкт-Петербург

2020

УДК 81`33
ББК 81.1
К63

Рецензенты:

канд. филол. наук Е.Л. Алексеева, канд. филол. наук О.А. Митрофанова

Редколлегия:

А.В. Добров, В.П. Захаров (председатель), М.В. Хохлова, А.В. Чижик

Ответственный редактор издания:

канд. филол. наук В.П. Захаров

К63 **Компьютерная лингвистика и вычислительные онтологии.** Выпуск 4 (Труды XXIII Международной объединенной научной конференции «Интернет и современное общество», IMS-2020, Санкт-Петербург, 17 – 20 июня 2020 г. Сборник научных статей). — СПб.: Университет ИТМО, 2020. — 132 с.

ISSN 2541-9781

ISBN 978 7 77 0 31 3

В сборник включены тексты статей, представленные на XXIII Международной объединенной научной конференции «Интернет и современное общество» (Internet and Modern Society - IMS). Работы прошли рецензирование и отобраны в результате конкурсной процедуры. Сборник снабжен авторским указателем.

Издание адресовано научным работникам, преподавателям, аспирантам и магистрантам, изучающим междисциплинарные проблемы влияния информационно-коммуникационных технологий на трансформацию социальных и политических отношений в современном обществе.

Информация о конференции «Интернет и современное общество» представлена на сайте объединенной конференции (<http://ims.ifmo.ru>). Все статьи и тезисы докладов конференции IMS публикуются в открытом доступе (лицензия Creative Commons — CC-BY 3.0 Unported). Сборники научных статей, издаваемые в рамках конференции IMS с 2011 года, размещаются в Научной электронной библиотеке (<http://elibrary.ru/>) и РИНЦ.

УДК 800(075.3)

ББК 81.1

 УНИВЕРСИТЕТ ИТМО

Университет ИТМО — ведущий вуз России в области информационных и фотонных технологий, один из немногих российских вузов, получивших в 2009 году статус национального исследовательского университета. С 2013 года Университет ИТМО — участник программы повышения конкурентоспособности российских университетов среди ведущих мировых научно-образовательных центров, известной как проект «5 в 100». Цель Университета ИТМО — становление исследовательского университета мирового уровня, предпринимательского по типу, ориентированного на интернационализацию всех направлений деятельности.

© Университет ИТМО, 2020

© Авторы, 2020

XXIII Международная объединенная научная конференция «Интернет и современное общество» (IMS-2020)

Санкт-Петербург, 17–20 июня 2020 г.

<http://ims.ifmo.ru>

Конференция «Интернет и современное общество» (Internet and Modern Society – IMS) проводится в Санкт-Петербурге ежегодно с 1998 года. С 2014 г. конференция проводится как международное научное мероприятие. С 2016 г. конференция проводится в рамках Недели технологий информационного общества в Санкт-Петербурге.

Организаторы конференции IMS-2020:

- Университет ИТМО
- Некоммерческое партнерство ПРИОР Северо-Запад

Конференция является «объединенной», т.к. научная программа конференции объединяет серию специализированных международных и российских научных конференций, симпозиумов, семинаров, круглых столов и других мероприятий, посвященных специальным вопросам развития технологий информационного общества. Отдельные специализированные и проблемно-ориентированные мероприятия проводятся в сотрудничестве с партнерскими организациями.

Основные мероприятия Недели технологий информационного общества Санкт-Петербурге в 2020 г.:

- Симпозиум молодых ученых «**Цифровые трансформации: перспективные социально-экономические и гуманитарные исследования**»: 17 июня (рабочий язык – русский). Организаторы: Университет ИТМО и Северо-Западный институт управления РАНХиГС - http://ims.ifmo.ru/ru/pages/24/molodezhnyy_simpozium.htm.
- Мероприятия объединенной конференции «**Интернет и современное общество**» (открытие, секции, круглые столы): 17 – 20 июня. Сайт конференции: <http://ims.ifmo.ru>. В программе объединенной конференции в 2020 году были проведены специализированные международные семинары, включающие две сессии – на русском и английском языках:
 - 19 июня: Международный семинар «**Киберпсихология**» (Internet Psychology – IntPsy-2020);
 - 19 – 20 июня: Международный семинар «**Электронное управление**» (E-Governance-2020);
 - 20 июня: Международный семинар «**Компьютерная лингвистика**» (Computational Linguistics – CompLing-2020).
- Международная конференция «**Digital Transformation & Global Society**» (DTGS-2020): 17 – 20 июня (рабочий язык – английский). Организаторы: Университет ИТМО и НИУ ВШЭ, Санкт-Петербург. Сайт конференции: <http://dtgs-conference.org>.

В 2020 году в связи с ограничениями, вызванными коронавирусной инфекцией, мероприятия конференций Недели технологий информационного общества проводились в онлайн-формате (ZOOM) с предварительным размещением всех презентаций и текстов статей в системе, доступной через личные кабинеты участников конференций и предоставляющей возможности сформулировать вопросы докладчикам.

Отбор докладов на конференции и текстов для публикации производится по результатам слепого рецензирования членами программного комитета с использованием международной системы сопровождения научных конференций EasyChair.org.

По результатам объединенной конференции IMS-2020 издаются три сборника научных трудов (сериальные издания), сборник тезисов на русском языке и сборник статей на английском языке:

- Информационное общество: образование, наука, культура и технологии будущего (ISSN 2587-8557), вып. 4;
- Государство и граждане в электронной среде (ISSN 2541-979X), вып. 4;
- Компьютерная лингвистика и вычислительные онтологии (ISSN 2541-9781), вып. 4;
- Интернет и современное общество: сборник тезисов докладов IMS-2020.

Статьи, представленные для докладов на английском языке и прошедшие рецензирование международным программным комитетом, публикуются в сборнике «IMS2020 Proceedings» (издание CEUR-WS: Free Open-Access Proceedings for Scientific Conferences and Workshops) на английском языке с индексацией в DBLP и Scopus.

ПРОГРАММНЫЙ КОМИТЕТ КОНФЕРЕНЦИИ

Председатель Программного комитета:

Васильев В.Н., докт. техн. наук, чл.-корр. РАН, ректор Университета ИТМО

Заместители председателя Программного комитета:

Борисов Н.В., докт. физ.-мат. наук, заведующий кафедрой информационных систем в искусстве и гуманитарных науках СПбГУ, председатель Оргкомитета конференции

Чугунов А.В., канд. политич. наук, директор Центра технологий электронного правительства ИДУ Университета ИТМО, генеральный директор НП ПРИОР Северо-Запад, ученый секретарь конференции

Члены Программного комитета:

Алехин А.Н., докт. мед. наук, Российский государственный педагогический университет имени А.И. Герцена

Антопольский А.Б., докт. техн. наук, академик РАЕН, ИНИОН РАН

Ащеулова Н.А., канд. соц. наук, Санкт-Петербургский филиал ИИЕТ РАН

Бакаев М.А., канд. техн. наук, Новосибирский государственный технический университет

Блинова О.В., канд. филол. наук, Санкт-Петербургский государственный университет

Богачева Н.В., канд. псих. наук, Первый Московский государственный медицинский университет имени И.М. Сеченова

Богдановская И.М., канд. псих. наук, Российский государственный педагогический университет имени А.И. Герцена

Болгов Р.В., канд. политич. наук, Санкт-Петербургский государственный университет

Борисов Н.В., докт. физ.-мат. наук, Санкт-Петербургский государственный университет

Войсунский А.Е., канд. психол. наук, Московский государственный университет имени М.В. Ломоносова

Горбунов-Посадов М.М., докт. физ.-мат. наук, Институт прикладной математики имени М.В. Келдыша РАН

Еникеева Е.В., компания «Яндекс»

Захаров В.П., канд. филол. наук, Санкт-Петербургский государственный университет

Комалова Л.Р., докт. филол. наук, ИНИОН РАН

Конюховский П.В., докт. экон. наук, Российский государственный педагогический университет имени А.И. Герцена

Королева Н.Н., докт. психол. наук, Российский государственный педагогический университет имени А. И. Герцена

Корытникова Н.В., канд. социол. наук, Харьковский национальный университет имени В.Н. Каразина

Митрофанова О.А., канд. филол. наук, Санкт-Петербургский государственный университет

Паничева П.В., канд. филол. наук, НИУ Высшая школа экономики

Проект Ю.Л., канд. психол. наук, Российский государственный педагогический университет имени А.И. Герцена

Прокудин Д.Е., докт. филос. наук, Санкт-Петербургский государственный университет

Райков А.Н., докт. техн. наук, Институт проблем управления РАН

Сморгунов Л.В., докт. филос. наук, Санкт-Петербургский государственный университет

Соколов А.В., канд. политич. наук, Ярославский государственный университет имени П.Г. Демидова

Толстикова И.И., канд. филос. наук, Университет ИТМО

Тропинова Е.А., канд. экон. наук, Санкт-Петербургский государственный университет

Федосов А.Ю., докт. пед. наук, Российский государственный социальный университет

Филатова О.Г., канд. филос. наук, Санкт-Петербургский государственный университет

Ходачек И.А., PhD, Северо-Западный институт управления РАНХиГС

Хохлова М.В., канд. филол. наук, Санкт-Петербургский государственный университет

Чугунов А.В., канд. политич. наук, Университет ИТМО

Faiola A., PhD, Professor and Head, Department of Biomedical and Health Information Sciences at the College of Applied Health Sciences, The University of Illinois at Chicago, USA

Mukhamediev R.I., Dr.sc.ing, Satbayev University, Kazakhstan

Petkevic V., Assoc. Prof., RNDr., CSc., Institute of Theoretical and Computational Linguistics,
Prague, Czech Republic

Reips U.-D., Professor and Head of Chair, Department of Psychology, University of Konstanz,
Germany.

Scrivner O., PhD, Research Scientist, Indiana University, Bloomington, USA

Sootla G., Prof. of Public Policy, School of Governance, Law and Society, Tallinn University, Estonia

Tao Y., Dr.Sc., Associate Professor, Shanxi Normal University, Xian, China.

Waterworth J., PhD, Professor of Informatics, Department of Informatics, Umea University, Sweden

ОРГАНИЗАЦИОННЫЙ КОМИТЕТ

Председатель оргкомитета:

Борисов Н.В., докт. физ.-мат. наук, заведующий кафедрой информационных систем в искусстве
и гуманитарных науках Санкт-Петербургского государственного университета

Заместитель председателя оргкомитета:

Прокудин Д.Е., докт. филос. наук, доцент СПбГУ, аналитик Центра юзабилити и смешанной
реальности Университет ИТМО

Члены оргкомитета

Дрожжин А.И. Университет ИТМО

Захаров В.П., СПбГУ

Королёва Н.Н., РГПУ им. А.И. Герцена

Михайлова В.В., СЗИУ РАНХиГС

Низомутдинов Б.А., Университет ИТМО, НП ПРИОР Северо-Запад

Смирнова П.В., Университет ИТМО (информационный менеджер конференции)

Толстикова И.И., Университет ИТМО

Чугунов А.В., Университет ИТМО, НП ПРИОР Северо-Запад (ученый секретарь конференции)

Предисловие редактора

В последние годы автоматизированная обработка текста и речи стала одной из фундаментальных задач не только в лингвистической науке, но и в жизнедеятельности современного информационного общества, в котором все большее число работающих занято производством, хранением и переработкой информации, особенно высшей ее формы – знаний. Создание и использование компьютерных лингвистических ресурсов востребовано в самых разных приложениях. За прошедшие годы в области компьютерной лингвистики были получены значительные научные и практические результаты. Достаточно упомянуть резкое повышение качества машинного перевода и повсеместное внедрение голосовых помощников. При этом и парадигма исследований в самой науке меняется очень быстро. Новые направления и точки роста в компьютерной лингвистике получили названия машинное обучение, глубокое обучение и нейронные сети, которые являются подмножествами искусственного интеллекта и основываются на мощных вычислительных ресурсах и огромных массивах данных, на которых системы самообучаются (то, что еще вчера называлось Big Data). И сегодня компьютерная лингвистика все теснее сливается с математикой.

Статьи, публикуемые в данном сборнике, представляют собой изложение докладов, сделанных на двух семинарах, «Компьютерная лингвистика и вычислительные онтологии» и «International Workshop Computational Linguistics (CompLing-2020)», являющихся частью Международной объединенной конференции «Интернет и современное общество» (IMS 2020). Специфика конференции 2020 г. заключалась в том, что она впервые проходила в режиме онлайн, что не помешало плодотворной работе и даже расширило круг участников. Часть докладов семинара CompLing-2020 печатается в этом сборнике на английском языке, а часть была отобрана для сборника IMS2020 Proceedings (электронное издание CEUR-WS: Free Open-Access Proceedings for Scientific Conferences and Workshops).

Сборник «Компьютерная лингвистика и вычислительные онтологии» представляет собой новое периодическое ежегодное издание, которое выходит уже четвертый раз. Что касается семинара с одноименным названием, то это уже тринадцатый семинар в составе данной конференции, первый из них состоялся в 2008 г. И считаю своим долгом напомнить, что у истоков этого направления конференции стоял Валерий Шлёмович Рубашкин.

Публикуемые здесь статьи отражают широкую тематику исследований по компьютерной лингвистике и многогранность задач в области автоматической обработки текста и речи и имеют как теоретическое, так и прикладное значение.

Значительная часть статей посвящена исследованиям в области корпусной лингвистики. При этом сам предмет исследования в каждой из них весьма разный. Так, в работе В. Бенко и К. Раусовой описывается извлечение из русскоязычного корпуса фраз на латыни (как части русскоязычного дискурса) и показывается, что по-настоящему это стало реальным только с появлением больших корпусов. Другие же статьи, как, например, работа М.В. Хохловой и Е.В. Еникеевой, посвященная методам машинного обучения применительно к задаче выделения глагольных и атрибутивных коллокаций, вводит в корпусную лингвистику новые методы исследования. И это хорошо демонстрирует широкий спектр направлений и методов современной корпусной лингвистики.

Во втором разделе сборника собраны работы, посвященные или примыкающие к семантике. И точно так же можно показать их разнообразие, сравнив, например исследование В.О. Кораблинова о подготовке набора данных для вопросно-ответного поиска по базе знаний с методами машинного обучения применительно к задаче генерации правил для автоматической проверки правописания, описанными в работе П.Я. Бахвалова.

Из материалов этого небольшого сборника хорошо видна динамика развития компьютерной лингвистики, стоящей перед лицом новых задач, связанных с ее междисциплинарностью, с выходом за пределы собственно лингвистики.

В заключение хочу поблагодарить всех авторов и рецензентов, обеспечивших высокий уровень статей сборника, а также оргкомитет конференции «Интернет и современное общество» за работу по организации конференции и по изданию сборника трудов нашего международного семинара.

В.П. Захаров

Применение деревьев решений для анализа сильных позиций текста в задаче атрибуции произведений Ф. М. Достоевского

А.А. Рогов, А.А. Лебедев, Р.В. Абрамов, Н.Д. Москин, К.А. Кулаков

Петрозаводский государственный университет

rogov@petsu.ru, perevodchik88@yandex.ru, monset008@gmail.com,
moskin@petsu.ru, kulakov@cs.karelia.ru

Аннотация

В работе рассматривается совокупность статей Ф. М. Достоевского и других авторов (М. М. Достоевский, Н. Н. Страхов, А. А. Головачев, И. Н. Шиль, А. Григорьев, А. У. Порецкий, Я. П. Полонский), опубликованных в журналах «Время» и «Эпоха» в период 1861-1865 гг. В текстах выделялись фрагменты размером 500, 700 и 1000 слов. При этом для увеличения объема выборки использовался шаг для отсчета начала следующего фрагмента: 100, 200 слов и т. п. На основе частеречного распределения фрагментов текстов были построены деревья решений, в узлах которых находятся условия ветвления, основанные на частоте встречаемости той или иной n -граммы (последовательности из n закодированных частей речи).

Анализ сильных позиций данных текстов (т. е. фрагментов, расположенных в начале или в конце текста) с помощью деревьев решений показывает возможность стилистической правки, которую вносил Ф. М. Достоевский в тексты изначальных авторов. Для проведения исследования использовалась информационная система СМАЛТ («Статистические методы анализа литературных текстов»), где была реализована автоматизированная разметка произведений с ручным контролем специалистов-филологов.

Ключевые слова: атрибуция текстов, корпусная лингвистика, Ф. М. Достоевский, сильные позиции текста, дерево решений, n -грамма, частеречное распределение

Библиографическая ссылка: Рогов А.А., Лебедев А.А., Абрамов Р.В., Москин Н.Д., Кулаков К.А. Применение деревьев решений для анализа сильных позиций текста в задаче атрибуции произведений Ф. М. Достоевского // Компьютерная лингвистика и вычислительные онтологии. Выпуск 4 (Труды XXIII Международной объединенной научной конференции «Интернет и современное общество», IMS-2020, Санкт-Петербург, 17 – 20 июня 2020 г. Сборник научных статей). — СПб.: Университет ИТМО, 2020. С. 118-127. DOI: 10.17586/0000-0000-2020-4-118-127

Введение

Чтобы определить особенности произведений конкретного автора на фоне всего массива текстов художественной литературы, необходимо представить комплексный анализ его текстов, выделить ключевые особенности авторского восприятия мира, а также определить, как эти особенности воплощаются в произведениях. Вряд ли можно поставить под сомнение тот факт, что «каждый писатель подбирает языковые средства не только в соответствии с содержанием и замыслом, но и в зависимости от своего видения мира, обусловленного его мировосприятием, социальным положением, личностными качествами и психологическими особенностями» [1, с. 50]. Однако реализация подобранных автором языковых средств в тексте литературного произведения может быть самой разной.

К классификации авторского стиля как лингвистического явления существует несколько подходов; в зависимости от выбранного варианта исследователь направляет свою работу в нужное ему русло. В частности, можно выделить многоаспектный подход, предложенный в работе [2] и опирающийся на принципы коммуникации (рассмотрение отдельных компонентов художественной системы писателя, прежде всего, рассмотрения языковых средств (как правило, лексических) в сочетании с анализом различных структурных и смысловых форм организации языкового материала для выявления характера соотнесенности на фоне других идиостилей); анализ авторского стиля на основании экстралингвистических и интралингвистических факторов [3]; когнитивный аспект в анализе индивидуально-авторского стиля (истоки подобного подхода к идиостилю лежат в развитии теоретических положений А. А. Потебни, Л. В. Щербы, Г. О. Винокура, В. В. Виноградова, Б. А. Ларина), когда обращение к концептосфере писателя представляет собой возможность лучше понять художественный мир автора (см., например, [4; 5]).

Мы, вслед за В. В. Виноградовым, отмечавшим, что «изучение литературного стиля должно быть комплексным и системным» [6, с. 198] понимаем идиостиль как систему формальных и содержательных характеристик, которые присущи произведениям того или иного автора и отражают уникальный, авторский способ языкового выражения (подробнее см. [7]). Воплощение подобного авторского способа может проявляться на разных уровнях текста, в том числе и в анализе сильных его позиций.

Вопрос, связанный с выделением сильных позиций текста и определением их роли в понимании и восприятии произведений, не просто является объектом давнего интереса лингвистов [8; 9], но и в последнее десятилетие решается в аспекте анализа индивидуально-авторского стиля и определения особенностей творчества того или иного писателя [10; 11]. Под универсальными сильными позициями текста традиционно понимаются:

- заглавие текста;
- инициальная позиция текста (его первое предложение, первый абзац, первое сложное синтаксическое целое);
- конечная позиция текста (последнее предложение текста, последний его абзац, последнее сложное синтаксическое целое).

При этом в зависимости от жанра текста возможно появление и других сильных позиций (например, эпиграфа в художественных произведениях, рифмы – в стихотворных текстах, слогана – в рекламных текстах и т. п.), однако составляющие текста, перечисленные выше, следует признать универсальными – то есть, характерными для любого произведения. Именно в сильных позициях автор сосредотачивает важные для себя смысловые доминанты, передаваемые речевым произведением.

Поскольку сильные позиции текста играют наиболее значимую роль в выражении авторской идеи, то очевиден дополнительный интерес филологов именно к этим элементам текста при решении вопросов, связанных с атрибуцией текстов [12; 13], установлением авторства анонимных текстов [14], а также в ходе лингвистического анализа материалов, на которые, помимо автора, мог оказать непосредственное влияние другой человек.

В частности, анализируя тексты статей, представленных в журналах «Время» и «Эпоха» (1861-1865 гг.), мы обнаруживаем немалый список публиковавшихся там авторов (Ф. М. Достоевский, М. М. Достоевский, Н. Н. Страхов, А. А. Головачев, И. Н. Шилль, А. Григорьев, А. У. Порецкий, Я. П. Полонский), и, рассматривая морфологическую структуру некоторых материалов, казалось бы, не принадлежащих непосредственно Фёдору Михайловичу Достоевскому, мы можем говорить о непосредственном его влиянии на общее построение и смысловое наполнение текста (что находит отражение в том числе и в перечисленных выше сильных позициях).

Использованная в работе методика исследования предусматривает работу с большими фрагментами текста (500, 700, 1000 слов), что делает малозначимой позицию заглавия текста; однако инициальная и конечная позиции текста при таком подходе становятся объектом дополнительного интереса и могут привлекать внимание филолога в тех случаях, когда частеречное распределение в них отличается от распределения по частям речи во всем остальном тексте. Выбор в качестве признаков только частей речи трех первых позиций текста опирается на исследования, проведенные Г. Хетсо [15], который учитывает такие признаки, как общее распределение частей речи в первых двух позициях предложения и сочетание частей речи в первых двух позициях предложения. В данном исследовании методика Г. Хетсо была дополнена – под сильной позицией в рамках исследования понимаются три слова в началах каждого предложения текста.

Целью проводимого эксперимента была автоматизация выявления наиболее важных, фундаментальных отличительных характеристик текстов одного автора от текстов других авторов в аспекте распределения частей речи в сильных позициях текста (началах предложений) в сочетании с их последующей точечной лингвистической интерпретацией. В перспективе это должно помочь решить проблему проверки принадлежности перу Ф. М. Достоевского ряда статей, которые входят в раздел *Dubia* полного собрания сочинений писателя. Использование математических методов в этом случае позволяет лингвисту не просто избавиться от необходимости вручную выполнять подсчеты, но и помогает определить наиболее значимые отличия, которые могут быть не замечены или проигнорированы как нерелевантные при традиционном подсчете. Для решения задачи атрибуции текстов хорошо зарекомендовали себя следующие математические методы [12; 13; 14]: нейронные сети, деревья решений, машина опорных векторов (SVG), метод *k*-средних, метод QSUM, байесовский классификатор, марковские цепи, метод главных компонент, дискриминантный анализ, генетические алгоритмы, статистические критерии (хи-квадрат Пирсона, критерий Стьюдента, Колмогорова-Смирнова) и др. Среди этих методов интеллектуального анализа данных деревья решений выделяются тем, что они просты в понимании и интерпретации, а также не требуют специальной предварительной обработки данных. Построение дерева решений позволяет однозначно определить, какие из признаков являются наиболее существенными, а потому требующими дополнительной интерпретации в аспекте лингвостилистического анализа – они будут вынесены в вершину дерева решений.

1. Построение и анализ деревьев решений

В табл. 1 представлен список из 19 анализируемых текстов (среди авторов: Ф. М. Достоевский, М. М. Достоевский, Н. Н. Страхов, А. А. Головачев, И. Н. Шилль, А. Григорьев, А. У. Порецкий, Я. П. Полонский). 12 из них принадлежат Ф. М. Достоевскому, остальные – другим авторам. Выбор текстов был осуществлен случайным образом с учетом распределения публикаций этих авторов на страницах журналов «Время» и «Эпоха» (1861-1865 гг.). Грамматическая разметка данных текстов учитывала 14 частей речи (существительное, прилагательное, числительное, местоимение, наречие, категория состояния, глагол, причастие, деепричастие, предлог, союз, частица, модальное слово, междометие), а также позволяла выделять цитаты, иностранные слова, вводные слова, сокращенные слова и неязыковые символы.

Для проведения разметки использовалась информационная система СМАЛТ («Статистические методы анализа литературных текстов»), разработанная в Петрозаводском государственном университете [16]. В рамках системы выполнялся импорт текстовых произведений с автоматизированным разбиением на абзацы, предложения и слова; редактирование подобранного автоматически в ходе импорта морфологического разбора (часть речи и другие параметры) или создание нового разбора филологом; построение фрагментов текстов с заданным размером слов и отступом

и с описанием морфологического разбора для детального анализа. Таким образом, в системе СМАЛТ была выполнена автоматизированная разметка произведений с ручным контролем специалистов-филологов.

Таблица 1. Исходные тексты для анализа

| Код | Название | Автор | Журнал | Год | № журнала | Количество слов |
|-----|--|--------------------|--------|------|-----------|-----------------|
| 2 | Пожары | Федор Достоевский | Время | 1862 | 1 | 1943 |
| 11 | Тарась Шевченко | Аполлон Григорьев | Время | 1861 | 4 | 1724 |
| 13 | Письмо к редактору | Полонский Я. П. | Время | 1863 | 3 | 2303 |
| 34 | Литературная истерика | Федор Достоевский | Время | 1861 | 7 | 2808 |
| 35 | Молодое перо | Федор Достоевский | Время | 1863 | 2 | 1872 |
| 40 | Подписка на 1863 годъ | Михаил Достоевский | Время | 1863 | 1 | 2541 |
| 42 | Ряд статей о русской литературе. Введение | Федор Достоевский | Время | 1861 | 1 | 12508 |
| 43 | Славянофилы, черногорцы и западники | Федор Достоевский | Время | 1862 | 9 | 2058 |
| 75 | Ряд статей... Г. -бов и вопрос об искусстве | Федор Достоевский | Время | 1861 | 2 | 11053 |
| 77 | Книжность и грамотность. Статья вторая | Федор Достоевский | Время | 1861 | 8 | 14210 |
| 78 | Последние литературные явления. Газета "День" | Федор Достоевский | Время | 1861 | 11 | 4323 |
| 82 | Необходимое литературное объяснение, по поводу ра... | Федор Достоевский | Время | 1863 | 1 | 3680 |
| 86 | Чтобы кончить. Последніе объясненія съ "Современн... | Федор Достоевский | Эпоха | 1864 | 9 | 1378 |
| 87 | Политическое обозрение | Головачев А.А. | Эпоха | 1864 | 8 | 10309 |
| 89 | Лермонтов и его направление. Статья вторая | Аполлон Григорьев | Время | 1862 | 11 | 7480 |
| 92 | Наши домашние дела | Порецкий А.У. | Эпоха | 1864 | 12 | 8602 |
| 96 | Голос за петербургского Дон-Кихота. (По поводу ст... | Федор Достоевский | Время | 1862 | 10 | 1334 |
| 97 | Примечание | Федор Достоевский | Эпоха | 1864 | 9 | 1599 |
| 116 | ДУРНЫЕ ПРИЗНАКИ | Страхов Н. Н. | Время | 1862 | 11 | 6331 |

Далее текст разбивался на фрагменты размером 500, 700 и 1000 слов. Если заключительный фрагмент оказывался меньше указанного размера, то он отбрасывался. Для увеличения объема выборки применялся шаг для отсчета начала следующего фрагмента: 100, 200 слов и т. п. Поскольку на первом этапе было установлено, какие части речи встречаются во фрагментах и в каком порядке, на втором этапе можно перейти к подсчету частоты встречаемости n -грамм (последовательностей из n закодированных частей речи) [17].

Эти частоты легли в основу построения деревьев решений [18]. Каждая вершина дерева определяет условие ветвления по одному из признаков, что позволяет в дальнейшем

классифицировать тексты на основе их частеречного распределения. Подобные модели широко используются в интеллектуальном анализе данных, и, в частности, при анализе текстов. Например, их применяла А. Р. Дубовик для исследования принадлежности русскоязычных текстов к тому или иному функциональному стилю (научный, художественный, деловой или публицистический) с опорой на ряд статистических параметров, таких как средняя длина слова, средняя длина предложения, частота встречаемости в текстах определенных n-грамм и др. [19]. Также деревья решений применялись в задаче разграничения фольклорных текстов и текстов, стилизованных под фольклор [20].

На рисунке 1 показан фрагмент графа, полученного в результате расчетов с шагом 100 слов и размером фрагмента 1000 слов для биграмм. Достоинства данного метода заключается в том, что он прост в понимании и интерпретации, а также не требует специальной подготовки данных.

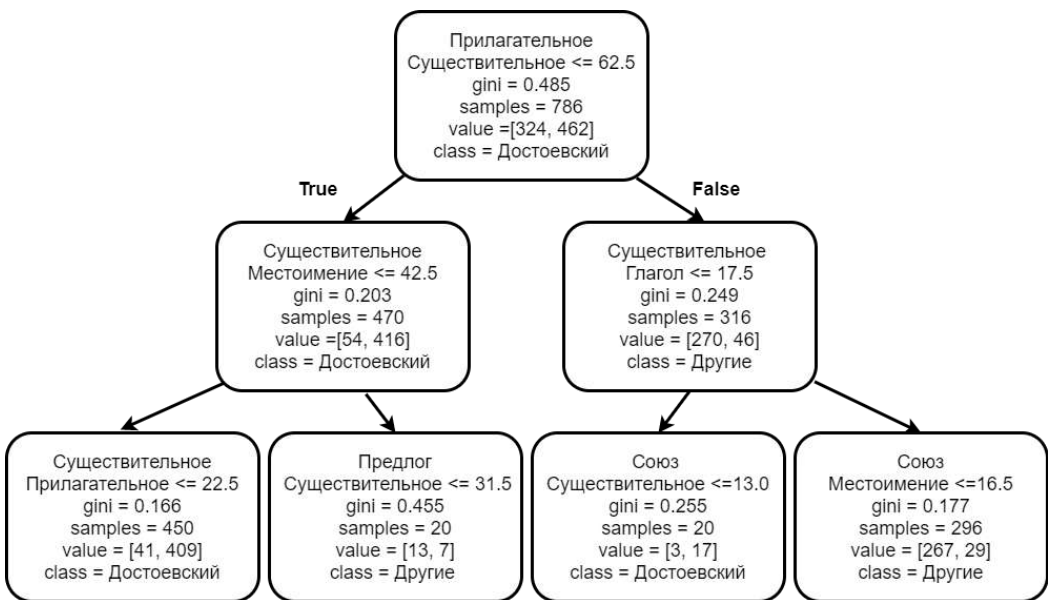


Рис. 1. Фрагмент дерева решений (биграмма, шаг 100 слов, фрагмент 1000 слов)

2. Результаты анализа и их интерпретация

Анализируя дерево решений, представленное на рисунке 1, можно заметить, что фрагменты текстов статьёй Ф. М. Достоевского, как правило, содержат частоту встречаемости биграммы «Прилагательное + Существительное» меньше 62,5. Анализ этих фрагментов показывает, что среди неверно классифицированных фрагментов статей, которые принадлежат другим авторам (не Фёдору Михайловичу Достоевскому), есть и те, что расположены в начале или в конце текста, то есть в его сильных позициях. Это может быть свидетельством стилистической правки, которую вносил Ф. М. Достоевский в данные фрагменты, что сказывалось на частеречном распределении текстов и потому не соответствовало стилю изначального автора статьи.

В таблице 2 приведены тексты с указанием нетипичных фрагментов (значительная часть которых располагается именно в сильных позициях текста). Данные фрагменты, будучи правильно интерпретированными, могут стать объектом дополнительного интереса со стороны литературоведов в контексте определения вопроса авторства статей

и стилистического взаимовлияния Ф. М. Достоевского и других авторов «Времени» и «Эпохи».

Таблица 2. Аномалии в классификации фрагментов статей

| Код | Название | Автор | Журнал | Год | № журна ла | Номера нетипичных фрагментов |
|-----|-----------------------|--------------------|--------|------|---------------|------------------------------------|
| 2 | Пожары | Федор Достоевский | Время | 1862 | 1 | 0-8 |
| 11 | Тарась Шевченко | Аполлон Григорьев | Время | 1861 | 4 | 0, 4 |
| 13 | Письмо к редактору | Полонский Я.П. | Время | 1863 | 3 | 10-12 |
| 35 | Молодое перо | Федор Достоевский | Время | 1863 | 2 | 0-2 |
| 40 | Подписка на 1863 годъ | Михаил Достоевский | Время | 1863 | 1 | 0-14 |
| 92 | Наши домашние дела | Порецкий А. У. | Эпоха | 1864 | 12 | 77, 78 |
| 116 | ДУРНЫЕ ПРИЗНАКИ | Страхов Н. Н. | Время | 1862 | 11 | 51, 52 |

К примеру, в тексте «Наши домашние дела» с кодом 92 («Эпоха», 1864, № 12), принадлежащем А. У. Порецкому, именно два последних фрагмента текста имеют иное частеречное распределение в сравнении со всем остальным текстом статьи (в частности, реже встречаются пары «Прилагательное + Существительное» и «Существительное + Местоимение»), что в целом нехарактерно для всех остальных 76 фрагментов данной статьи. В тексте «Письмо к редактору», автором которого является Я. П. Полонский, («Время», 1863, № 3, код 13) последние три фрагмента имеют другое распределение. В тексте «Тарась Шевченко» Аполлона Григорьева («Время», 1861, № 4, код 11) обращает на себя внимание первый фрагмент и середина текста. В тексте «ДУРНЫЕ ПРИЗНАКИ» Н. Н. Страхова («Время», 1862, № 11, код 116) также явно выделяются два последних фрагмента.

Вероятно, не всегда причины подобного рода изменений в распределении частей речи следует искать именно в стилистической правке со стороны Ф. М. Достоевского. К примеру, в упомянутом выше тексте «Письмо к редактору» причиной снижения случаев встречаемости комбинации частей речи «Существительное» + «Местоимение» может быть переход автора от размещаемых в начале текста пространных теоретических размышлений о художественных достоинствах повести Л. Н. Толстого «Казачья», к описанию совершаемых в произведении действий, выраженных глаголами (*люди живут как живет природа, умирают, рождаются, совокупляются, опять рождаются, дерутся, пьют, љдят, радуются и опять умирают*), равно как и цитированию чужого текста, что сказывается на частеречном распределении.

В некоторых текстах появление определенных пар частей речи может объясняться тематикой произведения. К примеру, в статье «Тарас Шевченко» (автор: Аполлон Григорьев) активизация во многих фрагментах текста пары «Прилагательное + Существительное» может объясняться характером описываемого явления (для характеристики произведений Шевченко неоднократно используется словосочетание *малороссійская литература*), а сам Григорьев очень любит пользоваться эпитетом великий, совмещая его с разными существительными: *великому таланту, великими представителями, великаго поэта, великой литературы, великому кобзарю* и т. п.).

В то же время, даже предположить причину изменения частеречного распределения удастся не всегда. К примеру, завершающие фрагменты текстов «Наши домашние дела» и «Дурные признаки», не отличаясь чем-то содержательно или тематически в сравнении с остальным содержанием статьи, классифицируются по-иному. Именно это несоответствие изначальному авторскому стилю может указывать на внесение стилистических правок неким иным лицом (возможно, самим Ф. М. Достоевским).

Заметим, что приведенный признак «частота биграммы «Прилагательное + Существительное» меньше 62,5» не является однозначным признаком стиля

Ф. М. Достоевского. Так, например, «Подписка на 1863 год» («Время», 1863, № 1, код 40) приписываемая М. М. Достоевскому целиком удовлетворяет этому параметру. Весь текст «Пожары» («Время», 1862, № 1, код 2) Ф. М. Достоевского не удовлетворяет этому признаку. Начало (три фрагмента) текста «Молодое перо» («Время», 1863, № 2, код 35) тоже не удовлетворяет этому параметру.

С учетом вышенаписанного текст статьи «Подписка на 1863 год» может стать объектом более пристального изучения литературоведов – весьма вероятно, что Фёдор Михайлович Достоевский мог оказать немалое влияние на Михаила Михайловича (как опосредованное, так и прямое – если предположить, что он составил какую-то часть текста вместо своего старшего брата). Не менее интересен и вопрос авторства текста «Пожары», ответ на который по мнению некоторых литературоведов не столь однозначен, как это может показаться – вполне возможно, что это был не Ф. М. Достоевский (см. статью Н. Г. Розенблюма «Петербургские пожары 1862 г. и Достоевский» [21]).

Выводы

Ф. М. Достоевский, будучи мастером работы не только с художественным, но и с публицистическим текстом, прекрасно представлял себе важность сильных позиций текста с точки зрения их воздействия на читателя, а потому мог уделять более пристальное внимание внесению правок в начальные и в конечные абзацы текстов чужих статей. Ведь именно эти элементы оказывают наибольшее воздействие на читателя – благодаря началу текста складывается первое впечатление о статье, а финал – подводит итоги и, как правило, хорошо откладывается в памяти читателя. Именно поэтому, комплексно решая вопросы, связанные с атрибуцией текстов в журналах «Время» и «Эпоха», в качестве одной из задач можно выделить специальный анализ данных элементов текста.

Приведенный метод предоставляет возможность выделять в тексте фрагменты, отличающиеся частеречными характеристиками. Анализ этих фрагментов позволяет исследователю ставить и решать различные задачи, начиная от авторства тех или иных фрагментов до автоматического реферирования текстов.

Работа выполнена при поддержке Российского фонда фундаментальных исследований, грант № 18-012-90026.

Литература

- [1] Гаспарян С.К., Князян А.Т. К вопросу об изучении индивидуального стиля автора // Филологические науки. 2004. № 4. С. 50 – 57.
- [2] Болотнова Н.С. и др. Коммуникативная стилистика художественного текста: лексическая структура и идиостиль / Н.С. Болотнова, И.И. Бабенко, А.А. Васильева, С.М. Карпенко, О.В. Орлова, С.В. Сыпченко, Р.Я. Тюрина. Томск. 2001. 331 с.
- [3] Шаркунова О.В. Идиостиль художественного текста как индивидуальное сочетание экстра- и интралингвистических параметров, основанных на референтных отношениях // Материалы международной заочной научно-практической конференции «Актуальные вопросы филологии, искусствоведения и культурологии». Новосибирск. 2011. С. 75 – 80.
- [4] Тарасова И.А. Категории когнитивной лингвистики в исследовании идиостиля // Вестник СамГУ. 2004. № 1 (31). С. 163 – 169.
- [5] Фокина Ю.М. Особенности репрезентации индивидуально-авторской концептосферы в англоязычной и русскоязычной прозе (на материале рассказов А.П. Чехова и Д. Джойса). Автореф. дис. ... канд. филол. наук. Саратов, 2010. 184 с.
- [6] Виноградов В.В. Проблема авторства и теория стилей // М.: Художественная литература. 1961. 613 с.

- [7] Лебедев А.А. Идиостиль П. А. Вяземского: синтаксический аспект // Петрозаводск: Изд-во ПетрГУ. 2013. 134 с.
- [8] Арнольд И.В. Значение сильной позиции для интерпретации художественного текста // Иностранные языки в школе. М. 1978. № 4. С. 23 – 31.
- [9] Гальперин И.Р. Текст, как объект лингвистического исследования // М.: Наука, 1981. 139 с.
- [10] Патроева Н.В., Лебедев А.А. Синтаксическая организация, размер и семантика инициальных предложений в лирике А. С. Пушкина // Вестник Томского государственного университета. Филология. Томск. 2018. № 53. С. 224 – 236.
- [11] Петрова К.В. Сильные позиции текста в автобиографии Джанет Уинтерсон // Вестник Новгородского государственного университета. Серия «Гуманитарные науки». Великий Новгород. 2015. № 4 (87), часть 1. С. 73 – 76.
- [12] Рогов А.А. и др. Математические методы атрибуции текстов / А.А. Рогов, А.В. Седов, Ю.В. Сидоров, Т.Г. Суровцова // Петрозаводск: Изд-во ПетрГУ. 2014. 95 с.
- [13] Stamataatos E. A Survey of Modern Authorship Attribution Methods // Journal of the American Society for Information Science and Technology. 2009. Vol. 60, № 3. P. 538 – 556.
- [14] Романов А.С. Методика и программный комплекс для идентификации автора неизвестного текста. Дис. ... канд. техн. наук. Томск, 2010. 149 с.
- [15] Kjetsaa G. Attributed to Dostoevsky: The Problem of attributing to Dostoevsky anonymous articles in Time and Epoch. Oslo: Solum Forlag A. S., 1986.
- [16] Рогов А.А., Кулаков К.А., Москин Н.Д. Программная поддержка в решении задачи атрибуции текстов // Программная инженерия. М.: Изд-во "Новые технологии", 2019. Т. 10, № 5. С. 234 – 240.
- [17] Котов А.А. и др. Лингвистические корпусы / А.А. Котов, З.И. Минеева, А.А. Рогов, А.В. Седов, Ю.В. Сидоров // Петрозаводск: Изд-во ПетрГУ, 2014.
- [18] Breiman L., etc. Classification and regression trees / L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone. Wadsworth, Belmont, Ca, 1984.
- [19] Дубовик А.Р. Автоматическое определение стилистической принадлежности текстов по их статистическим параметрам // Компьютерная лингвистика и вычислительные онтологии. Выпуск 1. Труды XX Международной научной конференции «Интернет и современное общество», IMS-2017 (Санкт-Петербург, 21 – 23 июня 2017 г.). СПб.: Университет ИТМО, 2017. С. 29 – 45.
- [20] Щеголева Л.В., Лебедев А.А., Москин Н.Д. Методы анализа данных в задаче разграничения фольклорных и авторских текстов // Вопросы языкознания. Москва, 2020. № 2. С. 61 – 74.
- [21] Розенблюм Н.Г. Петербургские пожары 1862 г. и Достоевский (запрещенные цензурой статьи журнала «Время») // Достоевский Ф. М. Новые материалы и исследования. Литературное наследство. М., 1973. Т. 86. С. 16 – 54.

Application of Decision Trees for Analyzing the Strong Positions of the Text in the Problem of Attribution of Works by F. M. Dostoevsky

A. Rogov, A. Lebedev, R. Abramov, N. Moskin, K. Kulakov

Petrozavodsk State University

The paper considers a set of articles by F. M. Dostoevsky and other authors (M. M. Dostoevsky, N. N. Strakhov, A. A. Golovachev, I. N. Shill, A. Grigoriev, A. U. Poretsky, Y. P. Polonsky) published in the magazines "Time" and "Epoch" in the period 1861-1865. Fragments of 500, 700 and 1000 words were selected in the texts. At the same time a step (100, 200 words and so on) was used to count the beginning of the next fragment to increase the sample size. Based on the

distribution of parts of speech of text fragments, decision trees were constructed, whose nodes contain branching conditions based on the frequency of occurrence of a particular n-gram (a sequence of n encoded parts of speech).

Analysis of the strong positions of these texts (i.e. fragments located at the beginning or end of the text) with the help of decision trees shows the possibility of stylistic editing, which was made by F. M. Dostoevsky in the texts of the original authors. The SMALT information system ("Statistical Methods of Analysis of Literary Texts") was used to conduct the study, where automated markup of texts with manual control of specialists of philology was implemented.

Keywords: text attribution, corpus linguistics, F. M. Dostoevsky, strong positions of the text, decision tree, n-gram, partial distribution

Reference for citation: Rogov A., Lebedev A., Abramov R., Moskin N., Kulakov K. Application of decision trees for analyzing the strong positions of the text in the problem of attribution of works by F. M. Dostoevsky // *Computer Linguistics and Computing Ontologies*. Vol. 4 (Proceedings of the XXIII International Joint Scientific Conference «Internet and Modern Society», IMS-2020, St. Petersburg, June 17-20, 2020). - St. Petersburg: ITMO University, 2020. P. 118 – 127. DOI: 10.17586/0000-0000-2020-4-118-127

References

- [1] Gasparyan S.K., Knyazyan A.T. K voprosu ob izuchenii individual'nogo stilya avtora // *Filologicheskie nauki*. 2004. № 4. S. 50 – 57. [In Russian].
- [2] Bolotnova N.S. i dr. Kommunikativnaya stilistika hudozhestvennogo teksta: leksicheskaya struktura i idiosil' / N.S. Bolotnova, I.I. Babenko, A.A. Vasil'eva, S.M. Karpenko, O.V. Orlova, S.V. Sypchenko, R.YA. Tyurina. Tomsk. 2001. 331 s. [In Russian].
- [3] SHarkunova O.V. Idiosil' hudozhestvennogo teksta kak individual'noe sochetanie ekstra- i intralingvisticheskikh parametrov, osnovannykh na referentnykh otnosheniyah // *Materialy mezhdunarodnoj zaochnoj nauchno-prakticheskoy konferencii «Aktual'nye voprosy filologii, iskusstvovedeniya i kulturologii»*. Novosibirsk. 2011. S. 75 – 80. [In Russian].
- [4] Tarasova I.A. Kategorii kognitivnoj lingvistiki v issledovanii idiosilya // *Vestnik SamGU*. 2004. № 1 (31). S. 163 – 169. [In Russian].
- [5] Fokina YU.M. Osobennosti reprezentatsii individual'no-avtorskoj konceptosfery v angloyazychnoj i russkoyazychnoj proze (na materiale rasskazov A.P. CHEkhova i D. Dzhosja). Avtoref. dis. ... kand. filol. nauk. Saratov, 2010. 184 s. [In Russian].
- [6] Vinogradov V.V. Problema avtorstva i teoriya stilej // *M.: Hudozhestvennaya literatura*. 1961. 613 s. [In Russian].
- [7] Lebedev A.A. Idiosil' P. A. Vyazemskogo: sintaksicheskij aspekt // *Petrozavodsk: Izd-vo PetrGU*. 2013. 134 s. [In Russian].
- [8] Arnol'd I.V. Znachenie sil'noj pozitsii dlya interpretatsii hudozhestvennogo teksta // *Inostrannye yazyki v shkole*. M. 1978. № 4. S. 23 – 31. [In Russian].
- [9] Gal'perin I.R. Tekst, kak ob'ekt lingvisticheskogo issledovaniya // *M.: Nauka*, 1981. 139 s. [In Russian].
- [10] Patroeva N.V., Lebedev A.A. Sintaksicheskaya organizatsiya, razmer i semantika inicial'nykh predlozhenij v lirike A. S. Pushkina // *Vestnik Tomskogo gosudarstvennogo universiteta. Filologiya*. Tomsk. 2018. № 53. S. 224 – 236. [In Russian].
- [11] Petrova K.V. Sil'nye pozitsii teksta v avtobiografii Dzhaneit Uinterson // *Vestnik Novgorodskogo gosudarstvennogo universiteta. Seriya «Gumanitarnye nauki»*. Velikij Novgorod. 2015. № 4 (87), chast' 1. S. 73 – 76. [In Russian].
- [12] Rogov A.A. i dr. Matematicheskie metody atribucii tekstov / A.A. Rogov, A.V. Sedov, YU.V. Sidorov, T.G. Surovcova // *Petrozavodsk: Izd-vo PetrGU*. 2014. 95 s. [In Russian].

- [13] Stamatatos E. A Survey of Modern Authorship Attribution Methods // Journal of the American Society for Information Science and Technology. 2009. Vol. 60, № 3. P. 538 – 556.
- [14] Romanov A.S. Metodika i programnyj kompleks dlya identifikacii avtora neizvestnogo teksta. Dis. ... kand. tekhn. nauk. Tomsk, 2010. 149 s. [In Russian].
- [15] Kjetsaa G. Attributed to Dostoevsky: The Problem of attributing to Dostoevsky anonymous articles in Time and Epoch. Oslo: Solum Forlag A. S., 1986.
- [16] Rogov A.A., Kulakov K.A., Moskin N.D. Programmnyaya podderzhka v reshenii zadachi atribucii tekstov // Programmnyaya inzheneriya. M.: Izd-vo "Novye tekhnologii", 2019. T. 10, № 5. S. 234 – 240. [In Russian].
- [17] Kotov A.A. i dr. Lingvisticheskie korpusy / A.A. Kotov, Z.I. Mineeva, A.A. Rogov, A.V. Sedov, YU.V. Sidorov // Petrozavodsk: Izd-vo PetrGU, 2014. [In Russian].
- [18] Breiman L., etc. Classification and regression trees / L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone. Wadsworth, Belmont, Ca, 1984.
- [19] Dubovik A.R. Avtomaticheskoe opredelenie stilisticheskoy prinadlezhnosti tekstov po ih statisticheskim parametram // Komp'yuternaya lingvistika i vychislitel'nye ontologii. Vypusk 1. Trudy XX Mezhdunarodnoj nauchnoj konferencii «Internet i sovremennoe obshchestvo», IMS-2017 (Sankt-Peterburg, 21 – 23 iyunya 2017 g.). SPb.: Universitet ITMO, 2017. S. 29 – 45. [In Russian].
- [20] SHCHegoleva L.V., Lebedev A.A., Moskin N.D. Metody analiza dannyh v zadache razgranicheniya fol'klornyh i avtorskih tekstov // Voprosy yazykoznanija. Moskva, 2020. № 2. S. 61 – 74. [In Russian].
- [21] Rozenblyum N.G. Peterburgskie pozhary 1862 g. i Dostoevskij (zapreshchennye cenzuroj stat'i zhurnala «Vremya») // Dostoevskij F. M. Novye materialy i issledovaniya. Literaturnoe nasledstvo. M., 1973. T. 86. S. 16 – 54. [In Russian].

Сведения об авторах

Абрамов Роман Владимирович, Петрозаводский государственный университет, студент, ORCID 0000-0001-9599-4906.

Бахвалов Павел Ярославович, Университет ИТМО, студент, ORCID 0000-0002-2700-1607.

Бенко Владимир (Benko Vladimír), PhD, научный сотрудник; Братиславский университет им. А.Коменского, доцент, ORCID 0000-0002-4600-5515.

Гребеников Александр Олегович, кандидат филологических наук, Санкт-Петербургский государственный университет, доцент, ORCID 0000-0003-2856-5049.

Джангольская Ольга Владимировна, Санкт-Петербургский государственный университет, магистрант, ORCID 0000-0002-0545-5763.

Добров Алексей Владимирович, кандидат филологических наук, Санкт-Петербургский государственный университет, старший преподаватель, ORCID 0000-0003-0245-5407.

Доброва Анастасия Евгеньевна, ООО «АИРЕ», лингвист, ORCID 0000-0002-8419-1005.

Еникеева Екатерина Владимировна, Компания «Яндекс», аналитик-лингвист, ORCID 0000-0003-4946-9697.

Захарова Алина Андреевна, Санкт-Петербургский государственный университет, бакалавр, ORCID 0000-0003-3196-9327.

Коган Марина Самуиловна, кандидат технических наук, Санкт-Петербургский политехнический университет Петра Великого, доцент, ORCID 0000-0002-7519-2161.

Кораблинов Владислав Олегович, Университет ИТМО, магистрант, ORCID 0000-0002-7983-212X.

Кузнецова Инга Вадимовна, Университет ИТМО, преподаватель, ORCID 0000-0002-4404-3884.

Кулаков Кирилл Александрович, кандидат физико-математических наук, Петрозаводский государственный университет, доцент, ORCID 0000-0002-0305-419X.

Лебедев Александр Александрович, кандидат филологических наук, Петрозаводский государственный университет, старший преподаватель, ORCID 0000-0001-9939-9389.

Марусенко Наталия Михайловна, кандидат филологических наук, Санкт-Петербургский государственный университет, доцент, ORCID 0000-0002-3347-1373.

Микони Станислав Витальевич, доктор технических наук, профессор, Санкт-Петербургский институт информатики и автоматизации РАН, ведущий научный сотрудник, ORCID 0000-0001-7153-6804.

Москин Николай Дмитриевич, кандидат технических наук, доцент, Петрозаводский государственный университет, доцент, ORCID 0000-0001-5556-5349.

Раусова Катарина (Rausová Katarína), Институт языкознания им. Л. Штура Словацкой академии наук, аспирант, ORCID 0000-0002-7855-2839.

Рогов Александр Александрович, доктор технических наук, профессор, Петрозаводский государственный университет, заведующий кафедрой, ORCID 0000-0002-8815-7920.

Смирнова Мария Олеговна, кандидат филологических наук, Санкт-Петербургский государственный университет, доцент, ORCID 0000-0001-5429-2051.

Соколова Елена Григорьевна, кандидат филологических наук, независимый исследователь.

Сомс Николай Леонидович, ООО «АИРЕ», генеральный директор, ORCID 0000-0002-0546-5101.

Толдова Светлана Юрьевна, кандидат филологических наук, Научно-исследовательский университет «Высшая школа экономики», доцент, ORCID 0000-0002-5777-9161.

Хохлова Мария Владимировна, кандидат филологических наук, доцент, Санкт-Петербургский государственный университет, доцент, ORCID 0000-0001-9085-0284.

Авторский указатель

| | | | |
|--------------------|-----|-----------------|-----|
| Benko V. | 11 | Коган М.С. | 29 |
| Dobrov A.V. | 63 | Кораблинов В.О. | 98 |
| Dobrova A.E. | 63 | Кузнецова И.В. | 29 |
| Dzhangolskaya O.V. | 63 | Кулаков К.А. | 118 |
| Rausová K. | 11 | Лебедев А.А. | 118 |
| Smirnova M.O. | 63 | Марусенко Н.М. | 21 |
| Soms N.L. | 63 | Микони С.В. | 109 |
| Zakharova A.A. | 73 | Москин Н.Д. | 118 |
| Абрамов Р.В. | 118 | Рогов А.А. | 118 |
| Бахвалов П.Я. | 82 | Соколова Е.Г. | 44 |
| Гребенников А.О. | 21 | Толдова С.Ю. | 44 |
| Еникеева Е.В. | 54 | Хохлова М.В. | 54 |

Содержание

| | |
|--|-----|
| Предисловие редактора..... | 7 |
| РАЗДЕЛ 1. | |
| КОРПУСНАЯ ЛИНГВИСТИКА | |
| Data-Driven Approach to Identification of Latin Phrases in Russian Web-Crawled Corpora Benko V., Rausová K. | 11 |
| Корпус русского рассказа начала XX века. Пример лингвостатистического анализа Гребенников А.О., Марусенко Н.М. | 21 |
| О возможности использования корпуса NOW в курсе английского для специальных целей для студентов специальности «Биотехнология» Кузнецова И.В., Коган М.С. | 29 |
| К вопросу о формировании набора отношений для корпуса с дискурсивной разметкой текста Соколова Е.Г., Толдова С.Ю. | 44 |
| Методы машинного обучения применительно к задаче выделения глагольных и атрибутивных коллокаций Хохлова М.В., Еникеева Е.В. | 54 |
| РАЗДЕЛ 2. | |
| СЕМАНТИЧЕСКИЙ АНАЛИЗ | |
| Ontological and Formal Grammatical Modeling of Tibetan Nominalized Verb Phrases Smirnova M.O., Dobrov A.V., Dobrova A.E., Soms N.L., Dzhangolskaya O.V. | 63 |
| Syntactic Disambiguation in Constructions with Attachment Ambiguity with Adjuncts by Means of Ontological Semantics Zakharova A.A. | 73 |
| Разработка и реализация методов генерации правил для автоматической проверки правописания Бахвалов П.Я. | 83 |
| Подготовка набора данных для вопросно-ответного поиска по базе знаний. Первый этап: сопоставление сущностей Кораблинов В.О. | 98 |
| Три подхода к определению понятий на основе собственных свойств модели Микони С.В. | 109 |
| Применение деревьев решений для анализа сильных позиций текста в задаче атрибуции произведений Ф. М. Достоевского Рогов А.А., Лебедев А.А., Абрамов Р.В., Москин Н.Д., Кулаков К.А. | 118 |
| Сведения об авторах | 128 |
| Авторский указатель..... | 130 |

Компьютерная лингвистика и вычислительные онтологии. Выпуск 4 (Труды XXIII Международной объединенной научной конференции «Интернет и современное общество», IMS-2020, Санкт-Петербург, 17 – 20 июня 2020 г. Сборник научных трудов). — СПб.: Университет ИТМО, 2020. — 132 с.

Сборник научных трудов

**Компьютерная лингвистика
и вычислительные онтологии**

Выпуск 4

Под редакцией В.П. Захарова
Дизайн обложки С.Н. Ушаков
Оригинал-макет Б.А. Низомутдинов, П.В. Смирнова
Редакционно-издательский отдел Университета ИТМО
Зав. РИО Н.Ф. Гусарова
Подписано к печати 25.11.2020
Заказ № 4364
Тираж 100 экз.

Университет ИТМО. 197101, Санкт-Петербург,
Кронверкский пр., 49 Лит. А.