

Письменный отзыв официального рецензента на диссертационную работу Якунина Кирилла Олеговича на тему «Разработка моделей и методов сбора, анализа и классификации медиа-публикаций на базе методов обработки естественных языков», представленную на соискание степени «доктор филологии» (Ph.D) по специальности 6D070400 – «Вычислительная техника и программное обеспечение»

№ п/п	Критерии	Обоснование позиции официального рецензента
1.	<p>Тема диссертации (на дату утверждения) соответствует направлениям развития науки и/или государственным программам</p>	<p>Диссертация работа выполнялась в рамках проекта программно-целевого финансирования (ППФ) КН МОН РК BR05236839 «Разработка информационных технологий и систем для стимулирования устойчивого развития личности как одна из основ развития цифрового Казахстана» в 2018-2020 годы (Институт информационных и вычислительных технологий КН МОН РК).</p>
2.	<p>Важность для науки</p>	<p>Вынесенные на защиту положения содержат явную научную новизну и были тщательно верифицированы в ходе экспериментов.</p>
3.	<p>Принцип самостоятельности</p>	<p>Личный вклад соискателя диссертации состоит в формулировке и обосновании темы исследования, постановке задач, проведении теоретических и экспериментальных исследований, разработке методов мультикритериальной оценки медиа-источников ММА на основе Байесовской модели агрегации гетерогенных данных, разработке информационной системы мониторинга медиа-пространства Казахстана на базе методов NLP с</p>

		<p>описанием программной архитектуры и основным функционалом, разработке выводов и рекомендаций, что позволяет оценить уровень самостоятельности в проведенных исследованиях, как высокий.</p>
<p>4. Принцип внутреннего единства</p>	<p>4.1 Обоснование актуальности диссертации: 1) Обоснование; 2) Частично обоснована; 3) Не обоснована.</p> <p>4.2 Содержание диссертации отражает тему диссертации: 1) <u>Отражает</u>; 2) Частично отражает; 3) Не отражает</p> <p>4.3 Цель и задачи соответствуют теме диссертации: 1) <u>соответствуют</u>; 2) частично соответствуют; 3) не соответствуют</p>	<p>Задача эффективной векторизации текстовых данных с целью последующей мультикритериальной классификации – одна из важнейших задач в области обработки естественных языков. Современные подходы к векторизации текстов обладают рядом недостатков – они либо слишком тривиальны (bag-of-words) либо наоборот требуют сложнейших моделей глубокого обучения и, следовательно, являются слабо интерпретируемыми (BERT, GPT).</p> <p>Любой новый подход к решению данного комплекса задач вызывает интерес и определяет актуальность настоящего диссертационного исследования. В частности, предлагаемые автором модели и методы, реализованные в рамках распределенной информационно-системы, позволяют решать задачу мультикритериального анализа информационных трендов в медиа-пространстве с небольшим объемом необходимой ручной разметки.</p> <p>Содержание диссертации отражает тему диссертации.</p> <p>Цели и задачи, поставленные в работе, соответствуют теме диссертации и конкретизируют направление исследований.</p>

	<p>4.4 Все разделы и положения диссертации логически взаимосвязаны:</p> <p>1) <u>ПОЛНОСТЬЮ</u> взаимосвязаны;</p> <p>2) взаимосвязь частичная;</p> <p>3) взаимосвязь отсутствует</p>	<p>Полученные в диссертационном исследовании научные результаты обладают внутренним единством, так как все они являются следствием решения одной единой задачи. Оценка внутреннего единства основана на логической связи последовательности элементов разработанного метода, выраженной в составе исследований и изложенной в диссертационной работе. Степень этого единства очень высока, так как каждый полученный результат является следствием предыдущего результата.</p>
<p>4.5 Предложенные автором новые решения (принципы, методы) аргументированы и оценены по сравнению с известными решениями:</p> <p>1) критический анализ есть;</p> <p>2) анализ частичный;</p> <p>3) анализ представляет собой не собственные мнения, а цитаты других авторов</p>	<p>Диссертационная работа является целостным научным исследованием, содержащим системный анализ состояния вопросов в исследуемой области, проработку актуальных направлений и обоснование достигнутых научных результатов. Соискателем раскрыты актуальность, конкретизированы проблемы, связанные с исследуемой темой диссертации.</p> <p>Предложен метод оценки межкурпусного тематического дисбаланса для решения ежедневных в ходе литературного обзора ограничений, в особенности необходимости больших объёмов ручной экспертной разметки и низкой интерпретируемости современных моделей.</p> <p>Разработан мультикритериальный метод оценки медиа-источников ММА. Метод основан на синтезе байесовской модели агрегации критериев, также описанной в данном разделе, на методе анализа иерархий (АИР), позволяющем проводить относительную оценку важности критериев на целевой показатель, а также на применение тех или иных тематических моделей для поиска информационных трендов в корпусе текстов.</p>	

5.	<p>Принцип научной новизны</p> <p>5.1 Научные результаты и положения являются новыми? 1) полностью новые; 2) частично новые (новыми являются 25-75%); 3) не новые (новыми являются менее 25%)</p> <p>5.2 Выводы диссертации являются новыми? 1) <u>полностью новые</u>; 2) частично новые (новыми являются 25-75%); 3) не новые (новыми являются менее 25%)</p> <p>5.3 Технические, технологические, экономические или управленческие решения являются новыми и обоснованными: 1) <u>полностью новые</u>; 2) частично новые (новыми являются 25-75%); 3) не новые (новыми являются менее 25%)</p>	<p>Основные научные результаты и положения включают следующие пункты:</p> <ol style="list-style-type: none"> 1 Исследован вопрос влияния открытого информационного источника на общество, выявлены основные направления влияния, сформирован перечень информативных признаков, на основе которых можно оценить это влияние. 2 Исследованы существующие подходы классификации документов и векторизации текстов, выявлены проблемы и слабые стороны текущих решений, сформированы рекомендации. 3 Разработан подход векторизации текстов на основе тематической модели. 4 Разработан метод оценки межкорпусного тематического дисбаланса, позволяющий автоматически или полуавтоматически получать веса топиков по отношению к заданному признаку. 5 Разработан метод мультикритериальной оценки медиа-источников ММА на базе байесовской модели агрегации. 6 Разработана распределенная информационная система на базе Open Source решений, позволяющая производить сбор (скрапинг), хранение, обработку текстовой информации, а также построение тематических моделей и классификаторов с возможностью визуализации полученных результатов. 7 Собран корпус, состоящий из более чем 6 миллионов публикаций из казахстанских и российских источников, включая как тексты публикаций, так и метаданные.
6.	Обоснованность основных выводов	<p>Все основные выводы основаны/не основаны на весомых с научной точки зрения доказательствах либо достаточно хорошо обоснованы (для qualitative)</p> <p>Полученные и представленные в диссертационной работе научные результаты опираются на строгое математическое изложение</p>

	research и направлений подготовки по искусству и гуманитарным наукам)	основных положений и применяемых методов и подтверждаются результатами проведенного моделирования. Достоверность результатов обеспечивалась использованием современных средств и методик проведения исследований. Это дает основание считать полученные результаты достаточно обоснованными и достоверными.
7. Основные положения, выносимые на защиту	<p>Необходимо ответить на следующие вопросы по каждому положению в отдельности:</p> <p>7.1 Доказано ли положение?</p> <p>1) <u>доказано</u>;</p> <p>2) скорее доказано;</p> <p>3) скорее не доказано;</p> <p>4) не доказано</p> <p>7.2 Является ли тривиальным?</p> <p>1) <u>да</u>;</p> <p>2) <u>нет</u></p> <p>7.3 Является ли новым?</p> <p>1) <u>да</u>;</p> <p>2) <u>нет</u></p> <p>7.4 Уровень для применения:</p> <p>1) <u>узкий</u>;</p> <p>2) <u>средний</u>;</p> <p>3) <u>широкий</u></p> <p>7.5 Доказано ли в статье?</p> <p>1) <u>да</u>;</p> <p>2) <u>нет</u></p>	<p>7.1 Основные положения, выносимые на защиту полностью доказаны вычислительными экспериментами для валидации предложенных моделей и методов. На базе этой системы были проведены вычислительные эксперименты для валидации разработанных моделей и методов.</p> <p>7.2 Основные положения, выносимые на защиту, не являются тривиальными, так как содержат решения, отличающиеся научной новизной и практической значимостью.</p> <p>7.3 Результаты, полученные автором и сформулированные в диссертации, являются новыми научными знаниями в области обработки естественных языков.</p> <p>Предложен метод векторизации текстовых документов с помощью тематической модели BigARTM; данный метод векторизации позволяет использовать данные, полученные из больших неразмеченных текстовых корпусов для получения эффективных векторных представлений, сопоставимых для ряда задач с векторизациями, полученным с помощью моделей глубокого обучения. Предложен метод оценки тематического межкорпусного дисбаланса для самообучения классификационной модели; данный метод</p>

		<p>представляет собой новый подход в области semi-supervised learning, который был верифицирован при решении ряда задачи классификации текстов.</p> <p>Предложена методика многофакторной оценки социальной значимости публикации; данная методика предлагает способ для решения тяжело формализуемой задачи оценки социальной значимости с учётом различных факторов. Методика была апробирована на размеченном корпусе достаточного объёма.</p> <p>Предложена методика многокритериальной оценки масс-медиа ММА на базе байесовской системы агрегации, метода анализа иерархий (АНР) и тематического моделирования; данная методика является новым подходом к решению задачи оценки масс-медиа документов (публикаций) и источников.</p> <p>7.4 Положения выносимые на защиту имеют широкий уровень применения.</p> <p>7.5 Основные положения достигнутых результатов опубликованы в открытой печати, опубликованы в ряде международных и научно-практических конференций, а также опубликованы в рейтинговых журналах (Q1, Q2, Q3).</p>
<p>8. Принципы достоверности Доверенность источников и предоставляемой информации</p>	<p>8.1 Выбор методологии – обоснован или методология достаточно подробно описана</p> <p>1) да; 2) нет</p> <p>8.2 Результаты диссертационной работы получены с использованием современных методов научных исследований и методик обработки и интерпретации данных с применением компьютерных технологий:</p> <p>1) да; 2) нет</p>	<p>8.1 Выбор методологии исследований проведен с учетом особенностей методов обработки естественных языков.</p> <p>8.2 В работе использованы новые подходы к решению комплекса задач. В частности, предлагаемые автором модели и методы, реализованные в рамках распределенной информационный системы, позволяют решать задачу мультикритериального анализа информационных трендов в</p>

	<p>8.3 Теоретические выводы, модели, выявленные взаимосвязи и закономерности доказаны и подтверждены экспериментальным исследованием (для направлений подготовки по педагогическим наукам результаты доказаны на основе педагогического эксперимента):</p> <p>1) да;</p> <p>2) нет</p> <p>8.4 Важные утверждения подтверждены/частично подтверждены/не подтверждены ссылками на актуальную и достоверную научную литературу</p> <p>8.5 Исползованные источники литературы <u>достаточно</u>/не достаточно для литературного обзора</p>	<p>медиапространстве. Разработана распределённая интеллектуальная информационная система Media Analytics, позволяющей проводить сбор, хранение и обработку текстовых данных.</p> <p>8.3 Соискателем в полной мере обоснованы теоретические выводы, методы, разработанные в ходе исследовательской работы.</p> <p>В рамках работы над диссертационным исследованием была разработана распределённая интеллектуальная информационная система Media Analytics, позволяющей проводить сбор, хранение и обработку текстовых данных. В частности, на базе этой системы реализован метод MMA (Мультикритериальный метод оценки медиа-источников), а также проведены вычислительные эксперименты для валидации предложенных моделей и методов. На базе этой системы были проведены вычислительные эксперименты для валидации разработанных моделей и методов. Автором проведена всесторонняя валидация разработанных моделей и методов, было проведено сравнение эффективности предложенных методов с современными моделями классификации.</p> <p>8.4 Важные утверждения, приводимые в диссертационной работе, подтверждаются ссылками на актуальную и достоверную научную литературу.</p> <p>8.5 Список источников, использованных в диссертационной работе, состоит из 92 наименований, которых достаточно для литературного обзора по данной теме.</p>
9.	<p>Принцип практической ценности</p>	<p>9.1 Диссертация имеет теоретическое значение:</p> <p>1) да;</p> <p>2) нет</p> <p>9.1 В научно-исследовательской работе внимание уделено исследованию, которое является актуальным</p>

	<p>9.2 Диссертация имеет практическое значение и существует высокая вероятность применения полученных результатов на практике:</p> <p>1) да; 2) нет</p>	<p>и имеет теоретическую и прикладную значимость. Научные результаты применимы для многокритериальной оценки медиасредства и автоматического поиска информационных трендов.</p>
	<p>9.3 Предложения для практики являются новыми?</p> <p>1) полностью новые; 2) <u>частично новые (новыми являются 25-75%)</u>; 3) не новые (новыми являются менее 25%) 4) низкое.</p>	<p>9.3 Разработанная система решает стандартные для систем медиа-мониторинга задачи, включая сбор текстовых данных, а также поиск, фильтрацию, визуализацию данных. Соискатель предлагает в работе другие варианты практического применения: применение автоматического тематического моделирования, оценку межкорпусного дисбаланса, а также применение предложенных в работе моделей классификации (включая метод ММА). Таким образом, предложения для практики являются в основном новыми.</p>
<p>10. Качество написания и оформления</p>	<p>Качество академического письма:</p> <p>1) <u>высокое</u>; 2) среднее; 3) ниже среднего; 4) низкое.</p>	<p>Диссертация написана грамотным научно-техническим, доступным для чтения, языком. В диссертационной работе имеются незначительные стилистические ошибки, не снижающие качества диссертационной работы.</p>

Заключение. Проведена качественная научно-исследовательская работа, по результатам которой получены актуальные, новые результаты, имеющие теоретическую и прикладную значимость, разработана методика мультифакторной оценки социальной значимости на базе многокритериальной методики оценки медиа ММА с использованием тематической векторизации текстовых документов. Разработана архитектура и программная реализация информационной системы для сбора, обработки и оценки масс-медиа текстов и социальных сетей Media Analytics,

включаяший предложенную методику оценки, а также ее оценку и верификацию. На основании вышеизложенного рекомендую присудить Якунину Кириллу Олеговичу степень доктора философии (Ph.D) специальности 6D070400 – «Вычислительная техника и программное обеспечение».

**Официальный рецензент, к.т.н., ассистент-профессор, декан факультета
Компьютерных технологий и кибербезопасности
АО «Международный Университет Информационных Технологий»**



Сейлова Н.А.

Подпись указанного лица **УАС**
Сейлова Н.А.
14.11.2022