

# АВТОМАТИЗАЦИЯ ОБРАБОТКИ ТЕКСТА

УДК 004.89 : 81'374

Е.В. Котельников, Е.В. Разова, А.В. Котельникова, С.В. Вычегжанин

## Современные словари оценочной лексики для анализа мнений на русском и английском языках\* (аналитический обзор)

*Рассматриваются способы создания словарей оценочной лексики на русском и английском языках с указанием их достоинств и недостатков. Анализируются 13 русскоязычных и 19 англоязычных словарей – приводятся их количественные характеристики и способы создания, вычисляются объединения и пересечения, определяется общая лексика, исследуется распределение по частям речи, указывается доля словосочетаний. Представлены современные области и методы применения словарей оценочной лексики.*

**Ключевые слова:** оценочная лексика, словари оценочной лексики, анализ тональности, анализ мнений

**DOI:** 10.36535/0548-0027-2020-12-3

### ВВЕДЕНИЕ

Анализ мнений (или анализ тональности, *sentiment analysis, opinion mining*) – это область компьютерной лингвистики, в которой исследуются мнения и оценки людей по отношению к различным объектам, таким как продукты, услуги, организации, персоны, события [1]. Под *тональностью* (или *полярностью*) понимают выраженную в тексте субъективность как позитивное или негативное отношение к некоторому объекту [2]. Тональность представляется в виде значения на определенной шкале, которая может быть бинарной (позитивное/негативное отношение), тернарной (добавляется нейтральное или противоречивое), *n*-арной или вещественной (например,  $[-1, 1]$ ).

Анализ мнений в текстах весьма востребован в настоящее время: существует широкий диапазон приложений, например, учет общественного мнения при принятии решений в государственном управлении, организация обратной связи в образовательном процессе, прогнозирование результатов выборов, построение рекомендательных систем, планирование ценообразования и др. [3].

Существуют три основных подхода к анализу мнений в текстах – на основе машинного обучения, на основе словарей и гибридный [2, 4]. В подходе на основе машинного обучения требуются качественно

размеченные обучающие данные и тратится значительное время на процедуру обучения; подход на основе словарей лишен указанных недостатков, но точность анализа часто оказывается недостаточно высокой; в гибридных системах комбинируются два рассмотренных выше подхода.

Ключевым элементом в последних двух подходах являются словари оценочной лексики. Точность анализа тональности в этом случае будет определяться качеством таких словарей. Существует много работ, посвященных созданию словарей оценочной лексики для анализа мнений в текстах [5–7], но проблеме исследования существующих словарей уделяется недостаточно внимания.

Цель настоящей статьи – представить аналитический обзор исследований в области создания, анализа и применения словарей русскоязычной и англоязычной оценочной лексики:

- 1) проанализировать все основные доступные на текущий момент словари оценочной лексики для двух языков – 13 русскоязычных и 19 англоязычных;
- 2) определить состав общей лексики для множества словарей – так называемое ядро оценочной лексики;
- 3) сделать выводы о сходстве и различии русскоязычной и англоязычной оценочной лексики;
- 4) предложить вариант описания автоматических методов создания словарей на основе пространства поиска;

\* Исследование выполнено при финансовой поддержке РФФИ в рамках научного проекта № 19-17-50117

5) сформулировать рекомендации по созданию словарей.

В первом разделе статьи приводятся определение и классификация словарей оценочной лексики; во втором – рассматриваются способы создания таких словарей; в третьем – описываются существующие русскоязычные и англоязычные словари; четвертый посвящен их совместному анализу; в пятом – перечислены области и методы применения словарей для анализа тональности текстов. Заключение содержит основные выводы и рекомендации.

## ОПРЕДЕЛЕНИЕ И КЛАССИФИКАЦИЯ СЛОВАРЕЙ ОЦЕНОЧНОЙ ЛЕКСИКИ

*Словарь оценочной лексики* (тональный словарь, sentiment lexicon, opinion lexicon) – это список оценочных слов и словосочетаний [8, с. 90]. Под *оценочными* понимаются слова и словосочетания, которые передают в тексте позитивное или негативное отношение к каким-либо объектам, например, *хороший, прекрасный, плохой, ужасный*. Понятие «словосочетание» в компьютерной лингвистике достаточно неопределенно [9, 10]. В настоящей работе под *словосочетанием* мы понимаем лингвистическую единицу, которая встречается, когда два и более слов используются совместно для выражения некоторого значения традиционным способом [11, с. 151]. Для краткости изложения, если специально не оговорено иное, под словами понимаем как отдельные слова, так и словосочетания.

Словари оценочной лексики можно классифицировать различными способами:

1) по составу – словари включают только отдельные слова или, также, словосочетания;

2) по шкале тональности – каждому элементу словаря приписывается только знак тональности (позитивная/негативная) или тональность представлена более детальной шкалой (например, действительные числа в диапазоне [-1, +1]);

3) по предметной области – словари могут быть универсальными или предметно-ориентированными;

4) по количеству языков – словари включают оценочную лексику для одного языка или являются многоязычными.

## СПОСОБЫ СОЗДАНИЯ СЛОВАРЕЙ

Существуют три основных способа построения словарей оценочной лексики (рис. 1):

1) ручной – словарь создается, в основном, с помощью разметки слов усилиями людей (аннотаторов);

2) автоматический – словарь строится преимущественно на основе применения различных методов машинной разметки;

3) гибридный – существенную роль при создании словаря играют как ручные методы, так и автоматические.

Это относится, главным образом, к процессу разметки слов – например, при ручном способе создание списка слов-кандидатов, как правило, осуществляется автоматически, а при автоматическом – начальное множество слов для расширения часто формируется разработчиками.

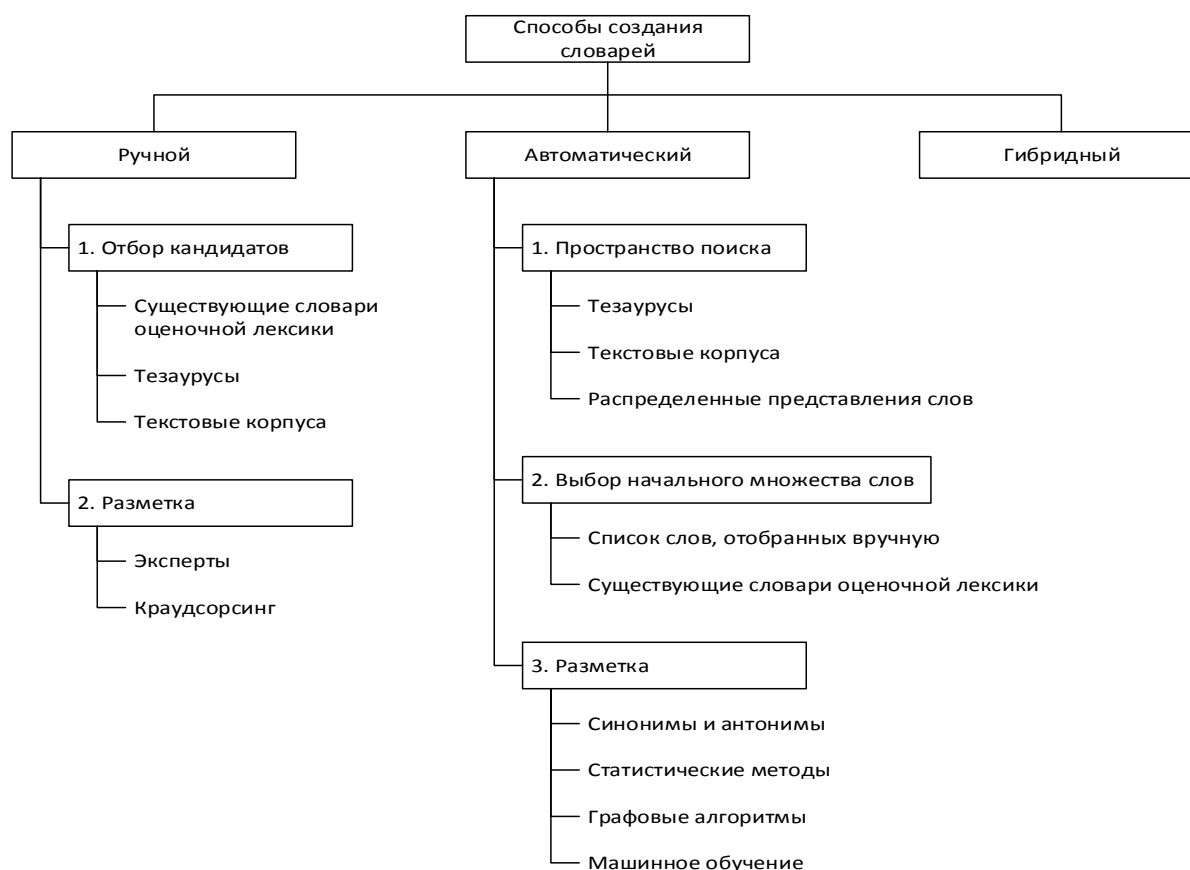


Рис. 1. Способы создания словарей оценочной лексики

## Ручной способ

Построение словаря при *ручном способе*, как правило, выполняется в два этапа:

первый – создание списка слов-кандидатов на вхождение в словарь оценочной лексики;

второй – разметка слов-кандидатов в соответствии с выбранной шкалой тональности.

На первом этапе формируется список слов и словосочетаний, которые, вероятно, являются оценочными. Как правило, такой список строится автоматически, с использованием трех основных методов: 1) на основе существующих словарей оценочной лексики; 2) на основе тезаурусов; 3) с использованием текстовых корпусов.

В первом из указанных методов в качестве кандидатов используются слова из существующих словарей оценочной лексики. Такой метод применялся при создании словаря VADER [12]. Может быть использован перевод существующих словарей на требуемый язык [13, 14]. Во втором методе список слов-кандидатов строится на основе тезаурусов, таких как WordNet для английского языка [15] и RuТез для русского [16]. В третьем методе из текстовых корпусов слова-кандидаты извлекаются на основе различных синтаксических и статистических методов. Данный метод использовался при формировании словарей SemEval-2015 English Twitter Lexicon [17] и SentiRusColl [18].

Наиболее часто применяется комбинация двух или всех трех методов. Например, существующие словари оценочной лексики и тезаурусы использовались при создании словарей MPQA [19]; существующие словари и корпуса применялись для формирования словарей SCL-OPP [20], SCL-NMA [21] и SO-CAL [22]. Все три метода использовались при подготовке словарей EmoLex [23], RuSentiLex [24] и LinisCrowd [14].

Второй этап формирования словарей при ручном способе – разметка отобранных слов-кандидатов. Существуют два основных метода выполнения процедуры разметки: 1) экспертная разметка; 2) разметка с использованием краудсорсинга. В первом методе одному или нескольким экспертам в предметной области (например, лингвистам) предлагается назначить метки всем словам-кандидатам в соответствии с заданной шкалой тональности. Такой метод применялся для создания словарей MPQA [19], SO-CAL [22], RuSentiLex [24] и SentiRusColl [18]. Во втором методе разметка слов осуществляется на основе краудсорсинга (crowdsourcing) – задействования для решения задачи большого количества людей с использованием специальных интернет-платформ, таких как Amazon Mechanical Turk или CrowdFlower [25]. Данный метод использовался при формировании таких словарей как VADER [12], EmoLex [23], SemEval-2015 English Twitter Lexicon [17] и LinisCrowd [14].

Особым приемом в рамках краудсорсинга является геймификация процедуры разметки, которая позволяет повысить заинтересованность и вовлеченность аннотаторов [26, 27].

При ручном способе создания словарей важен контроль качества разметки. В общем случае качество построения словаря (при любом способе формирования) можно оценить на основе его применения

для анализа тональности на заранее размеченных текстах в сравнении с существующими словарями [22, 28]. Также можно оценивать корреляцию разметки между независимыми группами аннотаторов [20]. При экспертной разметке, как правило, оценивается степень согласия между аннотаторами в соответствии с такими статистическими мерами как каппа Коэна (Cohen's  $\kappa$ ), каппа Флейса (Fleiss's  $\kappa$ ), пи Скотта (Scott's  $\Pi$ ) [23]. В случае краудсорсинга контроль включает предварительную оценку знаний аннотаторами языка, оценку их разметки на небольшом контрольном наборе данных, анализ полноты разметки, определение выбросов в разметке [12, 23].

## Автоматический способ

Процесс формирования словарей при автоматическом способе включает, как правило, три этапа (см. рис. 1):

первый – построение пространства поиска, в котором будет осуществляться разметка слов по тональности;

второй – выбор начального множества слов с известной тональностью;

третий – разметка слов в пространстве поиска на основе выбранного начального множества слов (бутстреппинг).

На первом этапе формируется (или используется существующее) пространство поиска, содержащее слова, требующие разметки. В таком пространстве определена метрика расстояния между словами, например, если пространство является векторным, то метрикой может служить косинусное или евклидово расстояние. Если пространство представляет собой связный граф, то расстояние определяется числом ребер в кратчайшем пути между вершинами, а в случае взвешенного графа – суммой весов ребер.

Существуют три основных метода при построении нового или использовании существующего пространства поиска: 1) на основе тезаурусов; 2) на основе корпусов; 3) с использованием распределенных представлений.

В первом из указанных методов формируется семантический граф понятий на основе существующих тезаурусов, таких как WordNet для английского языка [29, 30], RuТез для русского [24] и Wiktionary для множества языков [31, 32]. При этом может быть использован машинный перевод (как правило, с английского) для языков с недостаточной обеспеченностью такими лингвистическими ресурсами [32].

Во втором методе неявное пространство поиска образуется на основе корпуса текстов, снабженных разметкой по тональности. Информация о тональности текстов позволяет использовать такие корпуса для разметки слов на основе статистических методов или машинного обучения [33].

В третьем методе создаются или используются существующие распределенные представления слов, такие как word2vec [34] и GloVe [35]. Распределенные представления слов (word embeddings) – модели, в которых слова представлены в виде вещественных векторов фиксированной, обычно небольшой (несколько сотен), размерности [36, 37]. Строятся такие векторы на основе машинного обучения с использованием статистики совместной встречаемости слов в

неразмеченных текстовых корпусах. Цель обучения заключается в построении таких векторов, близость которых в пространстве распределенных представлений соответствует их семантической близости. В последнее время широкое распространение получили контекстные распределенные представления, такие как ELMo и BERT, в которых вектор слова зависит от текущего контекста [38]. Распределенные представления слов позволяют применять машинное обучение или просто функции расстояния для построения словарей оценочной лексики [5, 39].

На втором этапе выбирается начальное множество оценочных слов с известной тональностью. В качестве такого множества используется либо один из существующих словарей оценочной лексики [32, 40], либо небольшой список слов, отобранных вручную [29, 30, 33].

На третьем этапе осуществляется разметка слов в пространстве поиска с использованием начального множества оценочных слов. Для этого применяются следующие методы:

- расширение начального множества за счет синонимов, антонимов и других семантически связанных слов. Такой метод применим в случае использования пространства поиска на основе тезаурусов [40, 41];
- статистические методы, например, поточечная взаимная информация (pointwise mutual information, PMI), когда оценивается степень совместной встречаемости анализируемого слова и слов из начального множества [33];
- графовые алгоритмы распространения разметки, такие как graph propagation [42], label propagation [43] и random walk [29, 44]. Такие алгоритмы использовались при построении словарей в работах [5, 32];
- машинное обучение – на основе начального множества обучается классификатор, который затем применяется для определения тональности слов в пространстве поиска [29, 30, 45].

### Гибридный способ

В этом способе совместно используются ручные и автоматические методы разметки. Например, при создании словаря Stanford Sentiment Treebank на первом этапе осуществлялась разметка фраз при помощи краудсорсинга, а на втором этапе полученная разметка использовалась для автоматической разметки слов на основе обучения рекурсивных нейронных сетей и с применением распределенных представлений слов [46]. Для формирования словаря ProductSentiRus сначала два аннотатора разметили оценочные слова в предметной области отзывов о фильмах, затем был обучен классификатор, который применялся для автоматической разметки слов в предметных областях отзывов о книгах, играх, фотокамерах и телефонах.

### Достоинства и недостатки способов создания словарей

При использовании *ручного способа* можно выделить следующие достоинства:

- высокая точность – слова, вошедшие в словарь, с высокой вероятностью соответствуют указанной тональности [45];

- при качественном формировании словаря и мощном методе классификации, учитывающем лингвистические особенности, ручной способ может обеспечивать высокое качество анализа тональности [22].

Недостатки *ручного способа*:

- высокая трудоемкость работы по разметке [8];
- зависимость результатов от квалификации, мотивации и качества работы аннотаторов [23];
- низкая полнота [47], в частности, недостаточное количество слов, используемых в социальных медиа [12].

*Автоматический способ* обладает следующими достоинствами:

- высокая полнота (покрытие) – в рамках этого способа можно обеспечить попадание в словарь большого количества оценочных слов языка [45];
- низкая трудоемкость – минимальная ручная обработка [8].

Недостатки *автоматического способа*:

- низкая точность по сравнению с ручным способом – высока вероятность попадания в словарь слов, не являющихся оценочными;
- сложность создания универсального словаря в случае использования методов на основе корпусов;
- сложность создания предметно-ориентированного словаря в случае использования методов на основе тезаурусов;
- зависимость качества словаря от качества аннотированного корпуса при использовании методов на основе корпусов.

В целом, при создании словарей оценочной лексики на основе разных способов существует противоречие между точностью и полнотой (покрытием) [45]: ручной способ обеспечивает высокую точность определения оценочной лексики, но низкую полноту, в то время как автоматический – высокую полноту при низкой точности.

## СЛОВАРИ ОЦЕНОЧНОЙ ЛЕКСИКИ ДЛЯ РУССКОГО И АНГЛИЙСКОГО ЯЗЫКОВ

В настоящее время известно множество словарей оценочной лексики, большая их часть разработана для английского языка, но имеет переводные версии (в том числе, на русский). В настоящем разделе рассматриваются 18 англоязычных словарей, 8 русскоязычных и 3 словаря, имеющих версии как для английского, так и для русского языка.

### Англоязычные словари

1. General Inquirer (URL: <http://www.wjh.harvard.edu/~inquirer>) – один из первых словарей оценочной лексики, созданный в Гарвардском университете в 60-х гг. XX в. [48], содержит более десяти тысяч слов, размеченных вручную в соответствии со множеством синтаксических и семантических категорий, включая тональность. Свободного доступа к словарю нет.

2. LIWC (Linguistic Inquiry and Word Counts – URL: <http://liwc.wpengine.com>) является в настоящее время частью одноименной коммерческой системы анализа текстов [49]. Разметка слов в LIWC сильно

коррелирована с General Inquirer. Свободного доступа к словарю нет.

3. ANEW (Affective Norms for English Words) создан в 1999 г. [50] и содержит разметку по категориям тональности (affective valence), активности и контроля (dominance). Слова размечены вручную студентами-психологами.

4. Словарь Бинга Лью (Bing Liu's Opinion Lexicon или Hu&Liu's Lexicon – URL: <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html>) – результат многолетней работы, начавшейся еще в 2004 г. [41]. Его исходная версия была создана на основе расширения начального списка из 30 прилагательных синонимами и антонимами из тезауруса WordNet. В дальнейшем словарь расширялся, в том числе за счет анализа текстов социальных медиа, поэтому присутствуют слова с ошибками.

5. MPQA (Multi-Perspective Question Answering – Subjectivity Lexicon – URL: [http://mpqa.cs.pitt.edu/lexicons/subj\\_lexicon](http://mpqa.cs.pitt.edu/lexicons/subj_lexicon)) [19] является частью системы анализа мнений OpinionFinder (URL: [http://mpqa.cs.pitt.edu/opinionfinder/opinionfinder\\_2](http://mpqa.cs.pitt.edu/opinionfinder/opinionfinder_2)). Каждое слово в словаре MPQA имеет указание тональности (позитивная, негативная или нейтральная), а также степень тональности (сильная или слабая). При построении MPQA существующий словарь оценочной лексики [51] был расширен за счет тезауруса и словаря General Inquirer, а затем доразмечен вручную.

6. SO-CAL (Semantic Orientation CALCulator – URL: <https://github.com/sfu-discourse-lab/SO-CAL>) – это программный инструмент, определяющий тональность текстов [22]. Словарь, используемый в этом инструменте, в рамках настоящей статьи также обозначается SO-CAL. Он был получен путем экспертной разметки слов-кандидатов, собранных из корпусов отзывов, а также из словаря General Inquirer.

7. SentiWordNet (URL: <https://github.com/aesuli/sentiwordnet>) [29] создан в рамках автоматического способа на основе разметки понятий тезауруса WordNet с использованием машинного обучения и алгоритма случайного блуждания (random walk).

8. AFINN (URL: <https://www2.imm.dtu.dk/pubdb/pubs/6010-full.html>) – словарь (назван по имени разработчика) создавался автором с 2009 г. [52]. Словарь был дополнен нецензурными и сленговыми выражениями с целью получения лучшего результата при автоматическом анализе сообщений в социальных медиа. В настоящей статье используется версия AFINN-111.

9. Sentiment140-Lexicon (URL: <https://github.com/felipebravom/StaticTwitterSent/tree/master/extra/Sentiment140-Lexicon-v0.1>) [33] создан на основе одноименного корпуса, включающего 1,6 млн твитов с позитивными и негативными хештегами. Слова размечались с использованием метода поточечной взаимной информации (PMI).

10. Stanford Sentiment Treebank (URL: <https://nlp.stanford.edu/sentiment>) – это словарь оценочной лексики, сформированный на основе одноименного корпуса, содержащего предложения, извлеченные из отзывов о фильмах. Для каждого предложения построено частичное дерево синтаксического разбора (partial parse trees) [46]. Предложения были разбиты

на фразы, которые размечались по тональности при помощи краудсорсинга. Слова размечались на основе рекурсивных нейронных сетей с использованием деревьев синтаксического разбора и распределенных представлений слов.

11. ML-SentiCon (URL: <https://github.com/mauropelucchi/catalan-referendum/tree/master/ML-SentiCon>) – это автоматически созданные многоязычные (английский, испанский, каталонский, баскский и галисийский) многоуровневые оценочные словари [30]. Каждый словарь содержит 8 уровней, причем каждый вышележащий уровень включает все предыдущие, а также новые элементы. Многоуровневость словарей позволяет выбирать между количеством доступных слов и точностью оценок. Словари строились на основе машинного обучения и алгоритма PolarityRank с использованием WordNet.

12. VADER (Valence Aware Dictionary and sEntiment Reasoner – URL: <https://github.com/cjhutto/vaderSentiment>) – это название словаря и инструмента анализа мнений для социальных медиа на основе правил [12]. Исходный список оценочных слов из существующих словарей оценочной лексики (General Inquirer, LIWC и ANEW) был расширен смайликами, связанными с настроением, акронимами и часто используемым оценочным сленгом. Для разметки слов-кандидатов применялся краудсорсинг.

13. SCL-NMA (Sentiment Composition Lexicon for Negators, Modals, and Degree Adverbs – URL: <http://saifmohammad.com/WebPages/SCL.html#NMA>) [21] – представляет собой список слов и словосочетаний, включающих отрицания, модальные слова и наречия меры и степени. При создании словаря сначала были отобраны слова-кандидаты из General Inquirer, а также высокочастотные фразы из Британского национального корпуса, включающие слова из General Inquirer в комбинации с отрицаниями, модальными словами и наречиями меры и степени, которые затем были размечены при помощи краудсорсинга.

14. SCL-OPP (Sentiment Composition Lexicon for Opposing Polarity Phrases – URL: <http://saifmohammad.com/WebPages/SCL.html#OPP>) представляет собой список отдельных слов и фраз, включающих, по крайней мере, по одному позитивному и негативному слову, например, *счастливый инцидент* [20]. Слова-кандидаты отбирались из корпуса твитов (поэтому словарь содержит хештеги и слова с ошибками) с использованием словарей Бинга Лью, NRC Emotion lexicon (EmoLex), MPQA, ETSL и размечались с помощью краудсорсинга.

15. ETSL (SemEval-2015 English Twitter Sentiment Lexicon – URL: <https://saifmohammad.com/WebPages/SCL.html#ETSL>) – это список униграмм и биграмм с отрицанием [17], который использовался в качестве тестового множества на семинаре SemEval-2015 (задача 10, подзадача E) [53]. В качестве слов-кандидатов были отобраны высокочастотные термины (в том числе с ошибками в написании) из словаря хештегов и словаря Sentiment140-Lexicon, а разметка осуществлялась на основе краудсорсинга.

16. SocialSent (URL: <https://nlp.stanford.edu/projects/socialsent>) – это программный код и наборы данных, в том числе словари оценочной лексики, для проведе-

ния анализа мнений по конкретным предметным областям [5]. Алгоритм создания словарей SentProp включает два этапа: сначала формируется лексический граф на основе распределенных представлений слов, построенных с использованием предметно-ориентированных корпусов. Затем слова в лексическом графе размечаются на основе алгоритма случайного блуждания (random walk) и начального небольшого множества слов, специфичных для предметной области.

В SocialSent содержатся исторические словари оценочной лексики на английском языке. Для каждого десятилетия с 1850 по 2000 гг. построены пара словарей – один включает высокочастотные слова, второй – высокочастотные прилагательные. В настоящей статье использовалась пара словарей за последнее десятилетие.

17. SentiWords (URL: <https://hlt-nlp.fbk.eu/technologies/sentiwords>) создан в процессе исследования, каким образом на основе словаря SentiWordNet можно получить априорную оценку тональности слова, т.е. оценку, которая не зависит от различных семантических значений данного слова [45]. С этой целью было использовано машинное обучение и различные признаки, выводимые из характеристик слов SentiWordNet.

18. WordStat (URL: <https://provalisresearch.com/products/content-analysis-software/wordstat-dictionary/sentiment-dictionaries>) создан путем объединения отрицательных и положительных слов из словарей Harvard-IV, Regressive Imagery Dictionary и LIWC [54]. Затем список был автоматически расширен синонимами и связанными словами, а также различными морфологическими формами входящих в него слов.

## Русскоязычные словари

1. ProductSentiRus (URL: <http://panchenko.me/data/snlp/sentiment/ProductSentiRus.txt>) разработали Илья Четвёркин и Наталья Лукашевич для предметной области товаров (фильмы, книги, игры, цифровые фотокамеры, мобильные телефоны) [55]. Сначала было использовано машинное обучение для построения модели на основе множества вручную размеченных слов из отзывов о фильмах и набора статистических и лингвистических признаков слов. Затем построенная модель применялась для классификации оценочной лексики в других областях. В результате общий словарь оценочной лексики для всех пяти областей включает 5 000 слов. Слова в словаре отсортированы по мере убывания вероятности их оценочности, но не разделены на позитивные и негативные.

2. Словарь Блинова (URL: <https://github.com/kotelnikov-ev/BlinovSentimentLexicon>). Павел Блинов с соавторами [31] сформировали вручную список из 969 наиболее позитивных и 1 138 наиболее негативных слов из словаря ProductSentiRus, а затем автоматически расширили список синонимами и антонимами из русского Викисловаря (<https://ru.wiktionary.org/wiki>).

3. LinisCrowd (URL: <http://www.linis-crowd.org>). Олеся Кольцова с соавторами [14] создавали свой словарь при помощи краудсорсинга. Сначала они отобрали 7 546 слов на основе списка высокочастотных прилагательных, словаря ProductSentiRus, толкового словаря и перевода англоязычного словаря оце-

ночной лексики SentiStrength [56]. Затем не менее трех аннотаторов каждому слову присваивали оценки от -2 до +2. В настоящем исследовании позитивными и негативными считаются такие слова, которые получили большинство оценок соответствующей тональности.

4. Словарь Котельникова (URL: <https://github.com/kotelnikov-ev/KotelnikovSentimentLexicons>). Евгений Котельников с соавторами [57] сначала автоматически отобрали по 10 000 слов-кандидатов для каждой из пяти предметных областей (отзывы о ресторанах, автомобилях, фильмах, книгах и камерах), четыре аннотатора оценили каждое слово как позитивное, негативное, нейтральное или противоречивое, затем было создано два объединенных по предметным областям словаря: в первый вошли слова, относительно тональности которых были согласны три аннотатора из четырех (*Kotelnikov\_large*), во второй – слова, относительно тональности которых согласны все аннотаторы (*Kotelnikov\_small*).

5. Словарь Тутубалиной (URL: <https://www.ispras.ru/dcouncil/docs/diss/2016/tutubalina/dissertaciya-tutubalina.pdf>). Елена Тутубалина в своей диссертации [58] создала вручную словарь на основе строго позитивных и негативных отзывов пользователей об автомобилях (в отзывах учитывались только разделы преимуществ и недостатков). Словарь был расширен за счет добавления синонимов.

6. RuSentiLex (URL: <https://www.labinform.ru/pub/rusentilex>). Наталья Лукашевич и Анатолий Левчик [24] создали словарь RuSentiLex, в котором для каждого слова указывается тональность (позитивная, негативная, нейтральная) и источник (мнение, факт, чувство). Сначала были сгенерированы списки оценочных слов на основе тезауруса RuТез, существующих словарей оценочной лексики, новостных статей и Twitter, затем лингвисты анализировали полученные списки для формирования итогового словаря. Словарь содержит более десяти тысяч слов и выражений, имеющих оценку тональности.

В нашей работе использовалась версия словаря 2017 г. и только слова и сочетания с позитивной или негативной тональностью. Исследовались две версии – *RuSentiLex\_large*, включающая все позитивные и негативные элементы, и *RuSentiLex\_small*, включающая позитивные и негативные элементы, для которых источником является мнение.

7. SentiRusColl (URL: <https://github.com/kotelnikov-ev/SentiRusColl>) [18] – это оценочный словарь словосочетаний. Для его создания использовался корпус отзывов по десяти предметным областям (книги, фильмы, музыка, автомобили, компьютеры, бытовая техника, телефоны, банки, отели, рестораны), из него автоматически отбирались словосочетания-кандидаты, которые далее размечались тремя аннотаторами. В словарь внесены словосочетания, получившие большинство голосов. Для нашего исследования использовался вариант словаря SentiRusColl со стоп-словами, содержащий предлоги и союзы.

8. Карта слов (<https://kartaslov.ru>) – это онлайн-карта слов и выражений русского языка [59]. Оценочный словарь, разработанный в рамках этого проекта, содержит слова русского языка, снабжённые

метками тональности (позитивная, негативная или нейтральная) и силы эмоционально-оценочного заряда. При создании словаря использовался краудсорсинг: в процессе разметки пользователю предлагалось оценить то или иное слово как нейтральное, позитивное, негативное или ответить «не знаю». В настоящем исследовании используется версия данного словаря от ноября 2019 г. ([https://github.com/dkulagin/kartaslov/tree/master/dataset/emo\\_dict](https://github.com/dkulagin/kartaslov/tree/master/dataset/emo_dict)), которая содержит только отдельные слова.

## Многоязычные словари

1. Словарь Чена-Скиены (Chen-Skiena's lexicon – URL: <https://sites.google.com/site/datascienceslab/projects/multilingualsemiment>). Яньцин Чен (Yanqing Chen) и Стивен Скиена (Steven S. Skiena) [32] автоматически построили словари оценочной лексики для 136 языков (включая русский и английский). На основе тезаурусов WordNet и Wiktionary и с помощью машинного перевода (Google Translate) был построен многоязычный семантический граф, связывающий слова на разных языках. Затем, отталкиваясь от англоязычного словаря Бинга Лью, были сформированы оценочные словари для других языков на основе алгоритма распространения разметки. Варианты для английского и русского языков далее называются *Chen-Skiena\_en* и *Chen-Skiena\_ru*.

2. EmoLex (NRC Emotion Lexicon – URL: <https://saifmohammad.com/WebPages/NRC-Emotion-Lexicon.htm>) – словарь составлен Саифом Мохаммадом (Saif Mohammad) и Питером Тёрни (Peter Turney) с помощью краудсорсинга [23]. Для формирования множества слов-кандидатов использовались тезаурус Macquarie, а также словари General Inquirer и WordNet Affect. Словарь EmoLex содержит слова, соотношенные с позитивной и негативной тональностью, а также с эмоциями «гнев», «предвкушение», «отвращение», «страх», «радость», «грусть», «удивление» и «доверие». Для исследования были отобраны слова из версии словаря NRC-Emotion-Lexicon-Wordlevel-v0.92, имеющие позитивную или негативную тональность, в том числе словосочетания. Вариант для английского языка далее называется *Emolex\_en*.

В ноябре 2017 г. словарь был переведен на более чем 100 языков (в том числе, на русский) с помощью Google Translate. Для русского языка были отобраны слова и словосочетания, имеющие позитивную или негативную тональность. Вариант для русского языка далее называется *Emolex\_ru*.

3. SenticNet (<https://sentic.net/downloads>) – это проект по анализу мнений на уровне понятий, начатый в 2009 г. в MIT Media Laboratory. Для построения словаря оценочной лексики сначала был сформирован граф понятий с использованием нейронных сетей долгой краткосрочной памяти (Long short-term memory, LSTM), обучавшихся на основе распределенных представлений слов [39]. Затем для разметки понятий по тональности применялось специальное векторное пространство AffectiveSpace [60]. В нашей статье используется версия словаря SenticNet 5 (далее в статье этот словарь называется *SenticNet\_en*). В рамках SenticNet имеется проект BabelSenticNet [61],

содержащий словари для 40 языков, в том числе русского. Далее в нашей статье словарь для русского языка обозначается *SenticNet\_ru*.

## Характеристики словарей

С учетом отсутствия открытого доступа к словарям General Inquirer и LIWC, а также двух версий словарей RuSentiLex и Котельникова, в нашей работе далее исследовались 19 словарей для английского языка и 13 словарей для русского языка. В табл. 1 и 2 приведены характеристики рассмотренных словарей после следующей предобработки: все элементы были преобразованы к нижнему регистру, в словарях были оставлены только элементы, состоящие из букв (латинского алфавита для английских словарей и кириллицы для русских), знака дефиса и пробела. В этих таблицах содержится информация о размерах множеств позитивных и негативных слов; размере объединений и пересечений этих множеств; способе, использованном для создания словаря; шкале тональности (указывается, каким образом шкала делилась для получения позитивных и негативных значений тональности); диапазоне количества слов в элементах словаря; годе создания или последней модификации.

Объем словарей для английского языка варьируется сильнее, чем для русского: от 765 (AFINN) до 523 092 (Sentiment140-Lexicon) слов. Для русскоязычных словарей объем меняется от 1 115 (Котельников\_small) до 24 765 (SenticNet\_ru) слов.

Для десяти словарей английского языка и семи словарей русского языка множества позитивных и негативных слов имеют непустое пересечение. Видимо, в словарях Блинова и Тутубалиной это произошло из-за автоматического расширения исходных списков слов за счет синонимов, а в словаре EmoLex\_ru – вследствие автоматического перевода. В словаре Котельникова одно и то же слово может быть позитивным для одной области и негативным для другой, например, *непредсказуемый сюжет* – *непредсказуемые отказы*. В словаре RuSentiLex одинаковые слова могут иметь разный смысл и тональность, что указывается в ссылке на статью тезауруса RuТез, например, *легкий (покладистый)* – *легкий (поверхностный)*. В английских словарях пересечение связано с многозначностью слов или автоматическим режимом формирования словарей.

Для русскоязычных словарей среднее количество позитивных оценочных слов составляет 44%, негативных – 56%: негативная лексика более разнообразна. Для англоязычных словарей среднее количество позитивных оценочных слов составляет 58%, негативных – 42%: позитивная лексика более разнообразна. Однако если из рассмотрения исключить три самых больших англоязычных словаря, сформированных автоматически (SenticNet\_en, Sentiment140-Lexicon и Sentiment Treebank), то соотношение позитивной и негативной лексики в английских словарях становится в точности таким же, как и в русских словарях (44% – позитивная и 56% – негативная).

Интересно, что русскоязычные словари, как правило, строятся на основе ручного способа (9 из 13), а для англоязычных словарей способы распределились поровну.

Русскоязычные словари предпочитают бинарную шкалу тональности (10 из 13), в то время как в англоязычных бинарная шкала используется только в 5 из 19 словарей.

Шесть англоязычных словарей и восемь русскоязычных содержат только отдельные слова, остальные словари включают словосочетания. Наиболее длинные словосочетания в английских словарях (до 38 слов) имеются в словаре Sentiment Treebank, в русских словарях (до 8 слов) – в словаре EmoLex\_ru.

Из дальнейшего рассмотрения были исключены три самых больших англоязычных словаря, сформированных автоматически (SenticNet\_en, Sentiment140-Lexicon

и Sentiment Treebank), а также самый большой русскоязычный словарь SenticNet\_ru, поскольку они содержат слишком много некорректных элементов вследствие их автоматического формирования. Например, в Sentiment140-Lexicon входят такие словосочетания, как *they landed*, *describe what*, а в SenticNet\_ru – *амортизационная стойка*, *адвокатское сословие полтенца*. Кроме того, в дальнейшем при рассмотрении пересечений и объединений словарей мы не анализируем словарь ProductSentiRus, поскольку для него невозможно выполнить разделение на позитивные и негативные слова. Таким образом, далее рассматривается 16 англоязычных словарей и 11 русскоязычных.

Таблица 1

Характеристики англоязычных словарей оценочной лексики

№ п/п	Словарь	Поз. эл.	Нег. эл.	Объед.	Перес.	Способ	Шкала	Кол-во слов в элементах	Год	Ссылка
1	ANEW	419	346	765	0	Ручной	Непрерывная шкала [1; 9]: [1, 4] – нег., [6, 9] – поз.	1	1999	[50]
2	Словарь Бинга Лью	2 005	4 774	6 776	3	Автом.	Бинарная шкала {-1, +1}	1	2004	[41]
3	MPQA	2 304	4 152	6 450	6	Ручной	Бинарная шкала {-1, +1}	1	2005	[19]
4	SO-CAL	2 446	3 566	6 004	8	Ручной	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-3	2007	[22]
5	SentiWordNet	16 436	18 244	32 902	1778	Автом.	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-9	2010	[29]
6	EmoLex_en	2 312	3 324	5 555	81	Ручной	Бинарная шкала {-1, +1}	1	2011	[23]
7	AFINN	878	1 596	2 474	0	Ручной	Дискретная шкала [-5, 5]: [-5, -1] – нег., [1, 5] – поз.	1-3	2012	[52]
8	Sentiment140-Lexicon	328 188	194 904	523 092	0	Автом.	Непрерывная шкала [-6, 8]: [-6, 0] – нег., (0, 8] – поз.	1-2	2013	[33]
9	Sentiment Treebank	33 201	25 790	58 980	11	Гибрид.	Непрерывная шкала [0, 1]: [0; 0,4] – нег., (0,6; 1] – поз.	1-38	2013	[46]
10	ML-SentiCon	12 774	12 134	24 818	90	Автом.	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-8	2014	[30]
11	VADER	3 187	4 034	7 221	0	Ручной	Непрерывная шкала [-4, 4]: [-4, 0] – нег., (0, 4] – поз.	1-2	2014	[12]
12	Chen-Skiena_en	1 421	2 955	4 376	0	Автом.	Бинарная шкала {-1, +1}	1	2014	[32]
13	SCL-NMA	1 607	1 575	3 182	0	Ручной	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-4	2016	[21]
14	SCL-OPP	513	640	1 153	0	Ручной	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-3	2016	[20]
15	ETSL	654	496	1 150	0	Ручной	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-2	2016	[17]
16	SocialSent	1 255	1 213	2 463	5	Автом.	Непрерывная шкала [-3,9; 2,76]: [-3,9; -0,5] – нег., [0,5; 2,76] – поз.	1	2016	[5]
17	SenticNet_en	55 311	44 689	100 000	0	Автом.	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-5	2018	[39]
18	SentiWords	18 280	21 599	39 663	216	Автом.	Непрерывная шкала [-1, 1]: [-1, 0] – нег., (0, 1] – поз.	1-9	2018	[45]
19	WordStat	5 492	10 486	15 955	23	Автом.	Бинарная шкала {-1, +1}	1-4	2018	[54]
	<b>В среднем</b>	<b>25 720</b>	<b>18 764</b>	<b>44 367</b>	<b>117</b>					

## Характеристики русскоязычных словарей оценочной лексики

№ п/п	Словарь	Поз. эл.	Нег. эл.	Объед.	Перес.	Способ	Шкала тональности	Кол-во слов в элементах	Год	Ссылка
1	ProductSentiRus			5 000	0	Гибрид.	Нет разделения	1	2012	[55]
2	Словарь Блинова	1 864	2 145	3 839	170	Автом.	Бинарная шкала $\{-1, +1\}$	1	2013	[31]
3	Chen-Skiena_ru	1 246	1 630	2 876	0	Автом.	Бинарная шкала $\{-1, +1\}$	1	2014	[32]
4	LinisCrowd	566	1 940	2 506	0	Ручной	Дискретная шкала $[-2, 2]$ : $[-2, -1]$ – нег., $[1, 2]$ – поз.	1	2016	[14]
5	Котельников_large	1 046	2 211	3 247	10	Ручной	Бинарная шкала $\{-1, +1\}$	1	2016	[57]
6	Котельников_small	389	727	1 115	1	Ручной	Бинарная шкала $\{-1, +1\}$	1	2016	[57]
7	Словарь Тутубалиной	1 077	1 458	2 508	27	Ручной	Бинарная шкала $\{-1, +1\}$	1	2016	[58]
8	EmoLex_ru	2 085	2 805	4 750	140	Ручной	Бинарная шкала $\{-1, +1\}$	1–8	2017	[23]
9	RuSentiLex_large	3 433	9 485	12 784	134	Ручной	Бинарная шкала $\{-1, +1\}$	1–6	2017	[24]
10	RuSentiLex_small	2 686	5 719	8 331	74	Ручной	Бинарная шкала $\{-1, +1\}$	1–6	2017	[24]
11	SenticNet_ru	14 650	10 115	24 765	0	Автом.	Непрерывная шкала $[-1, 1]$ : $[-1, 0]$ – нег., $(0, 1]$ – поз.	1–6	2018	[39]
12	SentiRusColl	4 008	2 569	6 577	0	Ручной	Бинарная шкала $\{-1, +1\}$	2–7	2019	[18]
13	Карта слов	4 550	6 774	11 324	0	Ручной	Бинарная шкала $\{-1, +1\}$	1	2019	[59]
	<b>В среднем</b>	<b>3 133</b>	<b>3 965</b>	<b>6 894</b>	<b>43</b>					

## АНАЛИЗ СЛОВАРЕЙ

## Объединения и пересечения словарей

**Англоязычные словари.** Объединение всех 16 английских словарей включает 64 597 элементов: 27 980 позитивных (43,3%)<sup>1</sup>, 34 271 негативных (53,1%), а также 2 346 элементов, у которых не удалось однозначно определить тональность (3,6%). Пересечение всех словарей (обозначим его *En\_Intersection16*) содержит 4 слова: 1 позитивное (*pretty*) и 3 негативных (*hell, hurt, sick*). В множество общих слов для всех словарей входит по одному прилагательному, существительному, глаголу и наречию. В это множество не попали слова *good* и *bad*<sup>2</sup>.

Множество слов, которые встречаются хотя бы в 15 словарях из 16 – это объединение всех пересечений из 15 словарей (*En\_Intersection15*) (табл. 3). Это множество включает 41 слово: 15 позитивных и 26 негативных (22 прилагательных, 11 существительных, 5 глаголов 2 наречия и 1 междометие).

*En\_Intersection14* включает уже 122 слова: 51 позитивное и 71 негативное (72 прилагательных, 37 существительных, 9 глаголов, 2 междометия и 2 наречия).

Постепенное уменьшение *N* при объединении всех пересечений из *N* словарей позволяет формировать так называемое *ядро оценочной лексики*, т.е. та-

кое множество слов, относительно которых согласны все или почти все словари<sup>3</sup>.

На рис. 2 в каждой ячейке приведено отношение мощности пересечения словаря, указанного в строке, и словаря, указанного в столбце, к объему словаря в строке (в процентах) в виде тепловой карты. Например, для пары словарей (ANEW, Словарь Бинга Лью) значение 49,2% указывает на то, что такая доля лексики словаря ANEW присутствует в словаре Бинга Лью.

Словарь SentiWords базировался на словаре SentiWordNet, поэтому SentiWordNet полностью входит в SentiWords, а в словаре Chen-Skiena\_en оценочные слова были взяты из словаря Бинга Лью, поэтому Chen-Skiena\_en полностью входит в данный словарь.

Для словарей SentiWords и SentiWordNet средняя доля вхождений в них других словарей является наивысшей (64,2% и 54,5% соответственно). Это объясняется тем, что данные словари имеют самый большой объем (39 663 и 32 902 слова). Минимальную долю вхождений показывают словари наименьшего размера ANEW, SCL-OPP и ETSL.

Словари SCL-NMA, SCL-OPP, ETSL и VADER имеют наименьшее вхождение в самый большой словарь SentiWords (от 37,0% до 45,9%), так как указанные словари ориентированы на анализ социальных медиа и включают соответствующую специфическую лексику.

В целом совпадение лексики между английскими словарями оказывается относительно невысоким (30,0%).

<sup>1</sup> Тональность слова определялась голосованием по большинству словарей, в которые оно входит.

<sup>2</sup> Слова *good* и *bad* отсутствуют в словаре SocialSent, слово *bad* отсутствует в ANEW. Кроме того, слово *bad* входит в словарь ML-SentiCon как позитивное.

<sup>3</sup> Другой подход к формированию ядра предложен в работе [62], в которой ядро определялось на основе анализа областей концентрации оценочной лексики в пространстве распределенных представлений слов.

Пересечение 15 из 16 английских словарей (*En\_Intersection15*)

Позитивные слова	Негативные слова
<i>beautiful, beauty, cute, elegant, good, happy, love, lucky, nice, pretty, respect, safe, smile, sweet, wonderful</i>	<i>abuse, badly, bastard, blind, bloody, bother, cancer, damn, danger, death, dirty, disaster, dreadful, hate, hell, hopeless, hurt, lonely, lost, mad, pain, sad, sick, stupid, ugly, wrong</i>

Таблица 4

Пересечение 10 из 11 русских словарей (*Ru\_Intersection10*)

Позитивные слова	Негативные слова
<i>благоприятный, великолепный, волшебный, достойный, замечательный, красивый, красота, легендарный, превосходный, преимущество, прекрасный, привлекательный, приятный, роскошный, удобный, ценный, чудесный, энергичный, эффективный, яркий</i>	<i>бессмысленный, глупый, грязный, неприятный, неудачный, опасный, оскорбительный, трудный, тупой, ужасный</i>

Словари	ANEW	Словарь Бинга Лью	MPQA	SO-CAL	SentiWordNet	EmoLex_en	AFINN	ML-SentiCon	VADER	Chen-Skienna_en	SCL-NMA	SCL-OPP	ETSL	SocialSent	SentiWords	WordStat
ANEW	100	49,2	51,5	48,9	66,1	59,5	36,5	47,2	46,3	47,8	31,2	14,5	15,4	38,2	74,0	60,8
Словарь Бинга Лью	5,5	100	79,8	50,5	58,1	35,1	19,4	44,5	32,7	64,6	15,3	3,7	4,3	12,2	70,6	73,1
MPQA	6,1	83,8	100	57,2	65,9	39,5	19,0	49,1	32,3	55,9	19,9	4,4	4,6	13,9	79,6	74,8
SO-CAL	6,2	57,0	61,5	100	63,8	38,0	17,6	47,6	29,5	44,2	19,5	3,9	4,6	14,8	76,2	64,9
SentiWordNet	1,5	12,0	12,9	11,6	100	10,6	4,1	37,9	7,9	8,7	3,8	1,2	1,4	4,1	100	20,0
EmoLex_en	8,2	42,8	45,8	41,0	62,6	100	17,7	44,6	26,2	39,6	17,6	5,4	4,7	14,6	73,9	61,6
AFINN	11,3	53,1	49,5	42,6	54,5	39,7	100	43,1	98,0	48,7	23,2	8,1	10,3	22,7	64,3	68,8
ML-SentiCon	1,5	12,1	12,8	11,5	50,2	10,0	4,3	100	8,4	8,7	3,5	1,0	1,1	3,6	63,2	20,3
VADER	4,9	30,6	28,9	24,5	35,9	20,1	33,6	29,0	100	24,9	10,7	3,5	4,5	10,2	43,7	43,4
Chen-Skienna_en	8,4	100	82,3	60,6	65,2	50,2	27,6	49,2	41,1	100	22,3	5,7	6,4	18,7	77,8	82,2
SCL-NMA	7,5	32,7	40,4	36,7	39,7	30,8	18,0	27,6	24,2	30,7	100	6,3	6,4	16,8	44,7	43,8
SCL-OPP	9,6	21,9	24,8	20,5	34,6	26,1	17,4	21,0	22,2	21,6	17,3	100	16,0	21,9	37,0	28,1
ETSL	10,3	25,6	25,7	23,8	41,5	22,9	22,2	24,6	28,4	24,4	17,7	16,1	100	24,6	45,9	33,1
SocialSent	11,9	33,6	36,4	36,0	54,9	33,0	22,8	35,8	29,8	33,3	21,6	10,2	11,5	100	62,0	43,3
SentiWords	1,4	12,1	13,0	11,5	83,0	10,3	4,0	39,5	8,0	8,6	3,6	1,1	1,3	3,9	100	20,4
WordStat	2,9	31,0	30,3	24,4	41,3	21,5	10,7	31,6	19,6	22,6	8,7	2,0	2,4	6,7	50,6	100

Рис. 2. Отношение попарных пересечений словарей к размеру словарей в строке (%)

Словари	Словарь Блинова	Chen-Skiena_ru	LinisCrowd	Котельников_large	Котельников_small	Словарь Тутубалиной	EmoLex_ru	RuSentiLex_large	RuSentiLex_small	SentiRusColl	Карта слов
Словарь Блинова	100	17,3	21,5	32,3	14,8	29,1	21,0	43,5	36,2	0,0	38,3
Chen-Skiena_ru	23,1	100	16,5	15,8	6,7	11,7	32,9	35,3	24,3	0,0	39,4
LinisCrowd	33,0	19,0	100	31,3	15,6	26,1	31,6	74,3	56,2	0,0	60,2
Котельников_large	38,2	14,0	24,2	100	34,3	23,4	19,8	48,5	39,9	0,0	46,8
Котельников_small	50,9	17,2	35,2	100	100	34,4	26,0	59,6	50,9	0,0	53,4
Словарь Тутубалиной	44,5	13,4	26,1	30,3	15,3	100	24,1	52,7	48,0	0,0	41,6
EmoLex_ru	17,0	19,9	16,7	13,5	6,1	12,7	100	34,9	24,8	0,0	35,5
RuSentiLex_large	13,1	7,9	14,6	12,3	5,2	10,3	13,0	100	65,2	0,1	36,5
RuSentiLex_small	16,7	8,4	16,9	15,5	6,8	14,5	14,2	100	100	0,1	38,6
SentiRusColl	0,0	0,0	0,0	0,0	0,0	0,0	0,0	0,1	0,1	100	0,0
Карта слов	13,0	10,0	13,3	13,4	5,3	9,2	14,9	41,2	28,4	0,0	100

Рис 3. Отношение попарных пересечений словарей к размеру словарей в строке (%)

**Русскоязычные словари.** Объединение всех 11 русских словарей включает 32 902 элемента: 13 240 позитивных (40,3%)<sup>4</sup>, 19 323 негативных (58,7%), а также 339 элементов, у которых не удалось однозначно определить тональность (1,0%). Пересечение всех 11 словарей (обозначим его *Ru\_Intersection11*) пусто. Множество слов, которые встречаются хотя бы в 10 словарях из 11, – это объединение всех пересечений из 10 словарей (*Ru\_Intersection10*) (табл. 4). Это множество включает 30 слов: 20 позитивных (18 прилагательных и 2 существительных) и 10 негативных (10 прилагательных). Интересно, что в это множество не попали слова *хороший* и *плохой*<sup>5</sup>.

*Ru\_Intersection9*, т. е. множество слов, которые встречаются хотя бы в 9 словарях из 11, – это объе-

динение всех пересечений из 9 словарей, включает уже 134 слова: 71 позитивное (65 прилагательных и 6 существительных) и 63 негативных (51 прилагательное, 10 существительных, 1 наречие, 1 глагол). Следует отметить преобладание прилагательных для обоих языков, но в разной степени: например, в *En\_Intersection15* прилагательных 53,7%, в то время как в *Ru\_Intersection10* – 93,3%.

На рис. 3 в каждой ячейке приведено отношение мощности пересечения словаря, указанного в строке, и словаря, указанного в столбце, к размеру словаря в строке (в процентах) в виде тепловой карты.

Словари *RuSentiLex\_small* и *Котельников\_small* являются частью, соответственно, словарей *RuSentiLex\_large* и *Котельников\_large*.

Так же, как для английского языка, для самых крупных русскоязычных словарей *RuSentiLex\_large* и *Карта слов* средняя доля вхождений в них других словарей является наивысшей (49,0% и 39,0% соответственно). Наименьшую долю вхождений показывает *SentiRusColl*, так как он содержит только словосочетания.

Словари *Chen-Skiena\_ru* и *EmoLex\_ru* имеют наименьшее вхождение в самый большой словарь *RuSentiLex*, что объясняется формированием этих словарей при помощи машинного перевода.

<sup>4</sup> Тональность слова определялась голосованием по большинству словарей, в которые оно входит.

<sup>5</sup> *Хороший* отсутствует в словаре *EmoLex\_ru* (слово *good* было переведено как *хорошо*) и в словаре *SentiRusColl* (этот словарь содержит лишь словосочетания), *плохой* – в словаре *Котельников\_small* (один из аннотаторов отнес *плохой* к нейтральным словам для нескольких областей), *EmoLex\_ru* (слово *bad* было переведено как *плохо*) и в словаре *SentiRusColl* (этот словарь содержит лишь словосочетания). Кроме того, слово *плохой* входит в словарь *Блинова* как позитивное.

В целом, совпадения лексики между русскими словарями (без учета словарей-подмножеств RuSentiLex и словаря Котельникова) оказывается в среднем ниже (20,5%), чем между английскими словарями (30,0%).

### Части речи

**Англоязычные словари.** Части речи для английских словарей были получены при помощи библиотеки Stanza от Stanford NLP Group [63].

В табл. 5 показано распределение частей речи в отдельных словарях и в объединенном словаре. Во втором столбце приведен размер словаря, в третьем и четвертом – количество элементов словаря, являющихся отдельными словами, и их доля в словаре, далее доля существительных, глаголов, прилагательных и наречий среди элементов словаря, состоящих из одного слова.

В 14 из 16 словарей наблюдается преобладание существительных, лишь в SO-CAL и Chen-Skienna\_en количество прилагательных хоть и незначительно, но превышает количество существительных. В объединенном словаре 58,0% слов – это существительные. В большинстве словарей, в том числе и в объединенном словаре, количество прилагательных значительно больше количества глаголов, исключение составляют словари AFINN и ETSL. Интересно отметить довольно высокую долю наречий в MPQA, SO-CAL и словаре Бинга Лью.

По сравнению со словарями *En Intersection16*, *En Intersection15* и *En Intersection14*, значительную часть которых составляют прилагательные, их доля значительно уменьшилась. Таким образом, ядро оценочной лексики составляют прилагательные, но по мере расширения словаря существительные и глаголы начинают преобладать.

Словосочетания составляют 23,0% (14 862) от объема объединенного словаря. Из них 81,9% являются биграмами, 15,4% – триграммами. Наиболее частотными комбинациями по сочетанию частей речи являются существительное + существительное (30,6%, например, *air alert, animal disease, food poisoning*) и прилагательное + существительное (27,4%, например, *cool stuff, right direction, accidental injury*). Доли других комбинаций не превышают 4%.

**Русскоязычные словари.** В табл. 6 приведено распределение частей речи в отдельных словарях и в объединенном словаре; также показано соотношение отдельных слов и словосочетаний. Части речи были получены при помощи морфологического парсера *ru-morphu2* [64]. Заметим, что в табл. 6 отсутствует информация о словаре SentiRusColl, поскольку он содержит элементы, включающие не менее двух слов.

Некоторые словари предпочитают одни части речи, например, в словаре Тутубалиной, словаре Блинова и в словаре Котельникова много прилагательных; в переводных словарях (EmoLex\_ru и словаре Chen-Skienna\_ru), а также в RuSentiLex высока доля существительных. В словаре LinisCrowd относительно сбалансировано соотношение существительных и прилагательных, а в словаре Карта слов – соотношение существительных и глаголов. Можно отметить высокую долю наречий в словаре Котельникова.

В объединенном словаре преобладают существительные (41,2%); их доля значительно ниже, чем в англоязычных словарях (58,0%). Интересно, что доля прилагательных почти одинакова – 22,8% в русском и 23,0% в английском словарях. Глаголов в русскоязычных словарях значительно больше, чем в англоязычных – 29,0% против 11,8%.

Таблица 5

Распределение частей речи (%)

Словари	Объем словаря	Отдельные слова		Части речи отдельных слов				
		Кол-во	%	Сущ.	Глаг.	Прил.	Нареч.	Другие
ANEW	765	765	100,0%	66,5%	10,1%	22,6%	0,5%	0,3%
Словарь Бинга Лью	6 776	6 776	100,0%	37,9%	15,9%	34,0%	11,7%	0,5%
MPQA	6 450	6 450	100,0%	38,0%	13,5%	34,7%	13,2%	0,6%
SO-CAL	6 004	5 920	98,6%	34,4%	11,6%	39,4%	14,0%	0,6%
SentiWordNet	32 902	25 178	76,5%	54,7%	9,1%	26,5%	9,1%	0,6%
EmoLex en	5 555	5 555	100,0%	56,2%	14,9%	26,8%	1,7%	0,4%
AFINN	2 474	2 460	99,4%	35,1%	31,2%	30,7%	1,8%	1,2%
ML-SentiCon	24 818	18 874	76,0%	53,7%	8,9%	32,1%	4,7%	0,6%
VADER	7 221	7 217	99,9%	53,6%	16,6%	21,1%	7,4%	1,3%
Chen-Skienna_en	4 376	4 376	100,0%	37,6%	17,2%	38,5%	6,2%	0,5%
SCL-NMA	3 182	1 611	50,6%	43,8%	20,9%	33,3%	1,3%	0,7%
SCL-OPP	1 153	589	51,1%	46,9%	20,2%	24,6%	6,3%	2,0%
ETSL	1 150	902	78,4%	41,1%	24,1%	21,8%	3,7%	9,3%
SocialSent	2 463	2 463	100,0%	47,6%	13,1%	35,3%	2,5%	1,5%
SentiWords	39 663	30 784	77,6%	55,2%	8,5%	27,9%	8,0%	0,4%
WordStat	15 955	15 702	98,4%	51,2%	21,1%	21,4%	5,7%	0,6%
<b>Объединенный словарь</b>	64 597	49 735	77,0%	58,0%	11,8%	23,0%	6,4%	0,8%

Распределение частей речи (%)

Словари	Объем словаря	Отдельные слова		Части речи отдельных слов				
		Кол-во	%	Сущ.	Глаг.	Прил.	Нареч.	Другие
Словарь Блинова	3 839	3 839	100,0%	18,5%	19,5%	47,0%	8,2%	6,8%
Chen-Skienna_ru	2 876	2 876	100,0%	48,8%	20,6%	19,6%	6,3%	4,7%
LinisCrowd	2 506	2 506	100,0%	38,2%	18,1%	42,8%	0,4%	0,5%
Котельников_large	3 247	3 247	100,0%	26,4%	25,0%	35,0%	10,7%	2,9%
Котельников_small	1 115	1 115	100,0%	26,4%	11,8%	44,8%	13,5%	3,5%
Словарь Тутубалиной	2 508	2 508	100,0%	9,3%	2,8%	78,5%	1,6%	7,8%
EmoLex_ru	4 750	4 486	94,4%	53,1%	13,0%	25,9%	1,7%	6,3%
RuSentiLex_large	12 784	10 543	82,5%	47,9%	24,1%	26,3%	0,4%	1,3%
RuSentiLex_small	8 331	7 151	85,8%	48,1%	18,3%	31,6%	0,5%	1,5%
Карта слов	11 324	11 324	100,0%	40,9%	39,7%	17,3%	1,7%	0,4%
<b>Объединенный словарь</b>	32 902	23 836	72,4%	41,2%	29,0%	22,8%	3,5%	3,5%

Так же, как в англоязычных словарях, по мере расширения ядра оценочной лексики существительные начинают преобладать над прилагательными.

Словосочетания составляют 27,6% (9 066) от размера объединенного словаря (близко к англоязычным словарям – 23,0%). Из них 56,2% являются биграмами, 30,5% – триграммами, 10,3% – квадрограммами (доля триграмм и квадрограмм значительно выше, чем в английском языке). Наиболее частотными комбинациями по сочетанию частей речи являются прилагательное + существительное (20,4%, например, *канительное дело, огромный респект, экономный расход*) и глагол + существительное (7,2%, например, *испытать судьбу, подкладывать свинью, пожалеть время*).

## ПРИМЕНЕНИЕ СЛОВАРЕЙ ОЦЕНОЧНОЙ ЛЕКСИКИ

### Области применения

Словари оценочной лексики применяются для анализа тональности текстов в различных предметных областях – в онлайн-торговле, в рекомендательных системах, в экономике, в политических исследованиях, в медицине и образовании.

В онлайн-торговле анализ тональности используется для изучения мнений потребителей с целью выявления преимуществ и недостатков предлагаемых на рынке товаров и услуг в интересах производителей и для помощи в выборе другим потребителям, для объяснения и предсказания продаж. Например, с помощью словарей оценочной лексики анализируются мнения относительно мобильных телефонов [65], фотокамер [66], электроники [67], книг [68], кухонных приборов [28, 67], услуг провайдеров облачных сервисов [69], услуг энергетической компании [70].

Анализ тональности применяется в рекомендательных системах для выработки рекомендаций, со-

ответствующих интересам пользователей. Анализируются мнения пользователей о фильмах, банках, ресторанах, отелях [18, 28, 71].

В экономике, анализируя тональность текстовых сообщений с использованием словарей оценочной лексики, предсказывают направление движения стоимости акций на фондовом рынке [72]. В политических исследованиях анализ тональности используется для определения отношения избирателей к политическим партиям или кандидатам на выборах в президенты и законодательное собрание с целью предсказания результатов голосования [73]. В медицине с помощью словарей оценочной лексики анализируются мнения пациентов о врачах и лекарствах [74], о качестве обслуживания в медицинских учреждениях [75]. В сфере образования сообщения, полученные в качестве обратной связи от студентов, анализируются с целью оценки образовательных курсов и выступлений лекторов в учреждениях высшего образования [76].

### Методы анализа мнений с использованием словарей

Как указывалось во Введении, существуют три основных подхода к анализу мнений в текстах – на основе машинного обучения, на основе словарей и гибридный. Словари оценочной лексики используются в последних двух подходах.

**Методы на основе словарей.** В работе [73] словарный подход используется для анализа тональности твитов, посвященных выборам в законодательное собрание в Индии. Применяются словари оценочной лексики AFINN, Бинга Лью, EmoLex, Syuzhet, SentiWordNet, SenticNet, VADER. В качестве признаков рассматриваются N-граммы и смайлики. В результате экспериментов наименьшее значение средней абсолютной ошибки было получено с использованием словаря VADER.

В статье [76] предлагается система анализа полученных в качестве обратной связи от студентов сообщений OMFeedback, оценивающая выступления лекторов в учреждениях высшего образования. Для анализа тональности сообщений используется метод на основе словаря VADER.

Способ формирования словаря для определенной предметной области, разработанный в [75], применяется для решения задачи классификации по тональности отзывов пациентов о качестве обслуживания в медицинских учреждениях. Сформированный словарь позволяет получить более высокое качество классификации, чем словари VADER и AFINN.

В работе [70] анализируются твиты потребителей услуг энергетической компании. Классификация текстов по тональности осуществляется с использованием инструмента Sentimentr и словаря Бинга Лью. Авторы отмечают, что каждый из инструментов наиболее хорошо выделяет тексты определенных классов тональности. В предлагаемом методе сначала используется Sentimentr для выделения негативных твитов, а затем оставшиеся твиты классифицируются с помощью словаря Бинга Лью на позитивные и нейтральные. Такой подход позволяет повысить качество классификации текстов, принадлежащих специфической предметной области, по сравнению с использованием каждого словаря по отдельности.

Авторы работы [74] применяют словарь испаноязычной оценочной лексики iSOL, полученный машинным переводом словаря Бинга Лью, для классификации по тональности мнений о лекарствах из корпуса DOS и мнений пациентов о врачах.

В работе [66] предложен метод вычисления весов оценочных слов на основе генетического алгоритма и совместной кластеризации слов и документов. Для оценки качества анализа тональности используются текстовые коллекции семинара РОМИП-2011, содержащие отзывы о фильмах, книгах и фотокамерах. По результатам экспериментов предложенный метод позволил получить более высокие значения F1-меры по сравнению с SVM и словарным методом.

В статье [18] разработан способ создания универсального словаря оценочной лексики SentiRusColl из словосочетаний. Сформированный словарь показывает в среднем более высокое качество классификации текстов по тональности применительно к десяти предметным областям по сравнению со словарем RuSentiLex.

В целом, на основе анализа предыдущих работ, сложно сделать вывод о преимуществе того или иного англоязычного словаря оценочной лексики: отсутствуют исследования, в которых сравнивались бы основные существующие словари на базе единого подхода для одних и тех же текстовых корпусов. Результаты, приводимые в статьях, сильно зависят от используемого метода анализа тональности и предметной области. Тем не менее, на основе обзора публикаций, проведенного в рамках настоящего исследования, можно сделать следующие предварительные выводы относительно англоязычных словарей:

- практически во всех статьях используется словарь SentiWordNet;

- SentiWordNet показывает результаты, сопоставимые со словарем MPQA;
- словари Бинга Лью и VADER демонстрируют преимущество по точности над словарями SentiWordNet и MPQA.

Повторим, что эти выводы не окончательные и требуют уточнения на базе масштабного сравнительного анализа.

Для русскоязычных словарей проводилось значительно меньше исследований. Только в 2018 г. был выполнен сравнительный анализ существующих словарей оценочной лексики на основе применения машинного обучения для разных предметных областей [28], который показал преимущество словаря ProductSentiRus. Однако в сравнении не участвовал словарь Карта слов, а используемый метод машинного обучения мог не выявить всех преимуществ словарей при анализе тональности.

**Гибридные методы.** Авторы статьи [65] предложили метод анализа тональности, основанный на использовании словаря SentiWordNet и машинном обучении – методе k ближайших соседей, наивном байесовском классификаторе и случайном лесе. Для тестирования метода используется корпус отзывов о мобильных телефонах с сайта Amazon. Наилучшее значение точности было получено с использованием случайного леса.

В работе [69] разработан гибридный метод, основанный на словарном подходе и методах нечеткой логики. Эксперименты проводятся с использованием корпуса, содержащего мнения пользователей об услугах поставщиков облачных сервисов.

В статье [77] предложен метод расширения словаря оценочной лексики для аспектно-ориентированного анализа тональности. Авторы сравнили два словаря, составленных вручную, и словарь, созданный при помощи предложенного метода. В качестве классификатора применялся метод максимальной энтропии. Эксперименты проводились с использованием корпуса отзывов о ресторанах с семинара SentiRuEval-2015 и продемонстрировали преимущество предложенного метода.

В статье [78] предлагается метод анализа тональности текстов TextJSM, основанный на ДСМ-методе интеллектуального анализа данных и использующий словарь оценочной лексики из работы [57]. Предложенный метод показывает преимущество над традиционными методами машинного обучения.

В целом, в работах, посвященных гибридным методам анализа тональности, такие методы показывают более высокие результаты, чем словарный подход или машинное обучение по отдельности.

## ЗАКЛЮЧЕНИЕ

Проведенный анализ русскоязычных и англоязычных словарей оценочной лексики позволяет сделать следующие выводы.

1. Спектр англоязычных словарей шире, их размеры больше (даже после исключения автоматически построенных словарей с большим количеством ошибок объединенный англоязычный словарь в два раза превышает русскоязычный), в них чаще встре-

чаются словосочетания, а используемые шкалы тональности детальнее, чем у русскоязычных словарей.

2. Словари для обоих языков имеют ряд совпадающих характеристик: часто встречается ситуация, когда одно и то же слово рассматривается и как позитивное, и как негативное; негативная лексика более разнообразна, причем соотношение количества позитивных и негативных слов совпадает (44% и 56%); ядро оценочной лексики составляют прилагательные, а по мере расширения словаря начинают преобладать существительные; в словарях в среднем мало совпадающей лексики (30,0% для англоязычных и 20,5% для русскоязычных).

3. Невозможно однозначно говорить о преимуществе того или иного словаря оценочной лексики: отсутствуют масштабные исследования, в которых сравнивались бы основные существующие словари на базе единого подхода для одних и тех же текстовых корпусов.

Настоящий обзор позволяет сформулировать рекомендации для исследователей, имеющих потребность в словаре оценочной лексики. Такая потребность может быть удовлетворена либо за счет использования существующих словарей, либо с помощью разработки нового словаря. В первом случае рекомендуется использовать композиционный способ, т. е. формировать новый словарь на основе голосования существующих. При этом возможно контролировать баланс точности и полноты оценочной лексики – для высокой точности следует увеличивать порог количества словарей, проголосовавших за включение слова в новый словарь, а при снижении этого порога повышается полнота.

Если существующие словари не удовлетворяют исследователя (например, требуется словарь на другом языке или предметно-ориентированный словарь), то способ создания словаря будет зависеть от цели и имеющихся ресурсов:

- предметно-ориентированный словарь рекомендуется создавать с применением автоматического метода на основе корпусов;
- при наличии достаточных временных и/или финансовых ресурсов рекомендуется осуществить разметку словаря вручную;
- машинный перевод существующих словарей на другой язык, как правило, недостаточен для обеспечения высокой точности; требуется разметка или проверка вручную.

Таким образом, процесс исследования и разработки словарей оценочной лексики, в частности, русскоязычных, представляется незавершенным и требует дальнейших усилий научного сообщества. Одним из наиболее актуальных направлений является проведение масштабного исследования существующих словарей оценочной лексики на базе единого подхода, учитывающего преимущества словарных методов, в том числе с использованием композиции словарей.

## СПИСОК ЛИТЕРАТУРЫ

1. Liu B. *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. – Cambridge: Cambridge University Press, 2015.
2. Taboada M. *Sentiment Analysis: An Overview from Linguistics // Annual Review of Linguistics*. – 2016. – Vol. 2. – P. 325–347.
3. Yue L., Chen W., Li X., Zuo W., Yin M. *A survey of sentiment analysis in social media // Knowledge and Information Systems*. – 2018. – P. 1–47.
4. Poria S., Hazarika D., Majumder N., Mihalcea R. *Beneath the Tip of the Iceberg: Current Challenges and New Directions in Sentiment Analysis Research // Computing Research Repository*. – 2020. – arXiv: 2005.00357.
5. Hamilton W.L., Clark K., Leskovec J., Jurafsky D. *Inducing domain-specific sentiment lexicons from unlabeled corpora // Proceedings of Conference on Empirical Methods in Natural Language Processing*. – 2016. – P. 595–605.
6. Vo D.T., Zhang Y. *Don't count, predict! An automatic approach to learning sentiment lexicons for short text // Proceedings of 54th Annual Meeting of the Association for Computational Linguistics*. – 2016. – P. 219–224.
7. Wang L., Xia R. *Sentiment Lexicon Construction with Representation Learning Based on Hierarchical Sentiment Supervision // Proceedings of Conference on Empirical Methods in Natural Language Processing*. – 2017. – P. 502–510.
8. Liu B. *Sentiment analysis and opinion mining // Synthesis Lectures on Human Language Technologies*. – 2012. – Vol. 5(1). – P. 1–167.
9. Боярский К.К., Каневский Е.А. *Семантика устойчивых словосочетаний с глаголами // Научно-техническая информация. Сер. 2*. – 2019. – № 11. – С. 23–31.
10. *Multiword Units in Machine Translation and Translation Technology / eds. R. Mitkov, J. Monti, G.C. Pastor, V. Seretan*. – Amsterdam: John Benjamins Publishing Company, 2018.
11. Manning C.D., Schütze H. *Foundations of Statistical Natural Language Processing*. – The MIT Press, 1999. – 620 p.
12. Hutto C.J., Gilbert E. *VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text // Eighth International Conference on Weblogs and Social Media (ICWSM-14)*. Ann Arbor, MI, June 2014. – Palo Alto: The AAAI Press, 2014.
13. Abdaoui A., Azé J., Bringay S., Poncelet P. *FEEL: a French Expanded Emotion Lexicon // Language Resources & Evaluation*. – 2017. – Vol. 51(3). – P. 833–855.
14. Koltsova O.Yu., Alexeeva S.V., Kolcov S.N. *An Opinion Word Lexicon and a Training Dataset for Russian Sentiment Analysis of Social Media // Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference "Dialogue-2016"*. – 2016. – Vol. 15(22). – P. 277–287.
15. WordNet. *An electronic lexical database / ed. C. Fellbaum*. – Cambridge, MA: MIT Press; 1998.
16. Лукашевич Н.В. *Тезаурусы в задачах информационного поиска*. – М.: Изд-во МГУ, 2011.

17. Kiritchenko S., Zhu X., Mohammad S. Sentiment Analysis of Short Informal Texts // *Journal of Artificial Intelligence Research*. – 2014. – Vol. 50. – P. 723–762.
18. Kotelnikova A.V., Kotelnikov E.V. SentiRusColl: Russian Collocation Lexicon for Sentiment Analysis // *Artificial Intelligence and Natural Language Conference (AINL)*. Communications in Computer and Information Science (November 20–22, 2019, Tartu, Estonia). – Cham: Springer, 2019. – Vol. 1119. – P. 18–32.
19. Wilson T., Wiebe J., Hoffmann P. Recognizing contextual polarity in phrase-level sentiment analysis // *Proceedings of the 2005 Human Language Technology Conference and the Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP-05)*. – 2005. – P. 347–354.
20. Kiritchenko S., Mohammad S.M. Happy Accident: A Sentiment Composition Lexicon for Opposing Polarities Phrases // *Proceedings of the 10th edition of the Language Resources and Evaluation Conference (LREC)*. – Portorož, Slovenia, 2016. – P. 1157–1164.
21. Kiritchenko S., Mohammad S.M. The Effect of Negators, Modals, and Degree Adverbs on Sentiment Composition // *Proceedings of the 7th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis (WASSA)*. – San Diego, California, 2016. – P. 43–52.
22. Taboada M., Brooke J., Tofiloski M., Voll K., Stede M. Lexicon-based methods for sentiment analysis // *Computational Linguistics*. – 2011. – Vol. 37(2). – P. 267–307.
23. Mohammad S.M., Turney D.P. Crowdsourcing a word-emotion association lexicon // *Computational Intelligence*. – 2013. – Vol. 29(3). – P. 436–465.
24. Loukachevitch N., Levchik A. Creating a General Russian Sentiment Lexicon // *Proceedings of Language Resources and Evaluation Conference LREC-2016*. – 2016. – P. 1171–1176.
25. Bhatti S.S., Gao X., Chen G. General framework, opportunities and challenges for crowdsourcing techniques: A Comprehensive survey // *The Journal of Systems and Software*. – 2020. – Vol. 167.
26. Hong Y., Kwak H., Baek Y. Tower of babel: a crowdsourcing game building sentiment lexicons for resource-scarce languages // *Proceedings of the WWW 2013 companion*. – Rio de Janeiro, Brazil, 13–17 May 2013. – New York: Association for Computing Machinery, 2013. – P. 549–556.
27. Thisone C.C., Ghasemi A., Faltings B. Sentiment analysis using a novel human computation game // *Proceedings of the 3rd workshop on the people’s web meets NLP, Jeju Island, Republic of Korea, 8–14 July 2012*. – P. 1–9.
28. Kotelnikov E.V., Peskischeva T.A., Kotelnikova A.V., Razova E.V. A comparative study of publicly available Russian sentiment lexicons // *7th conference on Artificial Intelligence and Natural Language (AINL-2018)*. Communications in Computer and Information Science. – Cham: Springer, 2018. – Vol. 930. – P. 139–151.
29. Baccianella S., Esuli A., Sebastiani F. SentiWordNet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining // *Proceedings of the Seventh Conference on International Language Resources and Evaluation (LREC10)*. – 2010. – P. 2200–2204.
30. Cruz F.L., Troyano J.A., Pontes B., Ortega F.J. Building layered, multilingual sentiment lexicons at synset and lemma levels // *Expert Systems with Applications*. – 2014. – Vol. 41. – P. 5984–5994.
31. Blinov P.D., Klekovkina M.V., Kotelnikov E.V., Pestov O.A. Research of lexical approach and machine learning methods for sentiment analysis // *Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference “Dialogue-2013”*. – 2013. – Vol. 12(19). – P. 51–61.
32. Chen Y., Skiena S. Building Sentiment Lexicons for All Major Languages // *Proceedings of the 52<sup>nd</sup> Annual Meeting of the Association for Computational Linguistics*. – Baltimore, 2014. – P. 383–389.
33. Mohammad S.M., Kiritchenko S., Zhu X. NRC-Canada: Building the State-of-the-Art in Sentiment Analysis of Tweets // *Proceedings of the seventh international workshop on Semantic Evaluation – SemEval-2013 (June 2013, Atlanta, USA)*. – Madison: Omnipress, Inc., 2013. – P. 321–327.
34. Mikolov T., Chen K., Corrado G., Dean J. Distributed Representations of Words and Phrases and their Compositionality // *Proceedings of Conference on Neural Information Processing Systems*. – 2013. – P. 3111–3119.
35. Pennington J., Socher R., Manning C.D. GloVe: Global Vectors for Word Representation // *Proceedings of Conference on Empirical Methods in Natural Language Processing*. – 2014. – P. 1532–1543.
36. Almeida F., Xexeo G. Word Embeddings: A Survey // *Computing Research Repository*. – 2019. – arXiv:1901.09069.
37. Çano E., Morisio M. Word Embeddings for Sentiment Analysis: A Comprehensive Empirical Survey // *Computing Research Repository*. – 2019. – arXiv:1902.00753.
38. Liu Q., Kusner M.J., Blunsom P. A Survey on Contextual Embeddings // *Computing Research Repository*. – 2020. – arXiv:2003.07278v.
39. Cambria E., Poria S., Hazarika D., Kwok K. SenticNet 5: Discovering conceptual primitives for sentiment analysis by means of context embeddings // *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18)*. – 2018. – P. 1795–1802.
40. Loughran T., McDonald B. When is a liability not a liability? Textual Analysis, Dictionaries and 10-Ks // *The Journal of Finance*. – 2011. – Vol. 66(1). – P. 35–66.
41. Hu M., Liu B. Mining and Summarizing Customer Reviews // *Proceedings of the ACM SIGKDD International Conference on Knowledge, Discovery and Data Mining – KDD-2004 (Aug 22-25, 2004, Seattle, Washington, USA)*. – New York: Association for Computing Machinery, 2004. – P. 168–177.

42. Velikovich L., Blair-Goldensohn S., Hannan K., McDonald R. The viability of web-derived polarity lexicons // *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*. – 2010. – P. 777–785.
43. Zhu X., Ghahramani Z. Learning from labeled and unlabeled data with label propagation. Technical report, Technical Report CMUCALD-02-107. – Carnegie Mellon University, 2002.
44. Hassan A., Radev D.R. Identifying Text Polarity Using Random Walks // *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*. – 2010. – P. 395–403.
45. Gatti L., Guerini M., Turchi M. SentiWords: Deriving a high precision and high coverage lexicon for sentiment analysis // *IEEE Transactions on Affective Computing*. – 2016. – Vol. 7(4). – P. 409–421.
46. Socher R., Perelygin A., Wu J., Chuang J., Manning C., Ng A., Potts C. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank // *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP 2013)*. – 2013. – P. 1631–1642.
47. Pang B., Lee L., Vaithyanathan S. Thumbs up? Sentiment Classification using Machine Learning Techniques // *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*. – 2002. – P. 79–86.
48. Stone P.J., Dunphy D.C., Smith M.S., Ogilvie D.M. *The General Inquirer: A Computer Approach to Content Analysis*. – Cambridge, MA: MIT Press, 1966.
49. Pennebaker J.W., Boyd R.L., Jordan K., Blackburn K. *The development and psychometric properties of LIWC2015*. – Austin, TX: University of Texas at Austin, 2015.
50. Bradley M.M., Lang P.J. *Affective Norms for English Words (ANEW): Stimuli, instruction manual, and affective ratings (Tech. Report C-1)*. – Gainesville: University of Florida, Center for Research in Psychophysiology, 1999.
51. Riloff E., Wiebe J. Learning Extraction Patterns for Subjective Expressions // *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing (EMNLP-03)*. – Stroudsburg: Association for Computational Linguistics, 2003. – P. 105–112.
52. Nielsen F. A new ANEW: Evaluation of a word list for sentiment analysis in microblogs // *Proceedings of the ESWC2011 Workshop on Making Sense of Microposts: Big things come in small packages, Heraklion*. – 2012. – P. 93–98.
53. Rosenthal S., Nakov P., Kiritchenko S., Mohammad S., Ritter A., Stoyanov V. SemEval-2015 Task 10: Sentiment Analysis in Twitter // *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*. – 2015. – P. 451–463.
54. WordStat: content analysis and text mining software. – URL: <https://provalisresearch.com/products/content-analysis-software/worldstat-dictionary/sentiment-dictionaries> (дата обращения: 01.08.2020).
55. Chetviorkin I., Loukachevitch N. Extraction of Russian Sentiment Lexicon for Product Meta-Domain // *Proceedings of COLING 2012*. – Mumbai, 2012. – P. 593–610.
56. Thelwall M., Buckley K., Paltoglou G., Cai D., Kappas A.A. Sentiment strength detection in short informal text // *Journal of the American Society for Information Science and Technology*. – 2010. – Vol. 61(12). – P. 2544–2558.
57. Kotelnikov E., Bushmeleva N., Razova E., Peskischeva T., Pletneva M. Manually Created Sentiment Lexicons: Research and Development // *Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference “Dialogue-2016”*. – 2016. – Vol. 15(22). – P. 300–314.
58. Тугубалина Е.В. Методы извлечения и резюмирования критических отзывов пользователей о продукции: дис. ... канд. физ.-мат. наук. – М.: ИСП РАН, 2016. – 145 с.
59. Кулагин Д.И. Карта слов: переосмысление подхода к составлению онлайн-словарей в постмобильную эру // *Международная конференция «Диалог 2017» – Компьютерная лингвистика и интеллектуальные технологии (Москва, 31 мая – 3 июня 2017 г.)*. – URL: <http://www.dialog-21.ru/media/3974/kulagindi.pdf> (дата обращения: 01.08.2020).
60. Cambria E., Fu J., Bisio F., Poria S. AffectiveSpace 2: Enabling Affective Intuition for Concept-Level Sentiment Analysis // *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*. – 2015. – P. 508–514.
61. Vilares D., Peng H., Satapathy R., Cambria E. BabelSenticNet: A Commonsense Reasoning Framework for Multilingual Sentiment Analysis // *Proceedings of IEEE Symposium Series on Computational Intelligence*. – 2018. – P. 1292–1298.
62. Razova E.V., Kotelnikov E.V. Concentration Areas of Sentiment Lexica in the Word Embedding Space // *International Journal of Cognitive Informatics and Natural Intelligence*. – 2019. – Vol. 13(2). – P. 48–62.
63. Qi P., Zhang Y., Zhang Y., Bolton J., Manning C.D. Stanza: A Python Natural Language Processing Toolkit for Many Human Languages // *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, 2020*. – Stroudsburg: Association for Computational Linguistics, 2020.
64. Korobov M. Morphological Analyzer and Generator for Russian and Ukrainian Languages // *Proceedings of 3rd Conference on Analysis of Images, Social Networks and Texts (AIST)*. – 2015. – P. 320–332.
65. Hosel C., Roschke C., Thomanek R., Ritter M. Lexicon-Based Sentiment Analysis of Online Customer Ratings as a Quinary Classification Problem // *Communications in Computer and Information Science*. – 2019. – Vol. 1034. – P. 75–80.
66. Kotelnikov E.V., Pletneva M.V. Text Sentiment Classification based on Genetic Algorithm and Word and Document Co-clustering // *Journal of*

- Computer and Systems Sciences International. – 2016. – Vol. 55(1). – P. 106–114.
67. Han H., Zhang Y., Zhang J., Yang J., Zou X. Improving the performance of lexicon-based review sentiment analysis method by reducing additional introduced sentiment bias // PLOS ONE. – 2018. – Vol. 13(8). – P. 1–11.
  68. Khatun F., Chowdhury S., Tumpa Z., Rabby S., Hossain S., Abujar S. Sentiment Analysis of Amazon Book Review Data Using Lexicon Based Analysis // Advances in Intelligent Systems and Computing. – 2019. – Vol. 1108. – P.1303–1309.
  69. Alharbi J.R., Alhalabi W.S. Hybrid Approach for Sentiment Analysis of Twitter Posts Using a Dictionary-based Approach and Fuzzy Logic Methods: Study Case on Cloud Service Providers // International Journal on Semantic Web and Information Systems. – 2020. – Vol. 16(1). – P. 116–145.
  70. Ikoro V., Sharmina M., Malik K., Batista-Navarro R. Analyzing Sentiments Expressed on Twitter by UK Energy Company Consumers // 2018 Fifth International Conference on Social Networks Analysis, Management and Security (SNAMS). – 2018. – P. 95–98.
  71. Iqbal F., Maqbool J., Fung B., Batool R., Khattak A., Aleem S., Hung P. A Hybrid Framework for Sentiment Analysis using Genetic Algorithm based Feature Reduction // IEEE Access. – 2019. – Vol. 7. – P. 14637–14652.
  72. Vo D.T., Zhang Y. Don't count, predict! An automatic approach to learning sentiment lexicons for short text // Proceedings of 54th Annual Meeting of the Association for Computational Linguistics. – 2016. – P. 219–224.
  73. Bansal B., Srivastava S. Lexicon-based Twitter sentiment analysis for vote share prediction using emoji and N-gram features // International Journal of Web Based Communities. –2019. – Vol. 15(1). – P. 85–99.
  74. Jiménez-Zafra S.M., Martín-Valdivia M.T., Molina-González M.D., Ureña-López L.A. How do we talk about doctors and drugs? Sentiment analysis in forums expressing opinions for medical domain // Artificial Intelligence in Medicine. – 2019. – Vol. 93. – P. 50–57.
  75. Kumar C.S.P., Babu L.D.D. Evolving dictionary based sentiment scoring framework for patient authored text // Evolutionary Intelligence. – 2020.
  76. Wook M., Razali N., Ramli S., Wahab N., Hasbullah N., Zainudin N., Talib M. Opinion mining technique for developing student feedback analysis system using lexicon-based approach (OMFeedback) // Education and Information Technologies. – 2020. – Vol. 25. – P. 2549–2560.
  77. Tutubalina E., Nikolenko S. Constructing Aspect-Based Sentiment Lexicons with Topic Modeling // Proceedings of 5th Conference on Analysis of Images, Social Networks and Text. –2017. – P. 208–220.
  78. Котельников Е.В. Метод анализа тональности текстов TextJSM // Научно-техническая информация. Сер. 2. – 2018. – № 2. – С. 8–20.

*Материал поступил в редакцию 07.08.20.*

#### **Сведения об авторах**

**КОТЕЛЬНИКОВ Евгений Вячеславович** – доктор технических наук, доцент, профессор кафедры прикладной математики и информатики, Вятский государственный университет.  
e-mail: kotelnikov.ev@gmail.com

**РАЗОВА Елена Владимировна** – кандидат педагогических наук, доцент, доцент кафедры прикладной математики и информатики, Вятский государственный университет.  
e-mail: razova.ev@gmail.com

**КОТЕЛЬНИКОВА Анастасия Валерьевна** – кандидат педагогических наук, доцент кафедры прикладной математики и информатики, Вятский государственный университет.  
e-mail: kotelnikova.av@gmail.com

**ВЫЧЕГЖАНИН Сергей Владимирович** – инженер кафедры прикладной математики и информатики, Вятский государственный университет.  
e-mail: vychegzhaninsv@gmail.com