

Типы социально-сетевого дискурса в автоматической классификации текстов по тональности

Н.Е. Маслова

Московский государственный лингвистический университет

natalia.maslova277@gmail.com

Аннотация

В работе описывается опыт классификации тональности текстов социально-сетевого дискурса (ССД) экологической тематики. В рамках подхода машинного обучения с учителем используется модель нейронных сетей Doc2vec. Подтверждается эффективность применения метода нейронных сетей для обработки текстов на русском языке.

При классификации используется система параметров ССД, созданная Р.К. Потаповой. Вводятся понятия личностной и стратифицирующей деприваций, и подтверждается практическая ценность этого разделения.

В статье показано, как при минимальном аннотировании данных экспертом возможно достигнуть значительных показателей классификации неинституциональных текстов по тональности ($F1 = 95\%$ для текстов с стратифицированной депривацией).

Ключевые слова: автоматическое определение тональности текста, социально-сетевой дискурс, машинное обучение с учителем, неглубокая нейронная сеть, депривация

1. Введение

Проблема автоматического определения тональности текста (ТТ) приобретает свою актуальность в рамках изучения общественного мнения (англ. opinion mining). Анализ ТТ – сфера активных исследований последних десятилетий. На сегодняшний день необходимы инструменты для автоматической обработки огромных объемов текстов. Для поиска решений разработаны такие подходы, как: а) опирающиеся на эмотивную лексику (анализ по словарям и правилам); б) машинное обучение с учителем; в) машинное обучение без учителя.

Словарный метод классификации тональности получил достаточно широкое применение. Он лег в основу таких аналитических систем мониторинга СМИ, как Интегрум [1], Медиалогия [2], IQBuzz [3], PalitrumLab [4], SemanticForce [5].

Его суть сводится к тому, что создаются словари эмотивных слов с заведомо определенной тональностью и в исследуемом тексте определяется, слова из какого словаря (негативного или позитивного) в данном тексте преобладают. Обычно анализ с помощью словарей происходит по определенным правилам (учитываются синтагматические границы, отмеченные знаками препинания, особо обрабатываются указанные в поиске словосочетания, принимаются во внимание диминутивы и аугментативы, а также отрицание). Кроме того, разработана разновидность этого метода, исходящая из установки, что не все слова играют равную роль в формировании ТГ (теоретико-графовые модели). Такие модели выстраивают граф исследуемого текста, ранжируют его вершины, определяют «вес» каждого слова на основе тонального словаря и ранга слова-вершины [6].

Тем не менее, данный подход не лишен недостатков. Во-первых, экспрессивно-оценочную окраску слово получает, только становясь частью высказывания, т. е. элементом речи [7, 8]. Пока слово представляет собой элемент языка, оно не несет экспрессивности, даже если принадлежит к кругу эмотивной лексики. Во-вторых, возникают трудности с полисемией, омонимией и идиоматичностью. В-третьих, такой метод оставляет систему нечувствительной к новообразованиям (окказионализмам). Совершенно очевидно, что нарушение грамматических и орфографических правил представляет собой существенное препятствие для методов этого класса, на что указывается в работе [9]. Наконец, создание словарей требует постоянной глубокой работы экспертов-лингвистов.

Метод машинного обучения без учителя работает на основе того принципа, что термины, которые чаще встречаются в этом тексте и в то же время присутствуют в небольшом количестве текстов во всей коллекции, имеют наибольший вес в тексте. Выделив данные термины, а затем, определив их тональность, можно сделать вывод о тональности всего текста.

Данное исследование выполнено в рамках метода машинного обучения с учителем. Выбор обусловлен, в первую очередь, особенностями эмпирического материала, к которому применяется алгоритм анализа. Это обсуждения экологических проблем России в интернет-сообществах, т. е. сегмент неинституционального социально-сетевого дискурса. Данный вид дискурса относится к неформальному, а значит, в нем велика доля экспрессивно-оценочных окказионализмов, а также опечаток и ошибок. Кроме того, высказывания в форумах являются частью макрополилога (термин Р. К. Потаповой [10]) и, как правило, представляют собой реплики диалогичного, а не монологичного характера. В силу этого большинство высказываний — сравнительно малые тексты (1–9 предложений), что не позволяет применять машинное обучение без учителя.

2. Нейронные сети в классификации текстов по тональности

2.1. Методы

В последнее десятилетие в рамках такого подхода, как машинное обучение с учителем, большое распространение получило применение нейронных сетей. В

частности, Т. Миколовым была создана модель Doc2vec, представляющая собой неглубокую (один скрытый слой) нейронную сеть (более подробно о алгоритме ее работы см. [11]).

Модель Doc2vec представляет собой логическое развитие модели Word2vec, реализующей популярный метод «мешок слов» (англ. bag-of-words). Отличие первой заключается в том, что при создании вектора предложений учитывается порядок слов. В то время как модель Word2vec показывает высокие результаты применительно к очень коротким текстам, таким как сообщения в социальной сети Twitter [6, 12], модель Doc2vec работает эффективнее в рамках более длинных сообщений.

Для решения исследовательской задачи была разработана программа на языке Python¹. При этом использовалась библиотека gensim, разработанная Р. Рекуреком [13] с применением модели Doc2vec, созданной Т. Миколовым [11]. За основу был взят код, опубликованный в статьях [14, 15].

Векторы, созданные моделью Doc2vec, поступают в качестве параметров на вход классификатору stochastic gradient descent (SGD, стохастический градиентный спуск), который, в свою очередь, выносит решение о «знаке» каждого высказывания.

Таким образом, нейронные сети позволяют вести анализ текста, не привязываясь к словарному значению слова (т. е. к языку как системе), а опираясь на реальное словоупотребление (речь как реализацию коммуникативной деятельности), что предельно важно.

Особенности функционирования нейронных сетей имеют еще один плюс: они кроссдоменны, т. е. не зависят ни от тематики высказываний, ни от языка. Это подтверждается работами [6, 9-12]. В дальнейших исследованиях предполагается проверить работу созданного программного продукта, исследовав тексты не только экологической тематики, но и политической, экономической, конфессиональной и т. д.

2.2. Параметры формализованного описания дискурса

Разработанная программа уникальна тем, что позволяет учитывать параметры социально-сетевого дискурса (ССД), предложенные Р. К. Потаповой [10, 16]. С ее помощью можно исследовать влияние параметров ССД на качество автоматической классификации. Все параметры высказывания делятся на четыре группы:

- а) метаязыковые, или паспорт реплики — сайт, тема-стимул, дата, автор;
- б) типы ССД по содержанию:
 - тональность высказывания (положительная — отрицательная);
 - монотематический — политематический: т.е. содержащий информацию касательно одной или нескольких тематик;
 - информационно не насыщенный (низкоконтекстуальный) — информационно насыщенный (высоконтекстуальный), т. е. упоминающий известные события, включающий ссылки и т.д.;

¹ https://ru.wikipedia.org/wiki/Python_Software_Foundation_License

- не провоцирующий на полемику, конкретные действия, поступки — провоцирующий на полемику, конкретные действия, поступки.
- в) типы ССД по форме:
- тип дискурса (дистантный — прямой), в данном исследовании это всегда дистантный дискурс в силу особенностей интернет-коммуникации;
 - канал связи (прямой — опосредованный), в нашем исследовании — опосредованный, через компьютер или другой гаджет, позволяющий выйти в интернет;
 - в реальном времени (on-line) / отложенный (off-line);
 - монокронный — полихронный (есть ли существенная временная дистанция между текстом-стимулом и текстом-реакцией).
- г) типы ССД по функции:
- информирующий — содержащий точку зрения отправителя сообщения — побуждающий с определенной целевой установкой к совершению конкретных действий, поступков (в частности, к деструктивным, реализующийся по схеме «стимул — прагматическая реакция в виде конкретного деструктивного действия»);
 - рассчитанный на целевую ограниченную группу пользователей (в т. ч. на одного) — на неограниченное число пользователей.

Обращаем внимание на то, что составляющие каждой оппозиции перечисляются в порядке «невыраженный параметр — выраженный параметр». В дальнейшем эти понятия будут важны для анализов результатов классификации.

Эта система параметров является большим шагом в исследовании социально-сетевого дискурса (ССД). Она позволяет четко представить и, самое главное, зафиксировать в цифровом, усваиваемом компьютером виде структуру социально сетевого дискурса. В ней отражаются такие стороны ССД, как содержание (что говорится), форма (как говорится) и функции (для чего говорится).

В данном исследовании также сформулирован такой параметр типов социально-сетевого дискурса по содержанию, как вид депривации: личностная и стратифицирующая. Под *личностной депривацией* понимается такое негативное психическое состояние, которое вызвано лишением возможности удовлетворения определенных потребностей только у отдельных индивидов, не входящих в одну социальную группу. В свою очередь, под *стратифицированной депривацией* понимается осознание индивидом отсутствия возможности не только у него одного удовлетворить свои насущные потребности, но и у большинства других людей, принадлежащих к одной и той же социальной группе. Выделение такого параметра обосновано тем, что, находясь в состоянии стратифицирующей депривации, человек более склонен к организованным агрессивным действиям. Это объясняется тем, что в таком положении индивид видит корень проблемы не в себе, а в устройстве общества и, следовательно, считает невозможным разрешить конфликт только своими силами. Самое важное заключается в том, что именно параметр видов депривации, как будет показано в разделе 2.3, окажется одним из наилучших

для классификации текстов по тональности, а значит, обладает наибольшей практической ценностью.

Так, дискурс агитирования добровольцев на весенние лесопосадки будет описываться в рамках системы параметров следующим образом: тональность — положительная, депривация отсутствует, информационно насыщенный, монотематический, непровоцирующий; монокронный, в режиме реального времени, опосредованный, дистантный (последние три параметра тождественны для всех высказываний виртуального ССД, поэтому в других примерах они не будут указаны); рассчитанный на неограниченную группу пользователей, функция — побуждающий на определенные действия, проблематика — экология : леса : посадки. Например: *Эко-активисты, волонтеры, любители природы! ЭКА приглашает на посадку леса! 6 мая в 10.00 встречаемся на остановке "Птицефабрика "Вараксина". Перчатки и лопаты будут выданы. Информацию о количестве участвующих направить в личку Шивриной Александре до 10.00 4 мая 2016 года. Ждем старых и новых друзей ;) Репост приветствуется!* А вот критика противопожарных мер будет выглядеть так: тональность негативная, депривация стратифицированная, информационно ненасыщенный (конкретно для нижеприведенного примера), провоцирующий; монокронный; неограниченная аудитория, функция — сообщающий точку зрения говорящего, проблематика — экология : лес : пожары: *А пожары хоть потушены? Или тоже к первому сентября потушат? Вообще-то выгодное дельце с точки зрения освоения средств: что-то там восстанавливаешь, а оно сгорает. Ну, вроде как снег в Сочи завозишь, а он тает-тает.*

В качестве примера высказывания с личностной депривацией приведем следующее: *Это ТЫ, ИМЕННО ТЫ, "тупой трусливый ленивый совок". Ты же "безмозглый раб". Скулящий тут по заказу и не желающий видеть дальше своего носа!* Описание: негативный, личностная депривация, информационно ненасыщенный, монотематический, провоцирующий; монокронный; для ограниченной аудитории, функция – сообщающий точку зрения говорящего, проблематика — политические мировоззрения (это раскрывается в контексте дискуссии).

Как видно из вышеприведенных примеров, разработанная Р.К. Потаповой система параметров социально-сетевого дискурса позволяет характеризовать дискурс с разных ракурсов, что немаловажно для исследований компьютерной и прикладной лингвистики.

2.3. Алгоритм программного продукта

Несмотря на то, что нейронные сети не требуют предварительной обработки текстов, на первом этапе работы все же необходимо провести аннотирование определенной части выборки для тренировки нейронной модели и оценки точности классификации. Для этого лингвист аннотирует высказывания, характеризуя их по каждому из упомянутых в подразделе 2.2 параметров.

После создания аннотированной базы данных программа получает на вход книгу Excel. Особенностью программы является многоступенчатая классификация: сначала на основе результатов аннотирования высказывания разбиваются на группы положительных и негативных (каждому высказыванию

присваивается соответствующий ярлык), затем негативные получают по второму ярлыку, в соответствии с исследуемым параметром (например, ярлык «личностная депривация» или «стратифицирующая депривация»). Однако, как упомянуто в подразделе 2.2, для сравнения классификация была также проведена без дальнейшего разбиения на группы по выбранному параметру. Это необходимо для демонстрации того, какое значение играет разработанная Р. К. Потаповой система параметров ССД в автоматической классификации текстов по тональности.

Затем выборка разделяется на обучающую и тестовую группы в соотношении 4 : 1. Модель Doc2vec проходит обучение (формирует словарь выборки, строит векторы высказываний, сравнивает их с ярлыками), а затем применяет выработанные векторы на тестовой выборке. Классификация может осуществляться с помощью разных техник; в данном исследовании была взята логистическая регрессия. Далее происходит оценка эффективности работы классификатора (результаты работы классификатора сверяются с результатами ручного аннотирования). Результаты сравнения отображаются на матрице неточностей (confusion matrix, рис. 1).

В матрице неточностей строки означают количество настоящих, данных представителей классов, а столбцы — количество предсказанных классов. Пересечения строк и столбцов по главной диагонали показывают, сколько представителей каждого класса классификатор верно распознал. Так, на рисунке 1 ячейка 1*1 показывает, сколько низкоконтекстуальных высказываний классификатор распознал, ячейка 2*2 — высококонтекстуальных, а ячейка 3*3 — сколько положительных высказываний распознано верно. Как видно из рисунка, лучше всего были распознаны высококонтекстуальные высказывания. Ячейка на пересечении первой строки и второго столбца показывает, сколько низкоконтекстуальных высказываний было ошибочно принято классификатором за высококонтекстуальные и т. д.

2.4. Эксперимент

Разработанный программный продукт позволяет сравнивать как воздействие разных пар параметров на эффективность обучения модели Doc2vec, так и качество классификации с учетом параметров и без него.

В ходе эксперимента проводилась поочередная классификация текстов одной и той же выборки по тональности на основе каждого из перечисленных в предыдущем разделе параметров. Как уже было описано в разделе 2.3, сначала тексты разбиваются на положительные и негативные, затем негативные группируются на основе выбранного параметра.

Для положительных текстов подобного разбиения не проводилось, во-первых, в силу того, что многие параметры доступны только для негативных высказываний, а во-вторых, поскольку положительные высказывания не стоят в фокусе данного исследования. Так, высказывание положительной тональности не может содержать депривации или быть провоцирующим, но может быть нацеленным на ограниченную аудиторию (но в данном исследовании это не актуально).

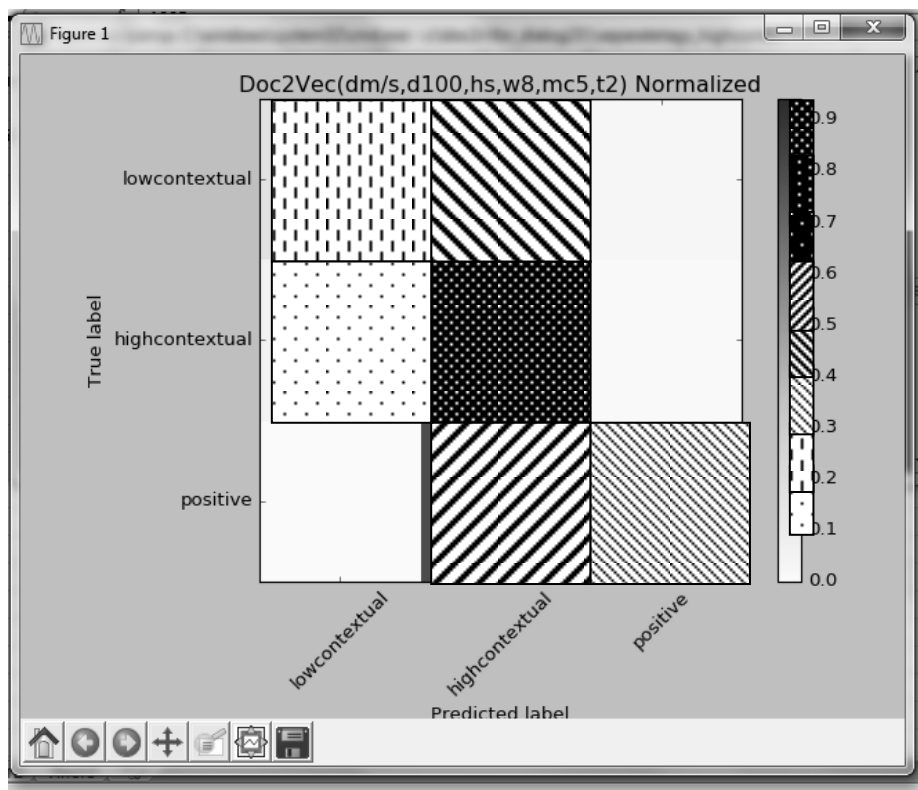


Рис. 1. Графическое представление матрицы неточностей

2.3. Результаты

Использование различных параметров при классификации ТТ с помощью модели Doc2vec показало, что указанные в разделе 2.2 параметры оказывают разное влияние на результаты классификации. Другими словами, распознавание негативных высказываний менялось при включении различных параметров ССД. Так, одним из наиболее подходящих параметров для корректного выявления негативных текстов является именно введенное впервые в нашем исследовании разделение по виду депривации «личностная/стратифицирующая депривация». При опоре на такой параметр программа определяет высказывания со стратифицированной депривацией с точностью 92%, т. е. лишь несколько положительных высказываний классифицируются как негативные при том, что истинно негативные высказывания распознаются корректно (см. рис. 2).

Также эффективно распознавались высказывания со следующими характеристиками: нацеленные на неограниченную аудиторию, провоцирующие и выражающие мнение говорящего (80–82%).

Стоит признать, что собранный нами эмпирический материал не обладает репрезентативностью и полнотой. Количество высказываний $n = 1330$ шт.

(около 500 Кб аннотированной базы данных), из них только 140 высказываний положительные. Трудность увеличения числа положительных высказываний обусловлена высоким уровнем цинизма пользователей Рунета по отношению как к экологическим проблемам страны, так и к мероприятиям, нацеленным на их решение. Безусловно, такой существенный перевес в сторону негативных высказываний может оказать влияние на качество работы классификатора. Но на данном этапе исследования было важно отработать алгоритм классификации, а затем уже наращивать выборку.

Следует отметить, что не всегда лучше распознавались высказывания с выраженным параметром (о невыраженных и выраженных параметрах см. раздел 2.2). На рисунке 2 представлены результаты, где лучше распознавались высказывания с выраженными параметрами, а на рисунке 3 – с невыраженными. Среди невыраженных лучше всего распознавались высказывания с такими характеристиками как монотематичность (91%) и контекстуальная ненасыщенность (83%).

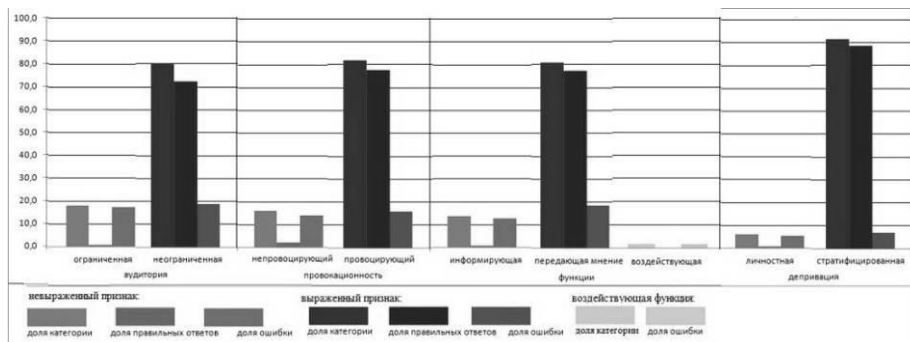


Рис. 2. Успешное распознавание высказываний с выраженным параметром

Не следует забывать, что на данном этапе работы исследователю не столько важно, чтобы программа корректно распознавала разбиение высказываний внутри параметров (например, истинно провоцирующие и истинно не провоцирующие), сколько важно верное разделение на положительные и негативные высказывания. Важно показать, как фокусирование программы на тех или иных признаках позволяет улучшить качество классификации текстов по тональности.

Несмотря на то, что при учете таких параметров, как вид депривации и тематичность высказывания, доля корректно распознанных положительных высказываний была близка к нулю, эти параметры признаются нами оптимальными для распознавания негативных высказываний, стоящих в фокусе внимания данного исследования. Доля правильно распознанных высказываний по отношению к всей выборке составляет 92% и 91% соответственно. Для сравнения: в исследовании [9] это значение составляет 72,8% для негативных высказываний. В табл. 1 приведены значения точности, полноты и меры F1 для классов «стратифицированная депривация» и «монотематичность».

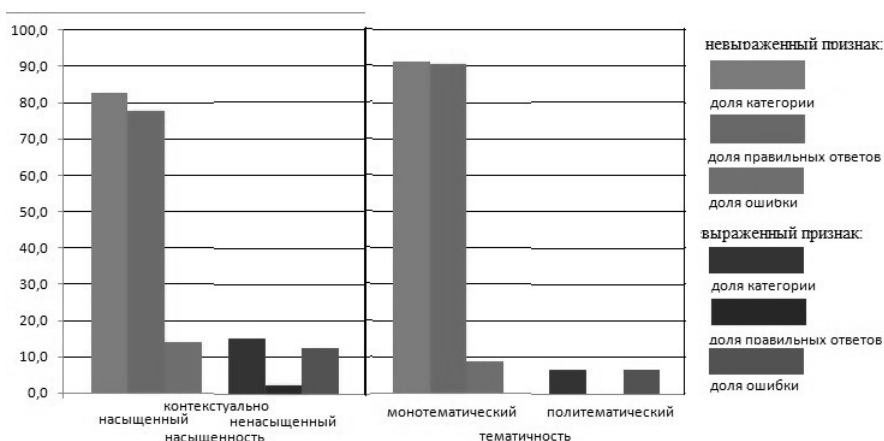


Рис. 3. Успешное распознавание высказываний с невыраженным параметром

Однако в таблице приведены характеристики только для указанных двух классов. Как видно из рисунков 2 и 3, не все категории успешно распознаются. Трудно распознаваемые категории снижают общую точность классификатора до 54%. Тем не менее, не стоит умалять результатов исследования. Если разработанная программа используется для мониторинга социальных сетей (а именно, для поиска пользователей, проявляющих вербальную агрессию) и за гипотезу H_0 принимается утверждение «данное высказывание агрессивно», то ошибка второго рода (неверное истолкование высказывания как агрессивного) будет предпочтительнее ошибки первого рода (пропуск настоящей вербальной атаки). Что же касается диагностирования агрессивных высказываний (которые входят в класс депривированных высказываний), то мера F1 для этого класса достаточно высока.

Таблица 1. Точность, полнота и F1 для оптимальных параметров ССД (%)

Параметр	Точность	Полнота	F1
Стратифицированная депривация	96,85	92,67	95,96
Монотематичность	99,09	91,29	94,83

Разработанный нами алгоритм позволяет проводить автоматическую классификацию, по точности сравнимую (а для отдельных параметров и превосходящую) другие исследования данной проблемы. Для сравнения, в работе [12], позиционирующей свой проект как передовой, максимальное значение $F_1 = 88,93$, в то время как в нашем исследовании $F_1 = 95,96$. Также, в работе [6] максимальный результат $F_1 = 86,58$.

3. Выводы

В данном исследовании развивается постулат о том, что при определении тональности текстов социально-сетевого дискурса следует опираться не на тональные словари (в силу активных словотворческих процессов в ССД), а на экспрессивность как характеристику высказывания в целом. Этого позволяют добиться нейронные сети.

Следует обратить внимание на то, что модели нейронных сетей позволяют при относительно небольшой предобработке текстов (создание аннотированной базы данных, на основе которой происходит машинное обучение) достигнуть большой точности при решении данной задачи. Т. е. при уже созданной в рамках данного исследования системе векторов слов автоматическая классификация будущих текстов по тональности будет проходить полностью без аннотирования текстов. Это ставит нейронные модели в выигрышное положение по сравнению с методом словарей и правил.

Также существенным преимуществом моделей нейронных сетей является то, что они работают, так сказать, с «живым текстом»: программа совершенно «не знает», как должно писаться то или иное слово, в каком значении оно употреблено, как изначально звучит тот или иной фразеологизм. Она считывает тексты с неисправленными опечатками, ошибками, с нетронутыми окказионализмами. Поэтому не возникает необходимости подключать модули исправления ошибок, что также облегчает работу исследователя. Стоит подчеркнуть, что ни разработка нейронных моделей, ни система параметров социально-сетевого дискурса (ССД) не являются заслугой авторов данного исследования. Заслугой является разработка кода, позволяющего учитывать параметры ССД при работе нейронной модели в автоматической классификации текстов по тональности.

Кроме того, нами сформулирован новый параметр, характеризующий ССД по содержанию, а именно вид депривации: личностная и стратифицирующая. Под депривацией нами понимается выражение неудовлетворенности индивидом сложившейся ситуацией. Разделение на личностную и стратифицирующую депривации позволяет отделить проблемы личного характера от проблем, представляющихся индивиду системными. Это разделение важно потому, что за выражением стратифицирующей депривации часто следуют попытки организовать отдельных представителей общества или целые группы для поиска решения проблемы. Такие случаи представляют особое значение для исследователей общественного мнения.

Использование различных параметров при классификации ТТ с помощью модели Doc2vec доказало, что указанные в разделе 2.2 параметры оказывают разное влияние на результаты классификации. Наиболее подходящими выбраны параметры «монотематический/политематический» и «вид депривации (личностная/стратифицирующая)», позволяющие достигнуть 91 и 92% точности при определении негативных высказываний со стратифицированной депривацией.

Таким образом, разработан инструмент, позволяющий эффективно автоматически классифицировать небольшие тексты неинституционального характера по тональности.

Исследование проводилось при поддержке Российского научного фонда (грант № 14-18-01059, руководитель проекта — Р. К. Потапова). Подана заявка на государственную регистрацию данной программы.

Литература

- [1] Интегрум, аналитическая система мониторинга, URL: www.integrum.ru.
- [2] Медиалогия, аналитическая система мониторинга, URL: www.mlg.ru.
- [3] IQBuzz, аналитическая система мониторинга, URL: www.iqbuzz.pro.
- [4] PalitrumLab, аналитическая система мониторинга, URL: www.palitrumlab.ru.
- [5] SemanticForce, аналитическая система мониторинга. www.semanticforce.net.
- [6] Tang D., Wei F., Yang N., Zhou M., Liu T., Qin B. Learning sentiment-specific word embedding for twitter sentiment classification // Proceeding of the 52th Annual Meeting of the Association for Computational Linguistics. ACL. 2014. pp. 1155-1166. <http://anthology.aclweb.org/P/P14/P14-1146.pdf>.
- [7] Бахтин М. М. Эстетика словесного творчества / Сост. С. Г. Бочаров, примеч. С. С. Аверинцев и С. Г. Бочаров. М.: Искусство. 1979.
- [8] Вольф Е. М. Функциональная семантика оценки. Изд. 2-е, доп. М.: Едиториал УРСС. 2002. 280 с.
- [9] Thelwall M., Buckley K., Paltoglou G., Cai D., & Kappas A. Sentiment strength detection in short informal text // Journal of the American Society for Information Science and Technology. 2010. 61(12). P. 2544–2558.
- [10] Потапова Р. К. Депривация как базовый механизм вербального и паравербального поведения человека (на материале социально-сетевой коммуникации) // Речевая коммуникация в информационном пространстве. Коллективная монография. Отв. ред. Потапова Р. К., М.: ЛЕНАНД. 2017. С. 17-36.
- [11] Mikolov T., Le Q. Distributed representations of sentences and documents // Proceedings of the 31st International Conference on Machine Learning, Beijing, China. 2014. JMLR: W&CP vol. 32. Copyright 2014 by the author(s). https://cs.stanford.edu/~quocle/paragraph_vector.pdf.
- [12] Mohammad S.M., Kiritchenko S., Zhu X. NRC-Canada: Building the state-of-the-art in sentiment analysis of tweets // Second Joint Conference on Lexical and Computational Semantics (*SEM), Vol. 2: Seventh International Workshop on Semantic Evaluation (SemEval 2013), Atlanta, Georgia, 2013. Pp. 321–327.
- [13] Rehurek R., Sojka P. Software framework for topic modelling with large corpora // Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks. ELRA. 2010. Pp. 45-50 <https://github.com/RaRe-Technologies/gensim#citing-gensim>.
- [14] Czerny M. Modern Methods for Sentiment Analysis. https://districtdatalabs.silvrback.com/modern-methods-for-sentiment-analysis#disqus_thread.
- [15] Word embeddings for fun and profit: document classification with gensim. https://github.com/RaRe-Technologies/movie-plots-by-gene/blob/5a2d9157f9bf1bf908794051597b7851333dcfca/ipynb_with_output/Document%20classification%20with%20word%20embeddings%20tutorial%20with%20output.ipynb#L1403.

- [16] Потапова Р. К. Социально-сетевой дискурс как объект междисциплинарного исследования // мат. 2-й международной конференции «Дискурс как социально-сетевая деятельность», М.: МГЛУ, 2014. с. 20-32

Social-network discourse types in automated sentiment classification

N Maslova

Moscow State Linguistic University

The article covers approaches to automated sentiment analysis task and describes a new algorithm that builds the state-of-the-art in this scientific problem. There are three approaches: rules- and dictionaries-based, unsupervised learning and supervised learning. Their pros and cons are compared.

A new programme was created under the supervised learning method with the help of neural network module Doc2vec. The programme's specialization is short informal texts of ecologic topic which are parts of macropolylogues in social network discourse (SND). The concept of macropolylogue along with the methodology of social network studies from linguistic point of view were introduced by Potapova R. K. Neural network models efficiently cope with the problems caused by the Russian language peculiarities (inflection) as well as by informal discourse (low level of grammar, word coinage, idiom modifications).

The programme's main advantage is that it classifies texts with consideration of the SND parameters system developed by R. Potapova. These parameters influence the classification accuracy in different ways.

The current research has expanded the alluded system by one more opposition: private vs. stratified deprivation. Their practical application is substantiated. Stratified deprivation along with the monotopic parameter are the best to foster sentiment classification.

The article shows how to reach excellent classification results (the best $F_1 = 95\%$) by minimal preliminary text annotation.

Keywords: automated sentiment analysis; social network discourse; supervised learning; deprivation.