

ИНФОРМАЦИОННЫЕ ТЕХНОЛОГИИ И СИСТЕМЫ

УДК 004.912

С. Ф. ЛИПНИЦКИЙ

АНАЛИЗ ТОНАЛЬНОСТИ ТЕКСТОВЫХ СООБЩЕНИЙ В ИНТЕРНЕТЕ НА ОСНОВЕ МОДЕЛИРОВАНИЯ ВЕРБАЛЬНЫХ АССОЦИАЦИЙ

Объединенный институт проблем информатики НАН Беларуси

(Поступила в редакцию 31.09.2013)

Введение. Актуальными направлениями исследований в области компьютерной лингвистики являются поиск тонально-окрашенной информации в Интернете и анализ ее тональности. Под тонально-окрашенными понимают текстовые сообщения, в которых выражено эмоциональное отношение их авторов к конкретным объектам или событиям. Своевременное выявление и аналитическая обработка такой информации позволяют минимизировать материальные потери или ущерб репутации, которые она может причинить персоне или организации, фигурирующей в тексте. Результаты анализа тонально-окрашенных интернет-сообщений полезны маркетологам, социологам и другим специалистам, изучающим общественное мнение с целью планирования своей деятельности.

Предложенный в данной статье подход к анализу тонально-окрашенной информации в отличие от существующих основан на использовании модели представления знаний на основе вербальных ассоциаций [1]. Эта модель позволяет учитывать наличие в текстах не только тонально-окрашенной лексики, но и семантических связей между словами, соответствующими ассоциативным отношениям между обозначаемыми ими сущностями предметной области. Благодаря этому модель обеспечивает формирование оценок тональности сообщений с использованием тонально-окрашенных корпусов текстов.

1. Вербально-ассоциативные сети. Вербально-ассоциативная сеть – это граф, вершинами которого являются слова, а ребрами – вербально-ассоциативные связи между ними [1]. Вершины сети помечены значениями информативности слов, а каждому ребру соответствует сила вербально-ассоциативной связи между словами. Такие сети могут быть построены для предложений, текстов и совокупностей текстов.

1.1. Тонально-окрашенные тематические корпуса текстов. Для оценки тональности сообщений в Интернете будем использовать совокупность тематических корпусов текстов, среди которых имеется корпус Ct с тонально-окрашенными текстами. Он состоит из совокупности тонально-окрашенных подкорпусов, каждому из которых соответствует некоторая оценка тональности. При n -бальной шкале оценок количество таких подкорпусов должно быть равно n , т. е. $Ct = \{Cp_i \mid i = \overline{1, n}\}$. Всякий i -й подкорпус Cp_i включает в себя текстовые документы одинаковой тональности, т. е. корпус Cp_i – пара $\langle Cp_i, Ev_i \rangle$, где Ev_i – оценка тональности каждого документа из множества Cp_i .

1.2. Вербально-ассоциативное отношение. Обозначим через W множество всех слов некоторого языка L . Тогда отношение толерантности Θ (рефлексивное и симметричное бинарное отношение) на множестве W назовем *вербально-ассоциативным*, если любая упорядоченная пара слов (a, b) из множества W является элементом отношения Θ тогда и только тогда, когда слова a

и b из этой пары содержатся хотя бы в одном предложении языка L . Если пара (a, b) любых слов из множества W является элементом отношения Θ , т. е. $(a, b) \in \Theta$, то (a, b) (по аналогии с [2]) будем называть *вербально-ассоциативной парой*.

1.3. **Словари тонально-окрашенной лексики.** Для поиска тонально-окрашенной информации и анализа ее тональности будем использовать следующие лингвистические словари.

Частотный словарь тонально-окрашенной лексики (табл. 1) $Disc_a = \{(a, n_{Cf}^a, n_{(Cp_1, Ev_1)}^a, n_{(Cp_2, Ev_2)}^a, \dots, n_{(Cp_n, Ev_n)}^a) \mid a \in W_{Cf}\}$. Здесь a – словоформа, n_{Cf}^a и $n_{(Cp_i, Ev_i)}^a$ ($i = \overline{1, n}$) – абсолютные частоты ее появления соответственно в полном корпусе текстов Cf и в i -м тонально-окрашенном подкорпусе, W_{Cf} – множество всех словоформ полного корпуса. (Полный корпус текстов Cf – это объединение всех тематических корпусов.)

Таблица 1. Фрагмент частотного словаря тонально-окрашенной лексики

Словоформа	n_{Cf}^a	$n_{(Cp_1, Ev_1)}^a$...	$n_{(Cp_n, Ev_n)}^a$
...				
Качественный	0204055	0056534	...	0014445
Качественного	0401657	0074526	...	0023747
...				

Частотный словарь вербально-ассоциативных пар слов (табл. 2), т. е. совокупность кортежей $Disc_{ab} = \{(a, b), n_{Cf}^{ab}, n_{(Cp_1, Ev_1)}^{ab}, n_{(Cp_2, Ev_2)}^{ab}, \dots, n_{(Cp_n, Ev_n)}^{ab}\} \mid a, b \in W_{Cf}, n_{Cf}^{ab} \neq 0, n_{(Cp_i, Ev_i)}^{ab} \neq 0, i = \overline{1, n}\}$, где $n_{Cf}^{ab}, n_{(Cp_i, Ev_i)}^{ab}$ – абсолютные частоты совместной встречаемости слов a и b в одном и том же предложении полного Cf и i -го тонально-окрашенного подкорпуса Cp_i ($i = \overline{1, n}$).

Таблица 2. Фрагмент частотного словаря вербально-ассоциативных пар слов

Пара словоформ	n_{Cf}^{ab}	$n_{(Cp_1, Ev_1)}^{ab}$...	$n_{(Cp_n, Ev_n)}^{ab}$
...				
(Качественный, отсутствует)	03020	00543	...	00121
(Предупреждение, избыточного)	04023	00623	...	00242
...				

Словарь словоизменительных парадигм [1] $Disc_{par} = \{(a, Par_a) \mid a \in W_{Cf}, a \in Par_a\}$, где Par_a – множество всех словоизменений слова a (включая a).

Словарь синонимичных словоформ [1] $Disc_{syn} = \{(a, Syn_a) \mid a \in W_{Cf}, a \in Syn_a\}$, где Syn_a – множество всех синонимов слова $a \in W_{Cf}$.

1.4. **Задачи интернет-мониторинга тонально-окрашенной информации.** В процессе интернет-мониторинга и анализа тонально-окрашенной информации решаются три основные задачи:

- документальный поиск веб-страниц;
- фактографический поиск на найденных страницах конкретной тонально-окрашенной информации, релевантной запросу;
- анализ, оценка найденной тонально-окрашенной информации и составление отчета по результатам мониторинга.

Этим основным этапам интернет-мониторинга должно предшествовать построение вербально-ассоциативных сетей для тонально-окрашенных текстов; множеств тонально-окрашенных текстов; кратких сообщений; предложений; веб-страниц.

1.5. **Вербально-ассоциативные сети тонально-окрашенных текстов.** Пусть T – некоторый непустой текст языка L . Обозначим через W_T множество всех слов текста T , а через Θ_T – сужение отношения Θ на множество W_T , т. е. $\Theta_T = \Theta \cap (W_T \times W_T)$. Обозначим через $G\Theta_T$ граф отношения Θ_T . Пометим каждую вершину a графа $G\Theta_T$ значением информативности I_T^a этого слова

(с учетом синонимии и словоизменения), а каждое ребро (a, b) – значением силы вербально-ассоциативной связи I_T^{ab} слов a и b (также учитывая синонимию и словоизменения). Обозначим полученный граф через Net_T , который назовем *вербально-ассоциативной сетью* текста T . Приведем формулы для вычисления значений параметров I_T^a и I_T^{ab} .

В [1] показано, что информативность слова в тексте вычисляется как отношение абсолютной частоты m_T^a встречаемости слова в тексте T к абсолютной частоте m_{Cf}^a его появления в полном корпусе текстов Cf (с учетом словоизменения и синонимии), т. е.

$$I_T^a = m_T^a / m_{Cf}^a.$$

Словоизменения зафиксированы в словаре словоизменительных парадигм Dic_{par} , а их синонимы – в словаре синонимичных словоформ Dic_{syn} . После нахождения синонимов каждого слова необходимо найти также их словоизменения в словаре Dic_{par} . С учетом этого обстоятельства общая формула для информативности I_T^a примет вид

$$I_T^a = \frac{n_T^a + \sum_{b \in Par_a, b \neq a} n_T^b + \sum_{c \in Syn_a, c \neq a} \left(n_T^c + \sum_{d \in Par_c, d \neq c} n_T^d \right)}{n_{Cf}^a + \sum_{b \in Par_a, b \neq a} n_{Cf}^b + \sum_{c \in Syn_a, c \neq a} \left(n_{Cf}^c + \sum_{d \in Par_c, d \neq c} n_{Cf}^d \right)}, \quad (1)$$

где n_T^a и n_{Cf}^a – абсолютные частоты встречаемости слова a в тексте T и полном корпусе текстов Cf (без учета словоизменения и синонимии); n_T^b и n_{Cf}^b – эти же частоты для словоформы b , входящей в состав парадигмы слова a ; n_T^c и n_{Cf}^c – аналогичные параметры для синонима c слова a ; n_T^d и n_{Cf}^d – абсолютные частоты встречаемости для словоформы d , входящей в состав парадигмы слова c .

Аналогичную формулу будем использовать для вычисления силы вербально-ассоциативной связи P_T^{ab} слов a и b :

$$P_T^{ab} = \frac{n_T^{ab} + \sum_{\substack{c \in Par_a, c \neq a \\ d \in Par_b, d \neq b}} n_T^{cd} + \sum_{\substack{r \in Syn_a, r \neq a \\ s \in Syn_b, s \neq b}} \left(n_T^{rs} + \sum_{\substack{p \in Par_r, p \neq r \\ q \in Par_s, q \neq s}} n_T^{pq} \right)}{n_{Cf}^{ab} + \sum_{\substack{c \in Par_a, c \neq a \\ d \in Par_b, d \neq b}} n_{Cf}^{cd} + \sum_{\substack{r \in Syn_a, r \neq a \\ s \in Syn_b, s \neq b}} \left(n_{Cf}^{rs} + \sum_{\substack{p \in Par_r, p \neq r \\ q \in Par_s, q \neq s}} n_{Cf}^{pq} \right)}, \quad (2)$$

где n_T^{ab} , n_T^{cd} , n_T^{rs} , n_{Cf}^{ab} , n_{Cf}^{cd} , n_{Cf}^{rs} – абсолютные частоты совместной встречаемости пар слов (a, b) , (c, d) , (r, s) и (p, q) в одном и том же предложении текста T и полного корпуса текстов Cf .

При практической реализации системы интернет-мониторинга вербально-ассоциативную сеть Net_T текста T представим в виде его поискового образа:

$$PO_T = \{(a, I_T^a), (b, I_T^b), \dots, ((c, d), P_T^{cd}), ((e, f), P_T^{ef}), \dots | I_T^a, I_T^b, \dots > I_T^0, P_T^{cd}, P_T^{ef}, \dots > P_T^{00}\}, \quad (3)$$

где a, b, c, d, e, f, \dots – слова текста T ; $(c, d), (e, f), \dots$ – вербально-ассоциативные пары слов; I_T^a, I_T^b, \dots – информативность слов a, b, \dots текста T ; $P_T^{cd}, P_T^{ef}, \dots$ – сила вербально-ассоциативной связи между словами; I_T^0 и P_T^{00} – пороговые значения информативности и силы связи.

1.6. Вербально-ассоциативные сети множеств тонально-окрашенных текстов. Пусть $\{T_i | i = 1, n_T\}$ – некоторое непустое множество тонально-окрашенных документов. Обозначим через T текст, полученный в результате объединения всех предложений каждого из текстов T_i . Тогда под вербально-ассоциативной сетью множества текстов $\{T_i | i = 1, n_T\}$ будем понимать сеть текста T . Параметры I_T^a и P_T^{ab} вычисляются по формулам (1) и (2), а поисковый образ текста T строится в соответствии с формулой (3).

1.7. Вербально-ассоциативные сети кратких сообщений и предложений. Вербально-ассоциативная сеть краткого сообщения ST – это подграф сети множества всех документов из тонально-окрашенного тематического корпуса текстов Ct . Вершины этого подграфа помечены значениями информативности соответствующих слов из корпуса текстов Ct , а ребра – показателями силы связи между словами из вербально-ассоциативной сети корпуса Ct .

Информативность слов краткого сообщения ST вычислим по формуле, аналогичной выражению (1):

$$I_{ST}^a = \frac{n_{Ct}^a + \sum_{b \in Par_a, b \neq a} n_{Ct}^b + \sum_{c \in Syn_a, c \neq a} \left(n_{Ct}^c + \sum_{d \in Par_c, d \neq c} n_{Ct}^d \right)}{n_{Cf}^a + \sum_{b \in Par_a, b \neq a} n_{Cf}^b + \sum_{c \in Syn_a, c \neq a} \left(n_{Cf}^c + \sum_{d \in Par_c, d \neq c} n_{Cf}^d \right)},$$

где n_{Ct}^a и n_{Cf}^a – абсолютные частоты встречаемости слова a из сообщения ST в корпусе текстов Ct и полном корпусе текстов Cf (без учета словоизменения и синонимии); n_{Ct}^b и n_{Cf}^b – эти же частоты для словоформы b , входящей в состав парадигмы слова a ; n_{Ct}^c и n_{Cf}^c – аналогичные параметры для синонима c слова a ; n_{Ct}^d и n_{Cf}^d – абсолютные частоты встречаемости для словоформы d , входящей в состав парадигмы слова c .

Для вычисления силы вербально-ассоциативной связи между словами краткого сообщения ST воспользуемся аналогом формулы (2):

$$P_{ST}^{ab} = \frac{n_{Ct}^{ab} + \sum_{\substack{c \in Par_a, c \neq a \\ d \in Par_b, d \neq b}} n_{Ct}^{cd} + \sum_{\substack{r \in Syn_a, r \neq a \\ s \in Syn_b, s \neq b}} \left(n_{Ct}^{rs} + \sum_{\substack{p \in Par_r, p \neq r \\ q \in Par_s, q \neq s}} n_{Ct}^{pq} \right)}{n_{Cf}^{ab} + \sum_{\substack{c \in Par_a, c \neq a \\ d \in Par_b, d \neq b}} n_{Cf}^{cd} + \sum_{\substack{r \in Syn_a, r \neq a \\ s \in Syn_b, s \neq b}} \left(n_{Cf}^{rs} + \sum_{\substack{p \in Par_r, p \neq r \\ q \in Par_s, q \neq s}} n_{Cf}^{pq} \right)},$$

где n_{Ct}^{ab} , n_{Ct}^{cd} , n_{Ct}^{rs} , n_{Cf}^{ab} , n_{Cf}^{cd} , n_{Cf}^{rs} – абсолютные частоты совместной встречаемости пар слов (a, b) , (c, d) , (r, s) и (p, q) из сообщения ST в одном и том же предложении корпуса текстов Ct и полного корпуса Cf .

Поисковый образ краткого сообщения ST имеет вид

$$ПО_{ST} = \{(a, I_{Ct}^a), (b, I_{Ct}^b), \dots, ((c, d), P_{Ct}^{cd}), ((e, f), P_{Ct}^{ef}), \dots | I_{Ct}^a, I_{Ct}^b, \dots > I_{Ct}^0; P_{Ct}^{cd}, P_{Ct}^{ef}, \dots > P_{Ct}^{00}\},$$

где a, b, c, d, e, f, \dots – слова сообщения ST ; $(c, d), (e, f), \dots$ – вербально-ассоциативные пары слов из ST ; $I_{ST}^a, I_{ST}^b, \dots$ – информативность слов a, b, \dots сообщения ST ; $P_T^{cd}, P_T^{ef}, \dots$ – сила вербально-ассоциативной связи между словами из ST ; I_T^0 и P_T^{00} – пороговые значения информативности и силы связи.

Предложение можно рассматривать как частный случай краткого сообщения. В связи с этим вербально-ассоциативная сеть и поисковый образ для него строятся так же, как и для сообщения ST .

1.8. Вербально-ассоциативные сети веб-страниц. Рассмотрим произвольную веб-страницу $S_{\text{веб}}$, вербально-ассоциативная сеть которой строится по аналогии с сетью тонально-окрашенного текста T . Значения информативности слов и вербально-ассоциативных пар в сети вычисляются по формулам (1) и (2) в предположении, что $T := S_{\text{веб}}$.

Поисковый образ страницы $S_{\text{веб}}$ аналогичен образу текста T :

$$ПО_{S_{\text{веб}}} = \{(a, I_{S_{\text{веб}}}^a); (b, I_{S_{\text{веб}}}^b), \dots, ((c, d), I_{S_{\text{веб}}}^{cd}), ((e, f), I_{S_{\text{веб}}}^{ef}), \dots | I_{S_{\text{веб}}}^a, I_{S_{\text{веб}}}^b, \dots > I_{S_{\text{веб}}}^0; I_{S_{\text{веб}}}^{cd}, I_{S_{\text{веб}}}^{ef}, \dots > I_{S_{\text{веб}}}^{00}\}.$$

2. Анализ тонально-окрашенной информации. Ему предшествуют поиск релевантных запросу веб-страниц с выдачей их адресов, а также поиск конкретных данных на этих страницах с представлением их информативных фрагментов.

2.1. Структура запроса на мониторинг тонально-окрашенной информации. Запрос пользователя на мониторинг тонально-окрашенной информации включает, как правило, два поля – название объекта мониторинга и его характеристики (табл. 3).

Таблица 3. Пример запроса на мониторинг тонально-окрашенной информации

Объект мониторинга	Характеристика
Компания X.	Информация о компании. Выпускаемая продукция. Отзывы о качестве продукции. Отзывы о руководстве компании.

В позиции объекта мониторинга тонально-окрашенной информации может быть указано физическое лицо, организация, средство массовой информации, событие и т. п. В качестве характеристик объекта мониторинга размещается любая текстовая информация о нем.

2.2. Этапы интернет-мониторинга. В общем случае мониторинг тонально-окрашенной информации реализуется в пять этапов.

На первом этапе индексируется первоначальный запрос пользователя.

На втором этапе корректируется поисковое предписание, полученное в результате индексирования этого запроса. Процесс коррекции реализуется следующим образом. По первоначальному поисковому предписанию проводится поиск релевантных документов в полном корпусе текстов. Найденное множество документов (динамический корпус текстов) индексируется и строится его вербально-ассоциативная сеть на основе полученного поискового образа. Эта сеть предъявляется пользователю, который ее корректирует путем исключения некоторых вершин. С использованием откорректированной вербально-ассоциативной сети строится новое поисковое предписание.

На третьем этапе проводится поиск веб-страниц по откорректированному поисковому предписанию. Найденные страницы анализируются на обновляемость. Из каждой обновленной веб-страницы извлекается вся текстовая информация.

На четвертом этапе выделяются информативные предложения на каждой найденной веб-странице. Строится контекст для всех информативных предложений. Полученные субтексты – это результат мониторинга.

На пятом этапе проводится анализ тональности каждого построенного субтекста с определением оценки тональности. Для этого выявляются наиболее релевантные тематические корпуса текстов. Оценка тональности субтекста совпадает с оценкой тональности релевантного ему корпуса.

Рассмотрим каждый из этапов интернет-мониторинга.

2.3. Индексирование запроса. Пусть $Z_1 = \{Z_1^{(1)}, Z_1^{(2)}\}$ – первоначальный запрос пользователя системы интернет-мониторинга, где $Z_1^{(1)}$ – наименование объекта мониторинга, а $Z_1^{(2)}$ – описание его характеристик.

Индексирование запроса Z_1 сводится к перечислению всех словоформ, присутствующих в цепочках $Z_1^{(1)}$ и $Z_1^{(2)}$, а также их синонимов и словоизменений. Полученное поисковое предписание представим в виде

$$\text{ПП}_1 = \{\text{ПП}_1^{(1)}, \text{ПП}_1^{(2)}\},$$

где $\text{ПП}_1^{(1)} = \{a \mid a \in Z_1 \cup \text{Par}_a \cup \text{Syn}_a\}$, $\text{ПП}_1^{(2)} = \{b \mid b \in Z_2 \cup \text{Par}_b \cup \text{Syn}_b\}$.

2.4. Коррекция поискового предписания. Исключим значения информативности слов и вербально-ассоциативные пары слов из поисковых образов всех документов полного корпуса текстов, т. е. преобразуем выражение (3) к виду

$$\text{ПО}_T = \{a, b, \dots \mid a, b, \dots \in T\}.$$

Сформируем динамический корпус текстов, релевантных запросу Z_1 . С этой целью проведем двухшаговый поиск текстовых документов: сначала по запросу $Z_1^{(1)}$ в полном корпусе текстов Cf , а затем по запросу $Z_1^{(2)}$ в множестве текстов, найденных на первом шаге по запросу $Z_1^{(1)}$. В качестве критерия выдачи будем использовать мощность пересечения поисковых образов документов и поисковых предписаний.

Определим для каждого текста $T \in Cf$ мощность пересечения множеств ПО_T и $\text{ПП}_1^{(1)}$, т. е. $m_1 = |\text{ПО}_T \cap \text{ПП}_1^{(1)}|$. Обозначим через $DZ_1^{(1)}$ множество документов полного корпуса текстов, найденных по поисковому предписанию $\text{ПП}_1^{(1)}$:

$$DZ_1 = \{T \mid T \in Cf, m_1 \geq m_1^0\},$$

где m_1^0 – пороговое значение мощности множества DZ_1 .

В множестве документов DZ_1 проведем поиск документов по поисковому предписанию $\text{ПП}_1^{(2)}$. Обозначим результат поиска (динамический корпус текстов) через DZ_2 :

$$DZ_2 = \{T \mid T \in DZ_1, m_2 \geq m_2^0\},$$

где $m_2 = |\text{ПО}_T \cap \text{ПП}_1^{(2)}|$.

Построим вербально-ассоциативную сеть корпуса DZ_2 как множества тонально-окрашенных текстов. Эту сеть предъявим пользователю для коррекции, т. е. для исключения из сети некоторых вершин, не отражающих содержание запроса. Используя откорректированную сеть, создадим поисковый образ динамического корпуса текстов, который приобретает статус откорректированного поискового предписания:

$$\text{ПП}_2 = \{(a, I_{DZ_2}^a), (b, I_{DZ_2}^b), \dots, ((c, d), P_{DZ_2}^{cd}), ((e, f), P_{DZ_2}^{ef}), \dots \mid I_{DZ_2}^a, I_{DZ_2}^b, \dots > I_{DZ_2}^0; P_{DZ_2}^{cd}, P_{DZ_2}^{ef}, \dots > P_{DZ_2}^{00}\},$$

где a, b, c, d, e, f, \dots – слова из поискового предписания ПП_2 ; $(c, d), (e, f), \dots$ – вербально-ассоциативные пары слов из ПП_2 ; $I_{DZ_2}^a, I_{DZ_2}^b, \dots$ – информативность слов a, b, \dots из ПП_2 в динамическом корпусе текстов DZ_2 ; $P_{DZ_2}^{cd}, P_{DZ_2}^{ef}, \dots$ – сила вербально-ассоциативной связи между словами из ПП_2 в корпусе DZ_2 ; $I_{DZ_2}^0$ и $P_{DZ_2}^{00}$ – пороговые значения информативности и силы связи.

2.5. Поиск веб-страниц. Обозначим через $\text{ПП}'_2 = \{a, b, \dots, (c, d), (e, f), \dots\}$ поисковое предписание, полученное из предписания ПП_2 путем исключения из него показателей информативности слов и силы вербально-ассоциативной связи между ними. Точно так же преобразуем и поисковые образы всех веб-страниц.

Поиск веб-страниц реализуем в два этапа. На первом этапе проведем поиск по предписанию $\text{ПП}'_2$. При поиске будем использовать в качестве критерия выдачи мощность пересечения множеств $m = |\text{ПО}'_s \cap \text{ПП}'_2|$, где $\text{ПО}'_s$ – преобразованный поисковый образ веб-страницы $s \in S_{\text{веб}}$ ($S_{\text{веб}}$ – множество всех веб-страниц). Тогда в качестве результата поиска по предписанию $\text{ПП}'_2$ имеем множество веб-страниц

$$S_{\text{ПП}'_2} = \{s \mid s \in S_{\text{веб}}, m \geq m^0\},$$

где m^0 – пороговое значение мощности множества $S_{\text{ПП}'_2}$.

На втором этапе результаты поиска ранжируем с использованием критерия выдачи, основанного на вычислении косинуса угла между векторами поискового предписания и поискового образа документа. Рассмотрим этот критерий применительно к ранжированию веб-страниц.

Обозначим через W множество всех слов и вербально-ассоциативных пар из полного корпуса текстов S_f . Построим n -мерное евклидово пространство E . Для этого лексикографически упорядочим все слова и вербально-ассоциативные пары из множества W , т. е. сформируем кортеж $W = \langle a_1, a_2, \dots, a_n \rangle$, где a_i – слово или вербально-ассоциативная пара слов из полного корпуса текстов. Для каждой веб-страницы $s \in S_{\text{веб}}$ построим вектор ее поискового образа в пространстве E : $\mathbf{F}_s = (I_{S_{\text{веб}}}^{a_1}, I_{S_{\text{веб}}}^{a_2}, \dots, I_{S_{\text{веб}}}^{a_n})$. Аналогично представим вектор поискового предписания $S_{\text{ПП}'_2}$: $\mathbf{F}_{\text{ПП}'_2} = (J_{S_{\text{веб}}}^{a_1}, J_{S_{\text{веб}}}^{a_2}, \dots, J_{S_{\text{веб}}}^{a_n})$.

При поиске тонально-окрашенной текстовой информации нужно учитывать то обстоятельство, что в Интернете индексируются не сами документы, а веб-страницы, на которых они расположены. Веб-страница может содержать более одного текста, а запрос пользователя, как правило, ориентирован на поиск именно текстовых документов. В связи с этим в векторе \mathbf{F}_s приравняем к нулю все координаты (значения информативности и силы вербально-ассоциативной связи), соответствующие всем нулевым компонентам вектора $\mathbf{F}_{\text{ПП}'_2}$. Обозначим полученный вектор через \mathbf{F}'_s .

Тогда для ранжирования результатов поиска по поисковому предписанию $\text{ПП}'_2$ воспользуемся критерием выдачи

$$\cos \varphi = \frac{\mathbf{F}'_s \mathbf{F}_{\text{ПП}'_2}}{|\mathbf{F}'_s| |\mathbf{F}_{\text{ПП}'_2}|} = \frac{\sum_{i=1}^n I_{S_{\text{веб}}}^{a_i} J_{S_{\text{веб}}}^{a_i}}{\sqrt{\sum_{i=1}^n (I_{S_{\text{веб}}}^{a_i})^2} \sqrt{\sum_{i=1}^n (J_{S_{\text{веб}}}^{a_i})^2}}.$$

2.6. **Поиск тонально-окрашенной информации на веб-страницах.** Фактографический поиск тонально-окрашенных сообщений на найденных веб-страницах сводится к выделению в текстах сообщений информативных фрагментов, релевантных поисковому предписанию $\text{ПП}'_2$. Процедура включает два этапа.

На первом этапе вычисляются информативность всех слов на веб-странице, которые являются элементами множества $\text{ПП}'_2$, а также сила вербально-ассоциативной связи для всех пар слов и информативность всех предложений.

На втором этапе выявляется контекстное окружение информативных предложений в виде информативных фрагментов. Совокупность этих фрагментов является результатом фактографического поиска. Рассмотрим более подробно этапы его реализации.

Информативность $I_{S_{\text{веб}}}^a$ каждого слова a на веб-странице $S_{\text{веб}}$ вычислим по формуле (1), считая, что $I_{S_{\text{веб}}}^a = 0$, если $a \in \text{ПП}'_2$. Силу вербально-ассоциативной связи $I_{S_{\text{веб}}}^{cd}$ для всех пар слов (c, d) вычислим по формуле (2). При этом предполагаем, что $T := S_{\text{веб}}$. Если хотя бы одно из слов c, d не является элементом множества $\text{ПП}'_2$, то полагаем $I_{S_{\text{веб}}}^{cd} = 0$.

Обозначим через $V_{S_{\text{веб}}}^\pi$ множество всех слов, а через $W_{S_{\text{веб}}}^\pi$ – множество всех вербально-ассоциативных пар слов в предложении π на веб-странице $S_{\text{веб}}$. Тогда для вычисления информативности $I_{S_{\text{веб}}}^\pi$ предложения используем формулу, аналогичную выражению (5) из [1]:

$$I_{S_{\text{веб}}}^\pi = \frac{\sum_{a \in V_{S_{\text{веб}}}^\pi, a \in \text{ПП}'_2} I_{S_{\text{веб}}}^a + \sum_{(c, d) \in W_{S_{\text{веб}}}^\pi; c, d \in \text{ПП}'_2} I_{S_{\text{веб}}}^{cd}}{\sqrt{\sum_{a \in V_{S_{\text{веб}}}^\pi, a \in \text{ПП}'_2} (I_{S_{\text{веб}}}^a)^2 + \sum_{(c, d) \in W_{S_{\text{веб}}}^\pi; c, d \in \text{ПП}'_2} (I_{S_{\text{веб}}}^{cd})^2}}.$$

При формировании контекстного окружения информативных предложений для вычисления силы вербально-ассоциативной связи между предложениями построим формулу на основе выражения для информативности $I_Q^{\pi\rho}$ из [1]:

$$I_{S_{\text{веб}}}^{\pi\rho} = \frac{\sum_{c \in \pi, d \in \rho} I_{S_{\text{веб}}}^{cd}}{\sqrt{\sum_{c \in \pi, d \in \rho} (I_{S_{\text{веб}}}^{cd})^2}},$$

где π и ρ – произвольные предложения на веб-странице $S_{\text{веб}}$.

Информативность сообщения Q , полученного в результате построения контекста информативного предложения, будем вычислять по формуле, аналогичной выражению I_Q^{Sub} [3]:

$$I_{S_{\text{веб}}}^Q = \frac{\sum_{\pi \in Q} I_{S_{\text{веб}}}^\pi}{\sqrt{\sum_{\pi \in Q} (I_{S_{\text{веб}}}^\pi)^2}}.$$

Если в тексте найдено несколько информативных предложений, то при выявлении их контекста каждое из неинформативных предложений v , расположенное между парой информативных, будем относить к контекстному окружению такого предложения μ из этой пары, для которого значение силы связи $I_{S_{\text{веб}}}^{\mu v}$ наибольшее.

2.7. **Оценка тональности интернет-сообщений.** Пусть по-прежнему Ct – тематический корпус текстов с тонально-окрашенной лексикой, состоящий из n (по числу оценок в шкале оценивания) подкорпусов, т. е. $Ct = \{Cp_i \mid i = \overline{1, n}\}$. Каждый подкорпус Cp_i состоит из текстов одинаковой тональности и представляет собой пару $\langle Cp_i, Ev_i \rangle$ (Ev_i – оценка тональности для всех текстов из множества Cp_i). Пусть также T – текстовое сообщение, полученное в результате формирования контекстного окружения к некоторому найденному на веб-странице $S_{\text{веб}}$ информативному предложению. Построим вектор \mathbf{F}_T сообщения T и векторы \mathbf{F}_{Cp_i} ($i = \overline{1, n}$) по аналогии с построением векторов \mathbf{F}_s и $\mathbf{F}_{\text{ПП}'_2}$. Для каждой пары $(\mathbf{F}_T, \mathbf{F}_{Cp_i})$ вычислим косинус угла между этими векторами по формуле

$$\cos(\widehat{\mathbf{F}_T, \mathbf{F}_{Cp_i}}) = \frac{\mathbf{F}_T \mathbf{F}_{Cp_i}}{|\mathbf{F}_T| |\mathbf{F}_{Cp_i}|}, i = \overline{1, n}.$$

Тогда сообщению T будет соответствовать оценка тональности $E\nu_i$ при такой величине i , при которой $\cos(\mathbf{F}_T, \mathbf{F}_{Cp_i})$ принимает наибольшее значение.

3. Формирование отчета по результатам интернет-мониторинга. В отчете по результатам интернет-мониторинга и анализа тонально-окрашенной информации каждому сообщению Q поставим в соответствие адрес веб-страницы, на которой оно расположено, а также его информативность $I_{S_{\text{веб}}}^Q$ и числовое значение оценки тональности $E\nu_i$ в принятой шкале (табл. 4).

Таблица 4. Пример отчета по результатам мониторинга и анализа тонально-окрашенной информации

Тонально-окрашенное сообщение	Адрес веб-страницы	Информативность сообщения (%)	Оценка тональности
Сильные дожди, идущие в Центральной Европе несколько дней, привели к наводнениям в ряде регионов Австрии, Чехии, Германии и Польши. Проводится эвакуация населения, затруднено транспортное сообщение.	http:// vsekommentarii.com/news/2013/06/02/9108226.htm	45	4 (из 10)
Недовольные оправдательным приговором 28 обвиняемым в трагедии в Порт-Саиде подожгли здание Египетского футбольного союза и клуб отдыха офицеров полиции.	http:// rss.novostimira.com/n_2361553.html	64	3 (из 10)

Заключение. Предложенный в статье подход может применяться для анализа тонально-окрашенных сообщений как в Интернете, так и в локальных базах данных. Благодаря использованию тонально-окрашенных тематических корпусов текстов и вербальных ассоциаций достигается универсальность программного обеспечения системы, т. е. его независимость от используемых входных языков и предметной области. Имеется возможность оперативного изменения шкалы оценок тональности анализируемых сообщений.

Литература

1. Липницкий С. Ф. // Информатика. 2011. № 4. С. 21–28.
2. Мартинович Г. А. Вербальные ассоциации и организация лексикона человека // [Электронный ресурс]. Филологические науки. 1989. № 3. С. 39–45. Режим доступа: [http:// lit. lib. ru/m/martinowich_g_a/02assfilnauk.shtml](http://lit.lib.ru/m/martinowich_g_a/02assfilnauk.shtml). Дата доступа: 7.08.2013.
3. Липницкий С. Ф., Мамчич А. А. // Весці НАН Беларусі. Сер. фіз.-тэхн. навук. 2011. № 1. С. 72–81.

S. F. LIPNITSKY

A PITCH ANALYSIS OF TEXT MESSAGES IN INTERNET BASED ON MODELING OF VERBAL ASSOCIATIONS

Summary

The problems of Internet monitoring of tone-colored information were formulated. An approach to search and analysis of tone-colored messages on the basis of verbal associations modeling was proposed. The questions of construction of verbal-associative networks for sentences, text and web pages were considered.