

ЧИСЛЕННЫЕ МЕРЫ “СПЛОЧЕННОСТИ” ИМЕННЫХ ГРУПП

Леонид Леоненко

Аннотация: Обсуждаются понятия, функции и алгоритмы, связанные с численными оценками подобия (или “аналогичности”) текстов натурального языка. В так называемой “теории подобия конечных последовательностей” (ТПКП) подобие двух текстов оценивается посредством длины максимальной общей для этих текстов подпоследовательности суб-текстов (например, слов, предложений, etc.) Если в сравниваемых текстах различные вхождения суб-текстов имеют разную значимость сравнительно с другими вхождениями, при оценке степени подобия учитывается не длина, а суммарный “вес” общей для них подпоследовательности суб-текстов (например, суммарный вес общих вхождений слов).

В данной статье основное внимание уделено мерам оценки структурного сходства текстов. Предполагается, что в тексте-“образце” его суб-тексты сгруппированы в “блоки” (именные группы, предложения, etc.). Принимается следующий принцип “сплоченности”: имена, образующие блок, обычно соседствуют в тексте; и перестановки имен внутри блока разрушают структуру текста в меньшей степени, чем чередование имен, принадлежащих разным блокам.

В статье формулируются понятия, численные меры и алгоритмы, оценивающие степени чередования суб-текстов, принадлежащих разным блокам (меры “сплоченности” блоков). В сочетании с мерами, учитывающими лексическое сходство двух текстов и относительную значимость их сходных суб-текстов, меры сплоченности позволяют адекватно оценивать подобие текстов. Статистические эксперименты показали, что методы ТПКП эффективны в областях автоматической проверки орфографии, идентификации сообщений в телекоммуникационных сетях, компьютерного тестирования знаний. Так, в последней области для компьютера оказывается возможным игнорировать несущественные ошибки в ответах тестируемых, учитывать сокращения и синонимы, разрешить, запретить или ограничить перестановки слов, и т.п.

Ключевые слова: аналогия, подобие текстов, алгоритмы оценки подобия, тестирование знаний

ACM Classification Keywords: I.2.6 Artificial Intelligence – Learning – Analogies; K.3.1 Computers and Education - Computer Uses in Education

Введение

Методы оценки подобия последовательностей составляют важный раздел современной информатики (см., напр., [Смит, 2006]). В работах [Леоненко и Поддубный, 1996], [Леоненко, 2002], [Леоненко, 2010] рассматривались различные численные меры подобия конечных последовательностей и приложения этих мер в задачах компьютерного тестирования знаний. Использование мер подобия позволяет оценивать ответы, вводимые в свободной форме на естественном (русском, украинском, английском и др.) языке.

В упомянутых работах рассматривались меры подобия, связанные главным образом с 1) «элементным» составом последовательностей (в частности, для текстов – с их лексическим составом); и 2) с учетом (неучетом) *линейного* порядка элементов в составе последовательности.

Очевидно, что при оценке подобия предложений и текстов необходимо учитывать не только их лексический состав, но и структуру. В этой статье я рассмотрю численные меры подобия,

предполагающие более «тонкую» структуру текста, чем просто линейный порядок слов. А именно, будут предложены алгоритмы, оценивающие степень чередования, "перемешивания" слов из разных "блоков" текста (именных групп, предложений, etc.).

В первом приближении под «текстом», как и в работе [Леоненко и Поддубный, 1996], понимается иерархически организованная система элементов произвольной природы, для которой значимым является линейный порядок на каждом уровне иерархии. Но теперь он будет *не единственно* значимым отношением на указанном уровне.

Будем считать, что в тексте по каким-либо основаниям выделяются особые **группы имен** (например, в русском языке это могут быть группа подлежащего и группа сказуемого в предложении; в алгоритмическом языке – группы, относящиеся к различным операторам программы, etc.). Принимается (сравн.: [Гладкий, 1973], [Добров и др., 2004]) следующий **«принцип сплоченности» для групп**:

СП (Cohesion Principle): Имена, образующие группу, обычно соседствуют в тексте; и перестановки имен внутри групп «разрушают» структуру текста в меньшей степени, чем чередование имен, принадлежащих разным группам.

Проблема, как именно выделять такие группы имен в тексте (в предложении натурального языка, или в конструкции языка программирования, или в формуле, etc.) обсуждается в [Леоненко, 2008].

Ниже предполагается, что заданы два текста **A** и **B**, и в тексте **A** выделены «сплоченные» группы имен. (Когда речь идет о компьютерном тестировании знаний, в качестве **A** выступает *правильный* ответ на тестовый вопрос, а в качестве **B** – *фактический* ответ учащегося). Необходимо оценить подобие, или «аналогичность», текстов **A** и **B**. При этом текст **A** – я буду называть *моделью*, а текст **B** – *прототипом* аналогии (следуя терминологии, принятой в [Уемов, 1971]). Предлагаемые ниже численные меры аналогичности могут применяться к произвольным текстам описанной выше структуры, *независимо* от способа выделения групп имен в тексте.

Меры "разобщенности" элементов заданной группы в тексте

Пусть даны два текста **A** и **B** – модель и прототип; и для модели **A** задано множество групп имен $\{G_i\}$ $i \in [1, n]$. Говоря неформально, для оценки подобия **A** и **B** будут проверяться следующие главные условия:

- 1) имеют ли **A** и **B** сходный лексический состав;
- 2) в какой степени группы имен прототипа **B**, соответствующие группам $\{G_i\}$ из **A**, «засорены» чужеродными именами.

Аналогично понятию «группы имен в тексте» можно ввести понятие «группы символов в слове» (например, группы букв, относящихся к корню слова в этнических языках, etc.). Всюду ниже термины "символ", "буква" и "слово" будут использоваться как *относительные*: "буква" может пониматься как слово в "слове следующего уровня" (т.е. в предложении), и.т.п.

Численные меры, оценивающие условие 2, базируются на следующем понятии расстояния между символами a и b из группы G_i в слове **W**:

Определение 1. Мерой «разобщенности» символов a и b в слове **W** называется число $\rho(a, b)$ элементов e_k , расположенных в **W** между a и b , таких что $e_k \in U \setminus G_0 \setminus G_i$.

Здесь **U** – некоторое «универсальное» множество символов, а G_0 – группа «нейтральных» символов, не влияющих на разобщенность.

Теорема 1. Мера разобщенности $\rho(a, b)$ является расстоянием на подмножестве элементов множества G_i , входящих в слово **W**.

Пусть $G_i[W]$ – последовательность всех символов группы G_i в слове W . Пусть длина $G_i[W]$ равна $v \geq 1$, а M – матрица размерности $v \times v$, элементы которой – разобщенности $\rho(a, b)$ всех элементов множества $G_i[W]$.

Определение 2. Средней разобщенностью $\mu(G_i[W])$ множества $G_i[W]$ называется число 0 при $v=1$, а при $v > 1$ – число:

$$\mu(G_i[W]) = \frac{\sum_{q=1}^v \sum_{r=1}^v M_{q,r}}{v(v-1)}$$

Определение 3. Спектром разобщенности элементов группы G_i в слове W назовем последовательность натуральных чисел $\langle L_1, \dots, L_{v-1} \rangle$, где $L_k = \rho(a_k, a_{k+1})$, $k \in [1, v-1]$.

Теорема 2.

$$\mu(G_i[W]) = \frac{2}{v(v-1)} \sum_{j=1}^{v-1} j(v-j)L_j$$

Можно получить различные оценки зависимости меры $\mu(G_i[W])$ от длины v последовательности $G_i[W]$, свойств спектра разобщенности, и других параметров. Простейшими являются оценки:

Оценка 1. Пусть $L(G_i[W])$ – спектр разобщенности, и пусть $minL$ и $maxL$ – соответственно минимальный и максимальный элементы $L(G_i[W])$. Тогда $minL \times \frac{v+1}{3} \leq \mu(G_i[W]) \leq maxL \times \frac{v+1}{3}$.

Оценка 2. Обозначим через p сумму элементов спектра $L(G_i[W])$, т.е. число всех «чужеродных» для G_i символов между a_1 и a_v в слове W . Тогда $\frac{2p}{v} \leq \mu(G_i[W]) \leq \frac{pv}{2(v-1)}$.

Эти оценки можно улучшать, если принимать те или иные гипотезы о структуре разобщающей a_1 и a_v подпоследовательности «чужеродных» для G_i символов.

Циклические "разобщенности"

Кроме меры $\mu(G_i[W])$, можно ввести меру $\mu_0(G_i[W])$ так называемой «циклической» разобщенности символов a и b из группы G_i в слове W . Неформально циклическая разобщенность – это «обычная» разобщенность букв в слове, полученном из W «сшиванием» его начала и конца (или «скручиванием W в кольцо»). Циклическая разобщенность применяется, главным образом, не к словам, а к предложениям и текстам (и позволяет учесть, например, те случаи, когда в предложении русского языка возможен перенос слов из начала в конец с сохранением смысла).

Пусть множества U , G_i и последовательность W определены как в предыдущем разделе; и пусть $\rho(a, b)$ – «обычная» разобщенность букв a и b из множества G_i в слове W . Пусть, далее, α , β и γ – подслова слова W (возможно, пустые), такие, что

$$W = \langle \alpha a \gamma b \beta \rangle$$

Пусть, наконец, α' и β' – инверсии слов α и β ; W' – слово, равное

$$W' = \langle a \alpha' \beta' b \rangle$$

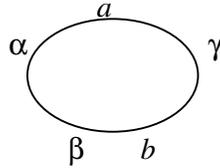
и $\rho'(a, b)$ – «обычная» разобщенность букв a и b в слове W' .

Определение 4. Мерой циклической разобщенности букв a и b в слове W будем называть число

$$\rho_0(a, b) = \min\{\rho(a, b), \rho'(a, b)\}$$

Например, если $c \notin G_i$, $c \notin G_0$, $d \in G_0$, то циклическая разобщенность a и b из G_i в словах $W_1 = \langle acdb \rangle$ и $W_2 = \langle accb \rangle$ равна 0 (так как здесь под слова α и β пусты); в словах $W_3 = \langle acbc \rangle$ и $W_4 = \langle cacb \rangle$ она равна 1; а в слове $W_5 = \langle caccbc \rangle$ – равна 2.

Неформальным оправданием термина “циклический” и определения 4 служит схема



из которой видно, в каком смысле инверсии слов α и β располагаются, подобно слову γ , “между” буквами a и b . Формально уместность указанного термина демонстрируется теоремой 3 (см. ниже).

Определение 5. Пусть $W = \langle e_1, \dots, e_i, a_1, \dots, a_\nu, e_j, \dots, e_n \rangle$, где $a_1 \in G_i$, $a_\nu \in G_i$, причем между a_1 и a_ν расположены все элементы G_i , входящие в слово W . Пусть, далее, W_1 – слово, полученное из W отбрасыванием элементов e_1, \dots, e_i и e_j, \dots, e_n ; а знак \oplus обозначает оператор конкатенации последовательностей. *Спектром циклической разобщенности* элементов множества G_i в слове W назовем последовательность натуральных чисел

$$\text{Lo}_{G_i[W]} = \{L_0\} \oplus \text{Lo}_{G_i[W_1]} \oplus \{L_\nu\} = \langle L_0, L_1, \dots, L_{\nu-1}, L_\nu \rangle$$

где $\text{Lo}_{G_i[W_1]}$ – спектр “обычной” разобщенности для G_i в слове W_1 ; L_0 и L_ν – количества тех букв e_k среди соответственно e_1, \dots, e_i и e_j, \dots, e_n , которые не принадлежат множеству G_0 .

Определение 6. В случае циклической разобщенности *длиной P разобщающей последовательности* для множества G_i в слове W будем называть сумму элементов спектра циклической разобщенности $\text{Lo}_{G_i[W]}$.

Очевидно, что если p – сумма элементов спектра обычной разобщенности, то $P = p + L_0 + L_\nu$. Следующая лемма также очевидна:

Лемма 1. Если $\text{Lo}_{G_i[W]} = \langle L_0, L_1, \dots, L_{\nu-1}, L_\nu \rangle$ – спектр циклической разобщенности, то для любых натуральных $k \geq 1$ и $t \in [k+1, \nu-k]$

$$\rho_0(a_k, a_t) = \min \left\{ \sum_{j=k}^{t-1} L_j, \sum_{j=0}^{k-1} L_j + \sum_{j=t+1}^{\nu} L_j \right\} = \min \{ \rho(a_k, a_t), P - \rho(a_k, a_t) \}$$

Оценка 3 (следствие леммы 1). При принятых выше обозначениях $\rho_0(a_k, a_t) \leq 0,5 * P$.

Теорема 3. Пусть слово W' – циклическая перестановка слова W , т.е. $W = \langle e_1, \dots, e_k, e_{k+1}, \dots, e_n \rangle$ и $W' = \langle e_{k+1}, \dots, e_n, e_1, \dots, e_k \rangle$. Пусть, далее, вхождения a и b некоторых букв слова W принадлежат множеству G_i , а a' и b' – вхождения в слово W' , в которые перешли a и b в результате циклической перестановки W . Тогда $\rho_0(a', b') = \rho_0(a, b)$.

Теорема 4. Мера циклической разобщенности $\rho_0(a, b)$ является расстоянием на подмножестве элементов множества G_i , входящих в последовательность W .

Определение 7. Пусть $\nu > 1$. *Матрицей циклической разобщенности* множества $G_i[W] = \{a_1, \dots, a_\nu\}$ назовем матрицу M^0 размерности $\nu \times \nu$, элементы которой $M_{q,r}^0$ суть $\rho_0(a_q, a_r)$ в последовательности W .

Определение 8. Средней циклической разобщенностью $\mu_0(\mathbf{G}_i[\mathbf{W}])$ множества $\mathbf{G}_i[\mathbf{W}]$ назовем число 0 при $v=1$, а при $v>1$ – число:

$$\mu_0(\mathbf{G}_i[\mathbf{W}]) = \frac{\sum_{q=1}^v \sum_{r=1}^v M_{q,r}^{\circ}}{v(v-1)}$$

где M° – матрица циклической разобщенности множества $\mathbf{G}_i[\mathbf{W}]$.

Следствием теоремы 4 является равенство (при $v>1$)

$$\mu_0(\mathbf{G}_i[\mathbf{W}]) = 2 \times \frac{\sum_{q=1}^{v-1} \sum_{r=q+1}^v M_{q,r}^{\circ}}{v(v-1)}$$

Из теоремы 3 следует: если слово \mathbf{W}' – циклическая перестановка слова \mathbf{W} , то $\mu_0(\mathbf{G}_i[\mathbf{W}]) = \mu_0(\mathbf{G}_i[\mathbf{W}'])$.

Теорема 5. Для любого слова \mathbf{W} имеет место неравенство $\mu_0(\mathbf{G}_i[\mathbf{W}]) \leq \mu(\mathbf{G}_i[\mathbf{W}])$.

Оценка 4. Пусть $minL$ и $maxL$ – соответственно минимальный и максимальный элементы спектра циклической разобщенности $L \oslash \mathbf{G}_i[\mathbf{W}]$. Тогда

$$minL \times \frac{v+1}{4} \leq \mu_0(\mathbf{G}_i[\mathbf{W}]) \leq maxL \times \frac{1}{4} \left(v+1 + \frac{1}{v-1} \right)$$

Оценки 3 и 4 можно улучшать, если принимать те или иные гипотезы о структуре спектра циклической разобщенности.

Понятие " φ -сплоченности"

Определение 9. Пусть φ – любая функция со следующими свойствами: 1) $\varphi(x)$ определена для любого $x \geq 0$; 2) $\varphi(0) = 1$; 3) $\varphi(x)$ невозрастает при $x \rightarrow \infty$. Тогда φ -сплоченностью элементов множества \mathbf{G}_i в слове \mathbf{W} назовем величину

$$C(\varphi, \mathbf{G}_i, \mathbf{W}) = \varphi(\mu(\mathbf{G}_i[\mathbf{W}])) ,$$

а циклической φ -сплоченностью – величину

$$C_0(\varphi, \mathbf{G}_i, \mathbf{W}) = \varphi(\mu_0(\mathbf{G}_i[\mathbf{W}]))$$

Ввиду свойства 3 функции φ и теоремы 5 предыдущего раздела $C_0(\varphi, \mathbf{G}_i, \mathbf{W}) \geq C(\varphi, \mathbf{G}_i, \mathbf{W})$.

Если трактовать слово \mathbf{W} как *предложение* (или иное корректное выражение) некоторого (натурального или искусственного) языка; буквы \mathbf{W} как *слова* этого выражения; а множества \mathbf{G}_i – как некоторые *части* \mathbf{W} (например, именные группы, входящие в \mathbf{W}); то средние разобщенности μ и μ_0 являются оценками “сплошности” размещения слов части \mathbf{G}_i в выражении \mathbf{W} . Циклическая либо нециклическая разобщенность используется в зависимости от того, допускает ли ситуация перенос части слов \mathbf{G}_i из начала в конец выражения без ущерба для “сплошности” \mathbf{G}_i . Например, для следующего предложения (из «Ночи перед Рождеством» Н.В. Гоголя):

\mathbf{W} = *Ведьма, увидевши себя вдруг в темноте, вскрикнула*

такие переносы допустимы, но лишь ограниченно: предложение

Вскрикнула ведьма, увидевши себя вдруг в темноте

вероятно, можно считать эквивалентным W ; однако

В темноте вскрикнула ведьма, увидевши себя вдруг

– уже нет.

Тем не менее, иногда разумно применять именно циклическую меру разобшенности. Так, в работе [Добров и др., 2004, с.66] обсуждается проблема распознавания в техническом тексте понятий, описанных формальной онтологией соответствующей предметной области. Авторы пишут: «распознавание “разорванных” терминов в тексте представляет достаточно сложную проблему», и приводят следующий пример: «...для понятия *вертикальный маневр* существует синоним *маневр в вертикальной плоскости*, который мог бы встретиться в предложении вида “*Маневр осуществляется самолетом... в вертикальной плоскости*”». Очевидно, что здесь циклическая сплоченность упомянутого термина, в отличие от нециклической, не нарушена.

Что касается φ -сплоченности, то она будет использоваться как особая мера сходства таких последовательностей слов из G_i , которые в различной степени “засорены чужеродными словами”. Этот тип сходства можно сопоставить, не претендуя на полноту аналогии, с понятием “узнаваемость”. Условие 2), которому удовлетворяет φ , означает, что когда между словами из G_i в W нет слов, принадлежащих другим частям W (кроме “нейтральных” слов из G_0), узнаваемость части G_i в W равна 1. В противном случае она может быть (не обязательно) меньше 1. Конкретный вид функции φ может выбираться в зависимости от принимаемой гипотезы, насколько узнаваемость части G_i искажается от ее “засорения”. Если, допустим, принимается, что даже одно “разобщающее” слово сильно портит узнаваемость G_i , рационально выбрать $\varphi = e^{-x}$. Если же считать, что G_i остается хотя бы “наполовину” узнаваемым, сколько бы чужеродных слов его не перемежали, можно выбрать $\varphi = (1+x)/(1+2x)$, пределом которой при $x \rightarrow \infty$ будет 1/2. Наконец, если принять, что при любой степени “засорения” узнаваемость снижается ровно вдвое, можно положить $\varphi(x)=1$ при $x=0$ и $\varphi(x)=0,5$ при $x>0$.

О применении мер φ -сплоченности

Приведем примеры применения рассмотренных мер к оценкам аналогичности структур различных предложений русского языка, сходных по лексическому составу. Предметная область – компьютерное тестирование знаний. Пусть «эталонным» ответом на вопрос

Каково наиболее значительное достижение в области математики конца XVII в.; и кто его автор?

считается предложение

Ньютон и Лейбниц открыли математический анализ.

Это предложение трактуется как модель вывода по аналогии; и в нем выделяются именные группы: G_1 – множество слов, подобных одному из слов {*Ньютон, Лейбниц*}; G_2 – множество слов, подобных для {*изобрел, открыл, придумал, создал, автор*}; G_3 – множество слов, подобных для {*математический, анализ, матанализ*}. Подобие трактуется в смысле [Леоненко и Поддубный, 1996]: задается некий численный уровень подобия, достаточный, чтобы, например, слова *Ляйбниц* или *Лейбницец* считались эквивалентными слову *Лейбниц*.

Пусть в качестве других возможных ответов на поставленный вопрос компьютер должен оценить предложения:

1. *Лейбниц и Ньютон – изобретатели матанализа.*
2. *Мат. анализ был создан Ньютоном и Лейбницец.*
3. *Матанализ придуман Лейбницец и Ньютоном.*

4. Лейбниц, а также Ньютон, открыли анализ.
5. Ньютон открыл матанализ, и Лейбниц также.
6. Мат. анализ Ньютон открыл, и Лейбниц.
7. Матанализ, Лейбниц, Ньютон.
8. Лейбниц придумал матанализ.
9. Анализ Ньютона математический Лейбниц придумал.
10. Математика Ньютона открыла Лейбницева анализ.
11. Ньютонов анализ создал математику Лейбница.
12. Открыт и Ньютон анализом Лейбница.
13. Ньютон придумал Лейбница и создал анализ.
14. Анализ изобретений Ньютона и Лейбница открыл мат.
15. Ньютон открыл математику, Лейбниц – анализ 'и'.

Обозначим эти шестнадцать предложений через W_0 (эталон), W_1, \dots, W_{15} . Видимо, не встретит возражений утверждение, что $W_1 - W_5$ являются приемлемыми ответами. Ответ W_6 стилистически неудачен, но верен. Ответ W_7 предельно лапидарен и верен. Ответ W_8 неполон; можно спорить, следует ли “засчитывать” его. При устном опросе преподаватель, возможно, тем или иным способом побудил бы студента дополнить ответ W_8 . Мы будем считать W_8 “пограничным” между допустимыми и недопустимыми ответами. Что касается ответа W_9 , то его можно истолковать как двусмысленный (или, если угодно, косноязычный) – и потому также отнести к “пограничным” ответам. Либеральный экзаменатор примет его, ригорист – нет (но и ригорист задумается, если ответ W_9 поступил от иностранца, плохо владеющего русским).

Ответы же $W_{10} - W_{15}$, по-видимому, являются совершенно неприемлемыми. Но все же “качество” их неприемлемости разное. Ответ W_{10} выражает – правильным русским языком – вполне ложное суждение. Но представим, что его грамматическая правильность лишь кажущаяся, поскольку W_{10} (как выше W_9) высказан иностранцем, не владеющим падежами. В этой ситуации W_{10} становится весьма похожим на W_9 , хотя все же менее приемлемым.

Постоянная “гипотеза иностранца” может показаться нарочитой и неоправданной. Однако ясно, что те или иные искажения могут появиться в содержательно верном ответе по разным причинам. Представим себе, например, вполне русскоязычного студента, который переименовывает слова в ответе, желая посмеяться над компьютерной тестирующей системой. Если конструкторы тестирующей системы хотят предусмотреть какую-то реакцию на подобные действия студентов, они могут зафиксировать степень допустимых искажений ответов. Теория подобия текстов [Леоненко и Поддубный, 1996] позволяет задать пределы буквенных искажений слов фразы. В нашем же случае мы хотим учесть те искажения, которые вносятся *перемещением слов, принадлежащих одним именным группам, внутрь других именных групп* предложения-ответа.

Допустим поэтому, что: **(A1)** Во всех рассматриваемых предложениях $W_0 - W_{15}$ падежные окончания, а также иные “незначительные” искажения слов (вроде *Лейбниц – Ляйбниц*) не следует принимать во внимание; **(A2)** Оттенки смысла слов-синонимов вроде *открыл, изобрел и создал* (эти оттенки в той или иной мере существенны при трактовке $W_0 - W_{15}$ как грамматически правильных фраз русского языка), также игнорируются. Наконец, учтем, что при вводе ответа с клавиатуры некоторые буквы слов могут быть непреднамеренно пропущены – и это же относится к коротким, особенно однобуквенным, словам. Поэтому примем: **(A3)** Пропуск/добавление одного-двух однобуквенных слов не должны сильно влиять на оценку ответа.

Ясно, что допущения (A1)–(A3) не следует трактовать как “общезначимые”. Вполне уместны (и неоднократно применялись) подходы, при которых, скажем, от тех или иных падежных окончаний слов существенно зависит способ распознавания смысла предложения (см., напр., [Гладун, 1987, с.100-119], [Добров и др., 2004, с.59]). Но при анализе предложения на “сплоченность” его именных групп принятие (A1)–(A3) оправданно, особенно если применять такой анализ *независимо* от других способов оценки предложения. Например, можно *сначала* отобрать приемлемый ответ на основе лексического подобия фраз, *не* предполагая при этом (A1)–(A3), и лишь *затем* оценить степень сплоченности этого ответа с учетом (A1)–(A3). Именно это предполагает метод оценки подобия текстов, описанный ниже.

При допущениях (A1)–(A3) предложение W_{11} : *Ньютонов анализ создал математику Лейбница* – может быть расценено как искаженное *Ньютон анализ создал математический, и Лейбниц*, что почти совпадает с W_6 или W_9 .

Таким образом, если принимать (A1)–(A3), то ответы W_{10} , W_{11} (и, как легко видеть, W_{12}) нужно оценить как хотя и неприемлемые, но все же находящиеся от множества “пограничных” ответов “на меньшем расстоянии”, чем, ответы W_{13} – W_{15} .

Нетрудно убедиться в том, что мера широкого невзвешенного подобия F_0 [Леоненко и Поддубный, 1996], оценивающая подобие фактического и эталонного ответа без учета порядка следования слов, припишет предложениям W_1, \dots, W_{15} следующие степени сходства с W_0 (с точностью до сотых):

$$\begin{aligned} W_1 - 1.0; & \quad W_2 - 0.86; & \quad W_3 - 1.0; & \quad W_4 - 1.0; & \quad W_5 - 0.83; & \quad W_6 - 1.0; & \quad W_7 - 0.6; \\ W_8 - 0.6; & \quad W_9 - 0.83; & & & & & \\ W_{10} - 0.83; & \quad W_{11} - 0.83; & \quad W_{12} - 1.0; & \quad W_{13} - 0.83; & \quad W_{14} - 0.86; & \quad W_{15} - 1.0. \end{aligned}$$

Таким образом, мера F_0 не отделяет удовлетворительным образом приемлемые ответы $\{W_1, \dots, W_7\}$ от “пограничных” $\{W_8, W_9\}$ и от неприемлемых $\{W_{10}, \dots, W_{15}\}$.

Ситуация практически не улучшается, если применить взвешенное подобие [Leonenko, 2002], [Leonenko, 2010]. Давайте припишем элементам эталона следующие относительные веса (в скобках):

Ньютон(4) и(0) Лейбниц(4) открыли(1) математический(3) анализ(8).

Заметим, что синонимичные словосочетания могут содержать разное количество слов, но их суммарные веса должны быть равны; таким образом, *матанализ* из группы G_3 будет иметь вес 11, равный сумме весов слов *математический* и *анализ*. Кроме того, при задании весов эталона можно указать множество слов, которые, хотя и не рассматриваются как синонимичные словам эталона, при появлении в ответе должны иметь вес 0 – таковы, для данного примера, слова *также, кроме, был, etc.*

Далее применим меру широкого (не учитывающего порядок слов) *взвешенного* подобия, которая принимает ожидаемый вес слов ответа, не вошедших в базис подобия, равным *среднему арифметическому слов* эталона, *имеющих ненулевой вес* (в нашем случае этот ожидаемый вес равен 4). Обозначим эту меру через $V_0(\mathbf{A}, \mathbf{B})$, где \mathbf{A} – текст-модель (эталон), и \mathbf{B} – текст-прототип. V_0 –взвешенные степени подобия пятнадцати предложений-ответов эталону W_0 будут такими:

$$\begin{aligned} W_1 - 1.0; & \quad W_2 - 1.0; & \quad W_3 - 1.0; & \quad W_4 - 0.85; & \quad W_5 - 1.0; & \quad W_6 - 1.0; & \quad W_7 - 0.95; \\ W_8 - 0.8; & \quad W_9 - 1.0; & & & & & \\ W_{10} - 1.0; & \quad W_{11} - 1.0; & \quad W_{12} - 0.85; & \quad W_{13} - 0.81; & \quad W_{14} - 0.83; & \quad W_{15} - 1.0. \end{aligned}$$

(Здесь при подсчете степеней для W_{13} и W_{14} учтено, что каждое из них содержит *одно* слово – например, для W_{14} им может считаться любое из слов *изобрел* или *открыл*, а также подобные им слова, – не входящее в базис подобия).

Улучшения, вносимые мерой V_0 сравнительно с F_0 , связаны только с относительным повышением веса ответа W_7 по сравнению с W_8 . Вместе с тем классы приемлемых, неприемлемых и “пограничных”

ответов по-прежнему не отделены.

Попробуем получить искомое отделение классов, используя понятие φ -сплоченности. Вычислим нециклические φ -сплоченности каждого из указанных шестнадцати предложений W_0, \dots, W_{15} , для каждой из выделенных в W_0 именных групп G_1, G_2 и G_3 , используя в качестве φ функцию $\varphi=1/(1+x)$. Мы получим, в частности: $C(\varphi, G_1, W_0) = 1$; $C(\varphi, G_1, W_1) = 1$; $C(\varphi, G_1, W_5) = 0.333$; и т.д.

В теории подобию текстов *базисом подобия* текстов A и B называется, если опустить детали, то мультимножество (для широкого подобию) или кортеж (для узкого подобию) их общих вхождений суб-текстов, на основании которого оценивается степень подобию A и B . Для оценки невзвешенного подобию важен размер (мощность) базиса, а для взвешенного – суммарный вес элементов базиса.

Пусть для любого предложения W_k с $k \in [1, 15]$ через $B(W_k)$ обозначен *базис широкого взвешенного подобию* фразы W_k относительно эталона W_0 . Скорректируем правило для вычисления степени широкого взвешенного подобию $V_0(W_0, W_k)$ следующим образом:

Если в базис подобию $B(W_k)$ входит группа слов из множества G_i , то при подсчете веса базиса суммарный вес этой группы следует умножить на φ -сплоченность множества G_i в предложении W_k .

Численную оценку подобию двух фраз, полученную в результате применения приведенного выше правила, учитывающего φ -сплоченности групп слов во фразе, обозначим через $H_0(A, B)$. Ясно, что в общем случае H_0 зависит от выбранной функции φ и типа сплоченности (циклической либо нециклической).

В нашем примере (нециклическая φ -сплоченность с $\varphi=1/(1+x)$) H_0 -взвешенные степени подобию пятнадцати предложений-ответов эталону W_0 окажутся следующими:

$$\begin{aligned} W_1 - 1.0; & \quad W_2 - 1.0; & \quad W_3 - 1.0; & \quad W_4 - 0.85; & \quad W_5 - 0.73; & \quad W_6 - 0.8; & \quad W_7 - 0.95; \\ W_8 - 0.8; & \quad W_9 - 0.53; & & & & & \\ W_{10} - 0.39; & \quad W_{11} - 0.43; & \quad W_{12} - 0.65; & \quad W_{13} - 0.63; & \quad W_{14} - 0.53; & \quad W_{15} - 0.46. \end{aligned}$$

Мы видим, что одновременный учет взвешенного лексического подобию и нециклической φ -сплоченности позволяют хорошо отделить классы допустимых $\{W_0, \dots, W_7\}$ и недопустимых $\{W_{10}, \dots, W_{15}\}$ ответов. Скажем, если выбрать уровень подобию, на котором ответ признается приемлемым, равным 0.7; то все содержательно правильные ответы будут признаны компьютером допустимыми, а содержательно неправильные – недопустимыми. При этом “пограничный” ответ W_8 будет отнесен к классу допустимых, а “пограничный” ответ W_9 признан недопустимым – что, по-видимому, приемлемо.

Если вместо нециклической использовать *циклическую* φ -сплоченность, то для предложения W_5 степень подобию эталону W_0 окажется равной 1. Но нетрудно показать, что в целом для совокупности предложений W_0, \dots, W_{15} циклические φ -сплоченности оказываются непригодными для различения классов допустимых и недопустимых ответов в множестве W_1, \dots, W_{15} (причем при различных выборах функции φ). Скажем, степени подобию эталону W_0 для ответов W_6 и W_{10} получаются равными; степень подобию W_6 меньше, чем у W_{14} , и т.п.

Заключение

Принцип сплоченности **СП** является одним из структурных принципов, нарушения которого снижают «целостность» восприятия текста человеком. Эта целостность обусловлена семантическими связями терминов текста (как для натуральных, так и для искусственных языков). Тем не менее сам принцип **СП** носит чисто синтаксический характер. В этом он сходен с принципом взвешенного лексического подобию.

Проведенные статистические эксперименты [Баранов, 2004] показали, что в определенных областях, – в частности, в сфере компьютерного тестирования знаний, – алгоритмы сравнения текстов, базирующиеся на описанных в данной статье мерах подобия, позволяют получить адекватные оценки подобия текстов. Оценки «открытых» ответов, выставляемые системой тестирования знаний CONTROL (использующей упомянутые алгоритмы), статистически не отличаются от оценок специалистов-преподавателей. Это дает возможность *игнорировать разного типа несущественные ошибки и искажения* ответов; и как следствие позволяет сделать «открытые» ответы равноправными с другими типами ответов при компьютерном тестировании (см. [Леоненко, 2010]).

Благодарности

Работа опубликована при финансовой поддержке проекта **ITHEA XXI** Института информационных теорий и приложений FOI ITHEA Болгария www.ithea.org и Ассоциации создателей и пользователей интеллектуальных систем ADUIS Украина www.aduis.com.ua.

Библиография

- [Leonenko, 2002] L. Leonenko. Analogical inferences in computer assisted knowledge testing systems // 6th Multi-Conference on Systemics, Cybernetics and Informatics. – Orlando, Florida, USA, 2002. Proc., Vol. XVIII, pp.371-376.
- [Leonenko, 2010] L. Leonenko. Analogies between Texts: Mathematical Models and Applications in Computer-assisted Knowledge Testing // Information Models of Knowledge. – Kiev, Ukraine – Sofia, Bulgaria: ITHEA, 2010, pp. 128 – 134.
- [Баранов, 2004] В. Ю. Баранов. Комп'ютерне тестування з інформатики: підсумки педагогічного експерименту в Одеській національній академії зв'язку // Теорія та методика навчання математики, фізики, інформатики. – Вип. 4. – Т. 3 – Кр.Ріг.: НМетАУ, 2004. – с. 6 –12.
- [Гладкий, 1973] А.В. Гладкий. Математические методы изучения естественных языков // Труды МИАН им. В.А.Стеклова. – 1973. – Том 133. – с. 95 – 108.
- [Гладун, 1987] В.П. Гладун. Планирование решений. – К.: Наукова думка, 1987. – 168 с.
- [Добров и др., 2004] Б.В. Добров, Н.В. Лукашевич, О.А. Невзорова, Б.Е. Федун. Методы и средства автоматизированного проектирования прикладной онтологии // Изв. РАН. Теория и системы управления. – 2004. – № 2. – с. 58–68.
- [Леоненко, 2008] Л.Л. Леоненко. Язык тернарного описания в оценках связности текстов // Сб. трудов VIII международной конф. "Интеллектуальный анализ информации". – К.: Просвіта, 2008, с. 286–295.
- [Леоненко, 2010b] Л. Леоненко. "Открытые" ответы в компьютерном тестировании знаний // Information Models of Knowledge. – Kiev, Ukraine – Sofia, Bulgaria: ITHEA, 2010, pp. 355 – 361.
- [Леоненко и Поддубный, 1996] Л.Л. Леоненко, Г.В. Поддубный. Теория подобия конечных последовательностей и ее приложения к распознаванию образов // Автоматика и телемеханика, 1996, № 8, с.119-131.
- [Смит, 2006] У. Смит. Методы и алгоритмы вычислений на строках. – М.: "И.Д. Вильямс", 2006. — 496 с.
- [Уемов, 1971] А.И. Уемов. Логические основы метода моделирования. М.: Мысль, 1971. – 311 с.

Информация об авторе



Леонид Леоненко – Одесская национальная академия связи им. А. С. Попова, доцент, ул. Конная, 22, кв. 6, Одесса, 65029, Украина; e-mail: Leonid.Leonenko@gmail.com
 Основные области научной деятельности: неклассическая логика, общая теория систем, компьютерное тестирование знаний