

Формирование единиц представления предметных знаний в задаче их оценки на основе открытых тестов

Емельянов Г. М., Михайлов Д. В., Козлов А. П.

Новгородский государственный университет
имени Ярослава Мудрого

10-я Международная конференция
«Интеллектуализация обработки информации» (ИОИ-2014),

4–11 октября 2014 г.

о. Крит, Греция

Предмет исследования

Методы и алгоритмы формирования знаний о синонимии.

Исследуемая проблема

Передача знаний, представляемых текстами на естественном языке (ЕЯ), между его носителями (экспертами и обучаемыми).

Основная цель исследований

Разработка и теоретическое обоснование методов и алгоритмов поиска оптимального плана передачи смысла между экспертами и обучаемыми в системе контроля знаний с применением открытых тестов.

Определение 2

Ситуация языкового употребления (СЯУ) — описание нового социального опыта (содержания совместных действий) средствами заданного ЕЯ.

Фиксируемый СЯУ S языковой контекст представляется тройкой:

$$S = (O, R, Ts), \quad (1)$$

где O — множество символов, обозначающих понятия действительности;

Ts — множество форм описания S в некоторой знаковой системе;

$R \subset O^n$, где $n \in 1, \dots, |O|$.

Замечание

В данной работе элементы множества Ts — семантически эквивалентные (СЭ) ЕЯ-фразы.

Пусть $Synt$ — сюръективная функция, определяемая синтаксисом языка.

Тогда для $\forall Ts_i \in Ts \exists Tr_i: Ts_i = Synt(Tr_i)$, Tr_i — помеченное дерево.

При этом если $O = M \cup V$, $M \cap V \neq \emptyset$, то для $\forall o_j \in M$ найдётся $o_k \in V$ такое, что понятию o_j соответствует дочерний узел с пометкой w_j , а понятию o_k — родительский узел с пометкой w_k в дереве Tr_i .

Представим СЯУ посредством формального контекста (ФК):

$$K = (G, M, I), \quad (2)$$

где $\forall g \in G$ — основа слова, синтаксически подчинённого другому слову из некоторой $Ts_i \in Ts$ в составе тройки (1).

Множество признаков M включает:

- признаки-указания на основы и флексии слов, синтаксически главных по отношению к словам с основами из G ;
- связи «основа–флексия» для синтаксически главного слова;
- комбинации флексий зависимых и главных слов.

Задача формирования оптимального плана передачи смысла

Требуется найти $I \subseteq G \times M$, определяющее фразы $Ts_i \in Ts$ минимальной символьной длины при максимизации числа слов, наиболее употребимых в различных фразах из Ts (с учётом синонимов).

Основные подзадачи:

- выделение буквенных инвариантов слов (основ) в составе фраз из Ts ;
- формирование критерия информативности слов в контексте СЯУ;
- выделение и классификация синтагматических связей, определяемых отношениями из R в составе модели (1).

Имеем:

$$Ts = \left\{ Ts_i : Ts_i = \odot_j w_{ij} \right\},$$

где \odot — операция конкатенации, а w_{ij} представляется последовательностью

$$W_{ij} = Wc_{ij} \odot Wf_{ij},$$

где Wc_{ij} составляют символы неизменной части (**основы**) слова w_{ij} ;

Wf_{ij} — изменяемой (флексивной) части w_{ij} .

Определение 3

Пусть J — множество индексов **основ** слов, составляющих фразы из Ts . Последовательность таких индексов для некоторой $Ts_i \in Ts$ назовём **моделью** её **линейной структуры** (МЛС), $Ls(Ts_i)$.

Пусть LS есть множество моделей линейных структур фраз из Ts на J .

Лемма 1

Индексы $j_1, j_2 \in J$ соответствуют словам-синонимам и могут быть заменены одним индексом из $(\mathbb{N} \setminus J)$, если $\exists \{Ls(Ts_1), Ls(Ts_2)\} \subseteq LS$:

$$Ls(Ts_1) = J_1 \odot \{j_1\} \odot J_2 \text{ и } Ls(Ts_2) = J_1 \odot \{j_2\} \odot J_2.$$

Пусть $h(j, \text{Ls}(Ts_i))$ — позиция индекса j в модели $\text{Ls}(Ts_i)$. Тогда множество синтагматических связей для $\text{Ls}(Ts_i)$ определяется как

$$D : Ts_i \rightarrow \left\{ \left(h(j, \text{Ls}(Ts_i)), h(k, \text{Ls}(Ts_i)) \right) : j \neq k \right\}. \quad (3)$$

При этом пара (j, k) содержательно соответствует одной синтагме.

Замечание

Согласно определению, синтагматические зависимости задаются на множестве флективных частей слов в составе фраз из Ts .

Далее отождествим с основой и флексией слова понятия «префикс» и «суффикс», принятые в информатике.

Определение 4

Пусть $\text{len}(j, k) = |h(j, \text{Ls}(Ts_i)) - h(k, \text{Ls}(Ts_i))|$.

Назовём указанную величину длиной связи, соответствующей паре (j, k) , относительно модели $\text{Ls}(Ts_i)$.

Пусть

LS' — множество моделей линейных структур фраз из определяющих СЯУ, преобразованное заменой индексов согласно Лемме 1;

J' — соответствующее преобразованное индексное множество J .

Введём в рассмотрение абсолютные частоты:

$N(j, LS')$ — встречаемости отдельного индекса в моделях линейных структур из LS' ;

$N((j, k), LS')$ — встречаемости связи (j, k) в моделях линейных структур из LS' независимо от $\text{len}(j, k)$;

$N(\text{len}(j, k), LS')$ — встречаемости связи (j, k) , имеющей длину $\text{len}(j, k)$, в моделях линейных структур из LS' .

Задача

На основе указанных частот найти набор минимальных текстовых единиц, необходимых и достаточных для формирования оптимального плана передачи смысла ситуации S , представляемой тройкой (1).

Пусть X — последовательность упорядоченных по убыванию значений $N(j, LS')$ для всех $j \in J'$.

Введём обозначения для используемых далее функций.

Функция	Возвращаемое значение
$\text{first}(X)$	первый элемент последовательности X
$\text{last}(X)$	последний элемент последовательности X
$\text{lrev}(X)$	исходная последовательность X без последнего элемента
$\text{rest}(X)$	исходная последовательность X без первого элемента

Разобьём последовательность X на кластеры с применением алгоритма, содержательно близкого алгоритмам класса FOREL.

Пусть $\text{mc}(H_i)$ — функция, вычисляющая центр масс кластера H_i .

При этом элементы X принадлежат **одному кластеру**, если

$$\begin{cases} |\text{mc}(X) - \text{first}(X)| < \frac{\text{mc}(X)}{4} \\ |\text{mc}(X) - \text{last}(X)| < \frac{\text{mc}(X)}{4} \end{cases} \quad (4)$$

Функцию, выдающую true/false в зависимости от выполнения условия (4), далее обозначим как $\text{good}(X)$.

Вход: X ; // упорядоченная числовая последовательность

Выход: $H_i, X_p, X_s : X_p \odot H_i \odot X_s = X$; // \odot — операция конкатенации

```
1:  $i := 1$ ;  
2:  $H_i := X$ ;  
3:  $X_p := \emptyset$ ;  
4:  $X_s := \emptyset$ ;  
5: если  $\text{good}(H_i) = \text{true}$  или  $|H_i| = 1$  то  
6:   вернуть  $H_i, X_p$  и  $X_s$ ;  
7: иначе если  $|\text{mc}(H_i) - \text{first}(H_i)| > |\text{mc}(H_i) - \text{last}(H_i)|$  то  
8:    $X_p := \{\text{first}(H_i)\} \odot X_p$ ;  
9:    $H_i := \text{rest}(H_i)$ ;  
10:  перейти к шагу 5;  
11: иначе если  $|\text{mc}(H_i) - \text{first}(H_i)| < |\text{mc}(H_i) - \text{last}(H_i)|$  то  
12:   $X_s := \{\text{last}(H_i)\} \odot X_s$ ;  
13:   $H_i := \text{lrev}(H_i)$ ;  
14:  перейти к шагу 5;  
15: иначе  
16:   $X_s := \{\text{last}(H_i)\} \odot X_s$ ; // Для разбивки исходной последовательности  
17:   $X_p := \{\text{first}(H_i)\} \odot X_p$ ; // на кластеры данный алгоритм применяется  
18:   $Tmp := \text{lrev}(H_i)$ ; // рекурсивно к  $X_p$  и  $X_s$  на его выходе.  
19:   $H_i := \text{rest}(Tmp)$ ; // Указанный процесс продолжается  
20:  перейти к шагу 5; // до тех пор, пока на очередном шаге  $X_p$  и  $X_s$   
21: конец если // не окажутся пустыми.
```

Пусть H_1, \dots, H_r , причём для $\forall i \neq j$ верно то, что

$$H_i \cap H_j = \emptyset, \text{ а } H_1 \odot H_2 \odot \dots \odot H_r = X.$$

Обозначим далее множество $\{j: N(j, LS') \in H_1\}$ как Cl .

Утверждение 1

Смысловый эталон СЯУ, представляемой тройкой (1), определяют те $Ts_i \in Ts$, для которых

$$Cl \cap Ls'(Ts_i) = Cl, \text{ а } |Ls'(Ts_i) \setminus Cl| \rightarrow \min,$$

где $Ls'(Ts_i) \in LS'$.

Данное условие *необходимо, но не достаточно* для отнесения некоторой $Ts_i \in Ts$ к фразам, определяющим смысловый эталон заданной СЯУ.

Следующий шаг — *выделение и кластеризация* связей, представляемых множеством (3), с последующим анализом индексов из $Ls'(Ts_i) \setminus Cl$.

Будем оценивать «силу» связи (независимо от взаимного расположения слов в линейном ряду фразы) посредством следующей *весовой функции*:

$$W_g((j, k), LS') = N((j, k), LS') \times \frac{N((j, k), LS')}{(N(j, LS') - N((j, k), LS')) + N((j, k), LS') + (N(k, LS') - N((j, k), LS'))}. \quad (5)$$

Пусть X^W — упорядоченная по убыванию последовательность значений функции (5) для индексных пар (j, k) , выделенных на моделях из LS' .

Разобьём X^W на кластеры

$$H_1^W, \dots, H_q^W : H_1^W \odot H_2^W \odot \dots \odot H_q^W = X^W$$

с применением алгоритма, представленного на [слайде 9](#).

При этом *связи, максимально значимые* для формирования оптимального плана передачи смысла заданной СЯУ, будут иметь *значения функции (5)*, вошедшие в кластер H_1^W .

Обозначим далее множество индексов в составе указанных связей как Cl_1 (по аналогии с Cl из [Утверждения 1](#)).

Рассмотрим связи, значения функции (5) которых не вошли в H_1^W .

Разобьём их на кластеры (слайд 9) по величине **среднеквадратического отклонения длины связи (СКОДС)** относительно LS' .

По определению СКОДС для пары (j, k) относительно LS' вычисляется по формуле

$$\sigma(\text{len}(j, k), LS') = \sqrt{E(\text{len}^2(j, k), LS') - E^2(\text{len}(j, k), LS')},$$

где $E(\text{len}(j, k), LS')$ — математическое ожидание длины связи,

$$\begin{aligned} E(\text{len}(j, k), LS') &= \sum_i \left(\frac{N(\text{len}_i(j, k), LS')}{N((j, k), LS')} \text{len}_i(j, k) \right) = \\ &= \sum_i \left(p(\text{len}_i(j, k), LS') \text{len}_i(j, k) \right). \end{aligned}$$

Гипотеза

Индекс $j \in J'$, соответствующий вершине, входит в одну из связей кластера наименьших СКОДС и одновременно в связь из некоторого другого кластера по указанной величине. При этом «индекс вершины» не входит в связи со значениями функции (5) из H_1^W .

Пусть Cl_2 — множество кандидатов на роль вершин деревьев фраз из Ts .

Утверждение 2

Смысловый эталон СЯУ, представляемой тройкой (1), определяют те $Ts_i \in Ts$, для которых помимо условия **Утверждения 1** верно то, что

$$\left| \left(Ls'(Ts_i) \setminus Cl \right) \setminus (Cl_1 \cup Cl_2) \right| \rightarrow \min$$

при минимальной длине суффикса для $\forall w_{ij} : \bigodot_j w_{ij} = Ts_i$.

Обозначим множество фраз $Ts_i \in Ts$, отобранных согласно условиям **Утверждений 1** и **2**, как Ts^* . Пусть

$$R_J = \left\{ ((j, k), Dir) : Dir \in \{\leftarrow, \rightarrow\}, \exists Ts_i \in Ts^* : \{j, k\} \subset Ls'(Ts_i) \right\}, \quad (6)$$

причём если $X^W \neq H_1^W$ и $|Ts^*| > 1$, то либо $(\{j, k\} \cap Cl_2) \neq \emptyset$, либо паре (j, k) соответствует связь со значением функции (5) из кластера H_1^W .

При этом связи из R_J задают минимальные семантико-синтаксические текстовые единицы в рамках **оптимального плана передачи смысла СЯУ**.

Выделение основ и флексий для слов в рамках СЯУ : ключевые процедуры и функции алгоритма

- `pref.show` (w_{ij}) возвращает текущее значение префикса слова w_{ij} ;
- `pref.inc` (w_{ij}) увеличивает длину префикса слова w_{ij} на 1;
- `prefs` объединяет словоформы в группы (списки) по сходству префикса, сортируя их при этом по убыванию длины;
- `pref.check` (Prf) для группы словоформ с общим префиксом Prf анализирует частоты (абсолютные) встречаемости букв на разных позициях относительно начала и конца слова.

При этом частота ν_p встречаемости первого слева символа и букв в составе Prf всегда максимальна. Относительно конца слова также производится поиск символов общего суффикса (включаются во флексивную часть) с частотой встречаемости ν_p .

Утверждение 3

Суммарная длина общих префикса и суффикса пары слов здесь должна составлять минимум треть длины слова, а разность длин у пары слов с общим префиксом (независимо от суффикса) всегда меньше половины длины меньшего слова.

Вход: T_s ;

Выход: $P_w = \bigcup_{i=1}^{|T_s|} P_{w_i}$; // $P_{w_i} = \left\{ (W_{c_{ij}}, W_{f_{ij}}) : W_{c_{ij}} \odot W_{f_{ij}} = W_{ij} \right\}$

1: $P_w := \emptyset$; // W_{ij} — последовательность символов слова w_{ij}

2: **для всех** W_{ij} : $\odot_j w_{ij} = T_{s_i}$, где $T_{s_i} \in T_s$

3: $W_{c_{ij}} := \{W_{ij}[1]\}$; $W_{f_{ij}} := \odot_{k=2}^{|W_{ij}|} W_{ij}[k]$;

4: **конец для** // инициализации основ и флексий

5: $\text{prefs}(\text{PrfsTmp})$;

6: **если** $\text{PrfsTmp} = \emptyset$ **то**

7: **вернуть** P_w и выйти из алгоритма;

8: **иначе**

9: **взять** очередной Prf из PrfsTmp ;

10: **если** $\text{pref.check}(\text{Prf}) = \text{true}$ **то**

11: $P_w := P_w \cup \left\{ (\text{Prf}, W_{f_{ij}}(\text{Prf})) \mid \text{pref.show}(w_{ij}) = \text{Prf} \right\}$;

12: $\text{PrfsTmp} := \text{PrfsTmp} \setminus \{\text{Prf}\}$;

13: **перейти к шагу 6**;

14: **иначе**

15: **для всех** w_{ij} : $\text{pref.show}(w_{ij}) = \text{Prf}$

16: $\text{pref.inc}(w_{ij})$;

17: **конец для**

18: **перейти к шагу 5**

19: **конец если**

20: **конец если**

Порядковый номер СЯУ, i	1	2	3	4	5	6
Число СЭ-фраз, задающих СЯУ	56	28	29	30	6	10
Минимальное число слов во фразе	5	8	11	10	10	11
Максимальное число слов во фразе	12	15	16	18	17	14

i Фразы максимальной длины из определяющих СЯУ

- 1 *Нежелательная переподгонка является причиной заниженности средней величины ошибки алгоритма на обучающей выборке.*
- 2 *Тренировочная выборка, на ней проявляется эффект заниженных значений средней ошибки, причиной же является переусложненная модель.*
- 3 *Контрольная выборка, принятие деревом решения на ней будет с большей вероятностью ошибки именно по причине переподгонки.*
- 4 *Оцениваемая частота, с которой алгоритм допускает ошибку на выборке, рассматриваемой как контрольная, может оказаться заниженной по причине переподгонки.*
- 5 *Распознавание обладает таким свойством, что его ошибка будет иметь заниженную оценку при неудачном выборе правила принятия решений.*
- 6 *Рост числа базовых классификаторов, который ведет к практически неограниченному увеличению обобщающей способности композиции алгоритмов.*

Порядковый номер СЯУ	1	2	3	4	5	6
Число связей со значениями функции (5), вошедшими в кластер H_1^W	4	9	14	11	6	13
Число найденных кластеров по СКОДС	5	6	5	7	1	5
Число фраз, представляющих эталон СЯУ	12	7	8	11	2	1
Общее число связей в рамках эталона СЯУ	26	28	39	43	12	19
в том числе истинных	21	17	23	24	10	14
ложных	5	11	16	19	2	5

Замечание

Для каждой найденной связи её *направление* здесь задаётся экспертом, причём только для связей, определенных им как *истинные*.

Совокупные знания системы по синтагматическим связям в рамках отдельной СЯУ могут быть представлены *булевым вектором*

$$(d_1, \dots, d_k, \bar{d}_{k+1}, \dots, \bar{d}_n), \quad (7)$$

где d_1, \dots, d_k соответствуют *истинным*, а $\bar{d}_{k+1}, \dots, \bar{d}_n$ — *ложным* связям.

Программная реализация и результаты экспериментов

Пример: исходное множество семантически эквивалентных фраз

Исходное множество семантически эквивалентных фраз

28:1

Insert

Indent

Нежелательное переобучение приводит к заниженности эмпирического риска.

Нежелательное переобучение, следствием которого является заниженность эмпирического риска.

Заниженность эмпирического риска является следствием нежелательного переобучения.

Заниженность эмпирического риска, являющаяся следствием нежелательного переобучения.

Эмпирический риск, заниженность которого является следствием нежелательного переобучения.

Эмпирический риск, заниженный вследствие нежелательного переобучения.

Эмпирический риск, к заниженности которого ведет нежелательное переобучение.

Риск, заниженный как следствие переобучения.

Эмпирический риск по причине, обусловленной нежелательным переобучением, может оказаться заниженным.

Эмпирический риск в силу обстоятельств, связанных с нежелательным переобучением, может оказаться заниженным.

Эмпирический риск по причине, вызванной нежелательным переобучением, может быть заниженным.

Эмпирический риск, к заниженности которого приводит нежелательное переобучение.

Нежелательное переобучение служит причиной заниженности эмпирического риска.

Заниженность эмпирического риска, причиной которой является нежелательное переобучение.

Заниженность эмпирического риска является результатом нежелательного переобучения.

Нежелательное переобучение, с которым связана заниженность эмпирического риска.

Эмпирический риск, с переобучением связана его заниженность.

Заниженность эмпирического риска связана с переобучением.

Заниженность эмпирического риска, являющаяся результатом нежелательного переобучения.

Нежелательное переобучение, результатом которого является заниженность эмпирического риска.

Нежелательное переобучение, результат которого есть заниженность эмпирического риска.

Нежелательное переобучение, приводящее к заниженности эмпирического риска.

Нежелательное переобучение, служащее причиной заниженности эмпирического риска.

Заниженность эмпирического риска относится к следствию нежелательного переобучения.

Заниженность эмпирического риска связана с нежелательным переобучением.

Нежелательное переобучение является причиной заниженности эмпирического риска.

Заниженность эмпирического риска, причиной которой служит нежелательное переобучение.

Нежелательная переподгонка приводит к заниженности эмпирического риска.

Исходные СЭ-фразы (продолжение) и фраза максимальной длины

Исходное множество семантически эквивалентных фраз

57:1

Insert

Indent

Modified

Заниженность эмпирического риска является следствием нежелательной переподгонки.

Заниженность эмпирического риска, являющаяся следствием нежелательной переподгонки.

Эмпирический риск, заниженность которого является следствием нежелательной переподгонки.

Эмпирический риск, заниженный вследствие нежелательной переподгонки.

Эмпирический риск, к заниженности которого ведет нежелательная переподгонка.

Риск, заниженный как следствие переподгонки.

Эмпирический риск по причине, обусловленной нежелательной переподгонкой, может оказаться заниженным.

Эмпирический риск в силу обстоятельств, связанных с нежелательной переподгонкой, может оказаться заниженным.

Эмпирический риск по причине, вызванной нежелательной переподгонкой, может быть заниженным.

Эмпирический риск, к заниженности которого приводит нежелательная переподгонка.

Нежелательная переподгонка служит причиной заниженности эмпирического риска.

Заниженность эмпирического риска, причиной которой является нежелательная переподгонка.

Заниженность эмпирического риска является результатом нежелательной переподгонки.

Нежелательная переподгонка, с которой связана заниженность эмпирического риска.

Эмпирический риск, с переподгонкой связана его заниженность.

Заниженность эмпирического риска связана с переподгонкой.

Заниженность эмпирического риска, являющаяся результатом нежелательной переподгонки.

Нежелательная переподгонка, результатом которой является заниженность эмпирического риска.

Нежелательная переподгонка, результат которой есть заниженность эмпирического риска.

Нежелательная переподгонка, приводящая к заниженности эмпирического риска.

Нежелательная переподгонка, служащая причиной заниженности эмпирического риска.

Заниженность эмпирического риска относится к следствию нежелательной переподгонки.

Заниженность эмпирического риска связана с нежелательной переподгонкой.

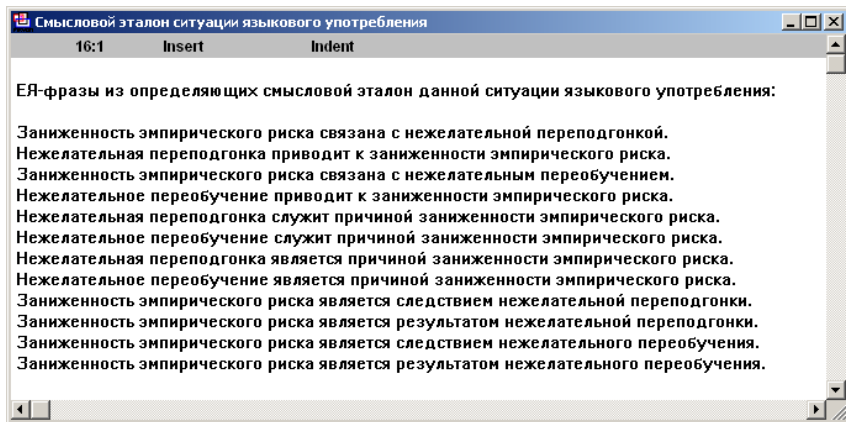
Нежелательная переподгонка является причиной заниженности эмпирического риска.

Заниженность эмпирического риска, причиной которой служит нежелательная переподгонка.

Нежелательная переподгонка является причиной заниженности средней величины ошибки алгоритма

на обучающей выборке.

Нежелательная переподгонка, следствием которой является заниженность эмпирического риска.



Выделенные кандидаты на роль вершин синтаксических деревьев:

« <i>привод/ведет</i> »	« <i>связан</i> »	« <i>с</i> »
« <i>результат/следстви</i> »	« <i>причин</i> »	« <i>к</i> »
« <i>есть/явля/служ</i> »	« <i>котор</i> »	

Связи в рамках смыслового эталона:

заниженн → риск
риск → эмпирическ
заниженн → эмпирическ
переподгонк,переобучени → нежелательн
результат,следстви → нежелательн
есть,явля,служ → результат,следстви
есть,явля,служ → причин
с → нежелательн
результат,следстви → переподгонк,переобучени
с → переподгонк,переобучени
связан → с
привод,ведет → переподгонк,переобучени
привод,ведет → нежелательн
связан → переподгонк,переобучени
есть,явля,служ → переподгонк,переобучени
есть,явля,служ → нежелательн
привод,ведет → к
есть,явля,служ → заниженн
к → заниженн
связан → заниженн
причин → заниженн

Ложные связи:

риск — причин
риск — к
есть,явля,служ — риск
риск — связан
риск — с

Кластеризация по значению весовой функции связи

Связи со значениями **весовой функции** из кластера H_1^W :

Основа для j	Основа для k	Dir(j, k)	Wg($(j, k), LS'$)
заниженн	риск	→	16,0556
эмпирическ	риск	←	15,0588
заниженн	эмпирическ	→	14,2222
нежелательн	переподгонк, переобучени	←	12,5000

Связи со значениями **весовой функции**, **не вошедшими** в кластер H_1^W :

Основа для j	Основа для k	Dir	Wg	$\sigma(\text{len}(j, k), LS')$
нежелательн	результат, следстви	←	1,0000	0,4330
результат, следстви	есть, явля, служ	←	1,1250	0,4714
есть, явля, служ	причин	→	1,1250	0,4714
с	нежелательн	→	0,5294	0,4714
переподгонк, переобучени	результат, следстви	←	1,3889	0,4899
с	переподгонк, переобучени	→	1,3889	0,4899
связан	с	→	5,0000	0,4899
переподгонк, переобучени	привод, ведет	←	0,2222	0,5000
нежелательн	привод, ведет	←	0,2667	0,5000
связан	переподгонк, переобучени	→	1,3889	0,8000
переподгонк, переобучени	есть, явля, служ	←	2,0000	0,8975
нежелательн	есть, явля, служ	←	2,4000	0,8975
привод, ведет	к	→	1,3333	1,0000
есть, явля, служ	заниженн	→	2,0000	1,2583
заниженн	к	←	0,5000	1,4142
заниженн	связан	←	1,3889	1,6733
причин	заниженн	→	1,3889	2,2450

Кластеризация по среднеквадратическому отклонению длины связи

Кластеры, выделенные по значению **среднеквадратического отклонения** длины связи:

№ кластера	1	2	3	4	5
Число связей, вошедших в кластер	36	10	8	5	1
Значение СКОДС для связи					
минимальное	0,0000	0,4330	0,7071	1,2000	2,2450
максимальное	0,0000	0,5000	1,0954	1,6733	2,2450

Связи в рамках эталона СЯУ из группируемых по значению **СКОДС**:

Основа для j	Основа для k	Dir	$\sigma(\text{len}(j, k), LS')$	№ кластера
нежелательн	результат, следстви	←	0,4330	2
результат, следстви	есть, явля, служ	←	0,4714	2
есть, явля, служ	причин	→	0,4714	2
с	нежелательн	→	0,4714	2
переподгонк, переобучени	результат, следстви	←	0,4899	2
с	переподгонк, переобучени	→	0,4899	2
связан	с	→	0,4899	2
переподгонк, переобучени	привод, ведет	←	0,5000	2
нежелательн	привод, ведет	←	0,5000	2
связан	переподгонк, переобучени	→	0,8000	3
переподгонк, переобучени	есть, явля, служ	←	0,8975	3
нежелательн	есть, явля, служ	←	0,8975	3
привод, ведет	к	→	1,0000	3
есть, явля, служ	заниженн	→	1,2583	4
заниженн	к	←	1,4142	4
заниженн	связан	←	1,6733	4
причин	заниженн	→	2,2450	5

№ СЯУ	1	2	3	4	5	6
l_1	56	28	29	30	6	10
n_1	12	15	16	18	17	14
l_2	12	7	8	11	2	1
n_2	7	10	12	13	10	13
$\text{vol}(n_1)$	$4,790 \cdot 10^8$	$1,308 \cdot 10^{12}$	$2,092 \cdot 10^{13}$	$6,402 \cdot 10^{15}$	$3,557 \cdot 10^{14}$	$8,718 \cdot 10^{10}$
vol_1	672	420	464	540	102	140
vol_2	84	70	96	143	20	13

Здесь:

n_1 — максимальное число слов во фразе по СЯУ в целом;

n_2 — во фразе из определяющих эталон;

$\text{vol}(n) = n!$ есть традиционная оценка для фразы длиной максимум в n слов;

vol_1 и vol_2 — оценки с применением предложенного метода выделения эталона СЯУ.

При этом:

$\text{vol}_1 = n_1 \cdot l_1$ есть оценка сверху, l_1 — число СЭ-фраз, задающих СЯУ;

$\text{vol}_2 = n_2 \cdot l_2$ есть оценка снизу, l_2 — число СЭ-фраз, определяющих эталон СЯУ.

- 1 *Ключевая особенность* изложенной методики формирования единиц представления экспертных знаний для разработки открытых тестов — *выделение лексоко-синтаксических связей слов* во фразе в рамках СЯУ *без привлечения внешних синтаксических анализаторов*.
- 2 Предложенный метод выделения смысловых эталонов даёт *минимум четырёхкратное сокращение объёма текстовых данных*, необходимых для *передачи* единицы *знаний* посредством ЕЯ без потери полезной составляющей между экспертами и обучаемыми *в открытых тестах*.
- 3 Показанная концепция СЯУ *позволяет* решать задачи *поиска систем зависимостей совместной встречаемости осмысленных фрагментов слов* в контексте связного текста на основе булевых векторов вида (7).
- 4 Отождествление компонент d_1, \dots, d_k вектора (7) с направленными лексико-синтаксическими связями *позволяет* формировать *правила принятия решений о наиболее вероятных направлениях* таких связей в заданных языковых контекстах.
- 5 Для *численной оценки* возможности *совместного появления* связей в контекстах, определяемых векторами указанного вида, может быть использована *оценочная функция*, аналогичная функции (5).

- 1 *Согласование данных* об основах и флексиях, выделяемых *по разным СЯУ* относительно фиксированной *предметной области*.
При этом объём баз знаний, формируемых на основе предложенного метода, может быть *дополнительно сокращён* в среднем *на 1,5%*.
- 2 *Статистика признаков* словоформ и определение *зависимых* слов *в составе связей* относительно ситуаций употребления *предметно-ограниченного* подмножества естественного языка:
 - интерпретация меры TF-IDF для оценки важности слова в контексте СЯУ;
 - в роли коллекции документов — совокупность СЯУ по заданной предметной области;
 - основополагающая гипотеза — зависимое слово имеет больший вес.
- 3 *Реконструкция* целостного *образа СЯУ* на основе вероятностей *совместной встречаемости* лексико-синтаксических *связей слов* в текстах заданного *предметно-ограниченного* ЕЯ-подмножества.