

Смысловые эталоны и передача знаний в задаче их оценки на основе открытых тестов

Михайлов Д. В., Емельянов Г. М.

Новгородский государственный университет
имени Ярослава Мудрого

11-я Международная конференция
«Распознавание образов и анализ изображений:
новые информационные технологии» (РОАИ-11-2013),

23–28 сентября 2013 г.

г. Самара

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 структуры знаний о синонимии, а также методов и алгоритмов их формирования и использования для совокупности задач:

- оценки схожести смыслов текстов предметно-ограниченного ЕЯ;
- автоматизации пополнения и компрессии баз языковых и предметных знаний;
- поиска наиболее рационального плана передачи заданного смысла между носителями заданного ЕЯ;
- согласования знаний, формируемых разными экспертами.

Определение 1.4

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 $n \in 1, \dots, |O|$.

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 подмножества, обозначаемые далее посредством соответствующего нижнего индекса и содержащие:

- указания на основу синтаксически главного слова (индекс 1);
- указания на флексию главного слова (индекс 2);
- связи «основа–флексия» для главного слова (индекс 3);
- сочетания флексий зависимого и главного слова (индекс 4). При этом после флексии главного слова через двоеточие указывается предлог (если такой имеется) для связи главного слова с зависимым;
- указания на флексию зависимого слова (индекс 5).

Задача: построить $I \subseteq G \times M$ анализом буквенного состава и отбором фраз $Ts_i \in Ts$ минимальной длины с наибольшим числом слов, наиболее употребимых в различных фразах из Ts (с учётом синонимов).

Пусть

W_{ij} — последовательность символов слова w_{ij} ,

Wc_{ij} — неизменная часть (основа),

Wf_{ij} — изменяемая (флексивная) часть,

\odot — обозначение для операций типа конкатенации.

Дано:

$Ts = \left\{ Ts_i : Ts_i = \odot_j W_{ij} \right\}$ — множество СЭ-фраз, задающих СЯУ.

Найти:

$Pw_i = \left\{ (Wc_{ij}, Wf_{ij}) : Wc_{ij} \odot Wf_{ij} = W_{ij} \right\}$ для всех $i = 1, \dots, |Ts|$.

Будем отождествлять с основой Wc_{ij} и флексией Wf_{ij} принятые в информатике понятия «префикс» и «суффикс».

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;
- $pref.s$ 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 . Суммарная длина общих префикса и суффикса при этом должна составлять минимум треть длины слова, а разность длин у пары слов с общим префиксом (независимо от суффикса) всегда меньше половины длины меньшего слова.

описание алгоритма

программная реализация

далее

Вход: Ts ;

Выход: $Pw = \bigcup_{i=1}^{|Ts|} Pw_i$;

1: $Pw := \emptyset$;

2: **для всех** W_{ij} : $\odot_j W_{ij} = Ts_i$, где $Ts_i \in Ts$

3: $Wc_{ij} := \{W_{ij}[1]\}$; $Wf_{ij} := \bigodot_{k=2}^{|W_{ij}|} W_{ij}[k]$;

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

5: $prfs(PrfsTmp)$;

6: **если** $PrfsTmp = \emptyset$ **то**

7: выдать Pw и выйти из алгоритма;

8: **иначе**

9: взять очередной Prf из $PrfsTmp$;

10: **если** $pref.check(Prf) = true$ **то**

11: $Pw := Pw \cup \left\{ (Prf, Wf_{ij}(Prf)) \mid pref.show(w_{ij}) = Prf \right\}$;

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

13: перейти к шагу 6;

14: **иначе**

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

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

17: **конец для**

18: перейти к шагу 5

19: **конец если**

20: **конец если**

Пусть T_s — множество СЭ-фраз, задающих некоторую СЯУ согласно (1),
 J — множество индексов неизменных частей слов фраз в составе T_s .

Определение 3.2

Последовательность индексов неизменных частей слов некоторой $T_{s_i} \in T_s$ назовём моделью её линейной структуры (МЛС), $L_s(T_{s_i})$.

Пусть $\{J_1, J_2\}$ — пара последовательностей индексов в $L_s(T_{s_i})$, где $J_1 = \{j_1^1, \dots, j_1^1\}$, $J_2 = \{j_1^2, \dots, j_1^2\}$, а парам (j_1^1, j_1^1) и (j_1^2, j_1^2) соответствуют синтаксические связи.

Для формирования эталона отбираются $T_{s_i} \in T_s$, в МЛС которых

$$(J_1 \subset J_2) \vee (J_2 \subset J_1) \vee (|J_1 \cap J_2| = 1) \vee (J_1 \cap J_2) = \emptyset, \quad (3)$$

а суммарная длина всех последовательностей указанного вида для всех связей, выявленных на T_{s_i} , **минимальна**.

Пусть $fr(w_j)$ — частота появления слова w_j во всех $Ts_i \in Ts$.

Тогда **наиболее информативные** в Ts слова составят кластер *Clust*:

- слово с максимальным значением данной частоты войдёт в *Clust*;
- для $\forall \{w_j, w_k\} \subset Clust$ и $\forall w_l \notin Clust$ верно то, что

$$\left(|fr(w_j) - fr(w_k)| < |fr(w_j) - fr(w_l)| \right) \wedge \\ \wedge \left(|fr(w_j) - fr(w_k)| < |fr(w_k) - fr(w_l)| \right) = \text{true} \quad (4)$$

Основу эталона составляют фразы с **максимумом** слов, вошедших в *Clust*.

При этом для слов из *Clust* учитываются различные порядки следования их во фразе и возможные синонимы.

Пусть LS — множество моделей линейных структур ЕЯ-фраз из Ts на J .

Лемма 5.1

Пара индексов $\{j_1, j_2\} \subset J$ соответствует словам-синонимам и могут быть заменены одним индексом из $(\mathbb{N} \setminus J)$, если $\exists \{Ls(Ts_1), Ls(Ts_2)\} \subseteq Ls$:

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

где $J_1 \subset J$, $J_2 \subset J$, а \odot есть операция типа конкатенации над J .

Пусть J_{Cl} — множество индексов слов кластера **наиболее информативных** относительно СЯУ, задаваемой множеством СЭ-фраз T_s ;

$freq((j, k), LS)$ — частота появления пары (j, k) в моделях из множества LS с учётом того, что $(j, k) \Leftrightarrow (k, j)$.

Тогда **эталон СЯУ** определяют фразы, МЛС которых входят в множество

$$LC = \bigcup_i LS_i : LS_i \subset LS, \exists \{Ts_i, Ts_j\} \in T_s : \\
 \begin{aligned}
 &LS(Ts_i) \in LS_i \\
 &|LS(Ts_i) \cap J_{Cl}| \rightarrow \max \\
 &\left((LS(Ts_j) \in LS_i) \wedge (Ts_j \neq Ts_i) \right) \rightarrow (LS(Ts_i) \cap J_{Cl}) \subset LS(Ts_j),
 \end{aligned}$$

а построение **признакового множества ФК (2) эталона СЯУ** требует:

- найти пары индексов $(j, k) : freq((j, k), LS) > 1$, отвечающие условию (3), для всех МЛС множества LC ;
- каждой найденной (j, k) задать направление синтаксической связи;
- из $\forall LS_i \subset LC$ исключить МЛС с индексами, не входящими ни в одну из найденных связей.

Поиск $Dir(j, k)$, $Dir \in \{\leftarrow, \rightarrow\}$, идёт в три этапа:

- проверка на признак **ложной связи**;
- попытка отождествления с **ранее выделенными связями**;
- при отсутствии ассоциации с известными связями — опрос эксперта.

Пусть $St(j)$, $St(k)$ и $St(l)$ — основы слов, отвечающие индексам j , k и l .

В заданной СЯУ для пары (j, k) связь ложная, если $j, k, l \in Ls(Ts_i)$ в некоторой $Ts_i \in Ts$ и имеется СЯУ, где связь $St(j)$ и $St(k)$ ложная, но есть связь либо между $St(j)$ и $St(l)$, либо между $St(k)$ и $St(l)$.

Пусть Lnk — множество **ранее найденных связей**, каждая представлена:

- идентификационным номером СЯУ (Id);
- основой главного слова (St_1);
- основой зависимого слова (St_2);
- списком пар «флексия главного слова–флексия зависимого» (FCm).

Паре (j, k) соответствует связь $((j, k), \rightarrow)$, если для некоторой СЯУ

$\exists (Id, St_1, St_2, FCm) \in Lnk$:

$$St(j) = St_1, St(k) = St_2, \text{ а } (Fl(j), Fl(k)) \in FCm.$$

Пример: исходное множество семантически эквивалентных фраз

Синонимичные перифразы

27:89

Insert

Indent

Modified

"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."

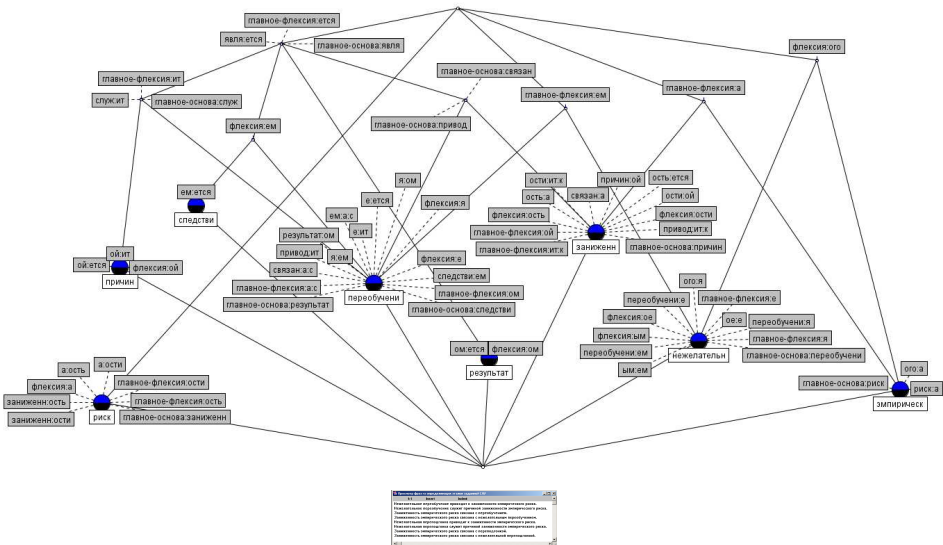
"Заниженность эмпирического риска относится к следствию нежелательного переобучения."

"Заниженность эмпирического риска связана с нежелательным переобучением."

"Нежелательное переобучение является причиной заниженности эмпирического риска."

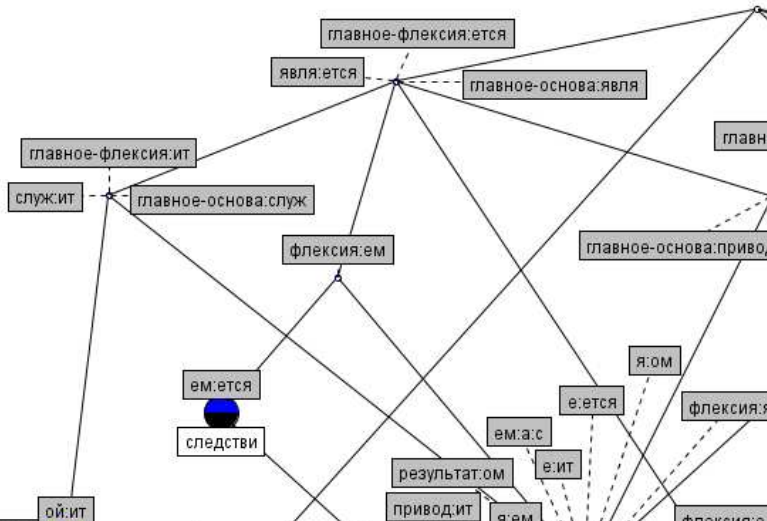
"Заниженность эмпирического риска, причиной которой служит нежелательное переобучение."

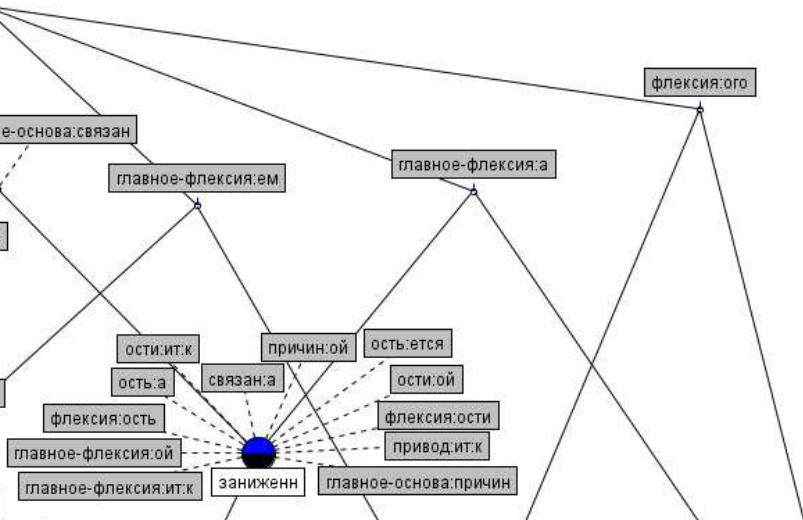
Результат: ЕЯ-фразы смыслового эталона и его формальный контекст

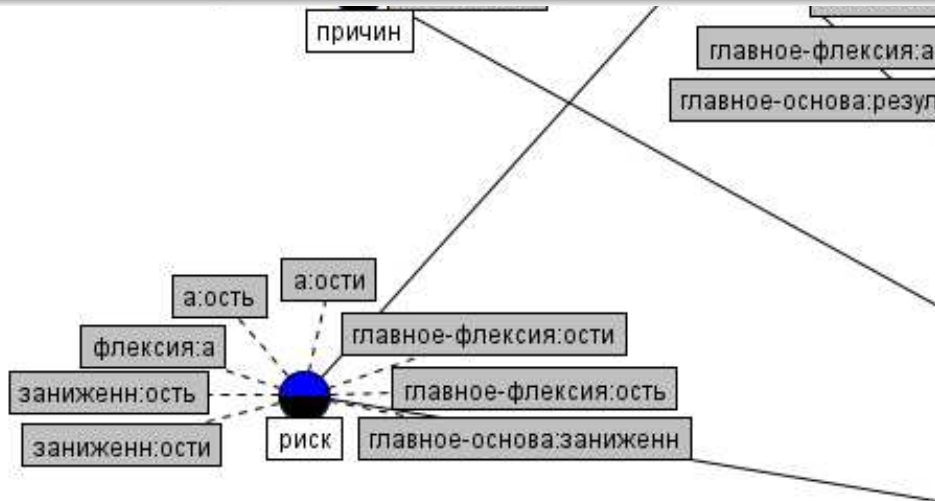


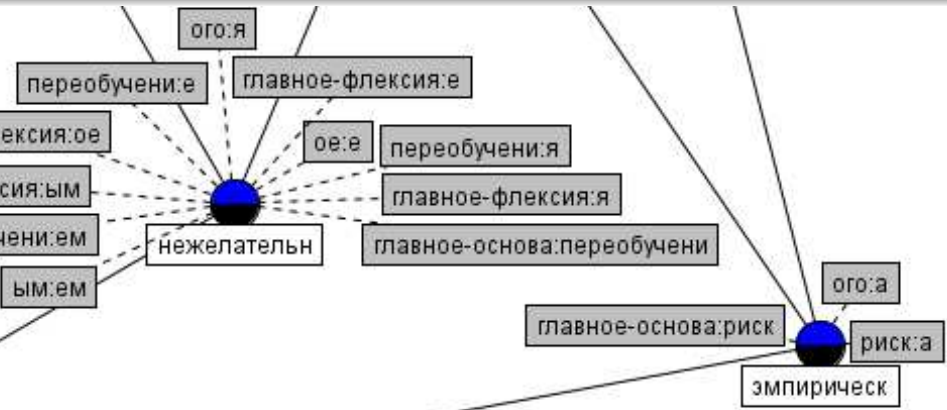
назад

далее

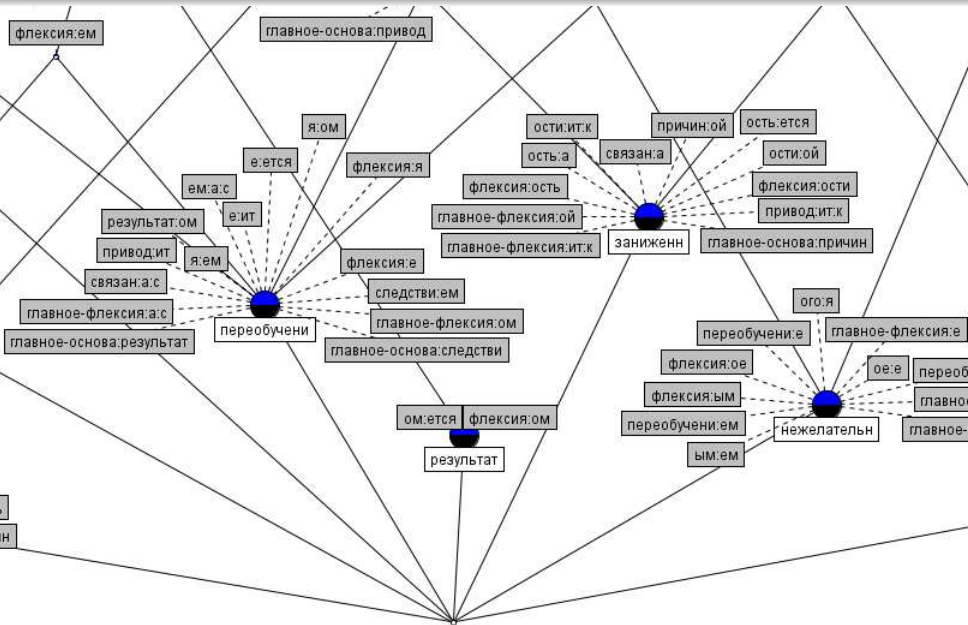








Результат: ЕЯ-фразы смыслового эталона и его формальный контекст



Просмотр фраз из определяющих эталон заданной СЯУ

1:1 Insert Indent

Нежелательное переобучение приводит к заниженности эмпирического риска.
Нежелательное переобучение служит причиной заниженности эмпирического риска.
Заниженность эмпирического риска связана с переобучением.
Заниженность эмпирического риска связана с нежелательным переобучением.
Нежелательная переподгонка приводит к заниженности эмпирического риска.
Нежелательная переподгонка служит причиной заниженности эмпирического риска.
Заниженность эмпирического риска связана с переподгонкой.
Заниженность эмпирического риска связана с нежелательной переподгонкой.

Порядковый номер СЯУ, i	1	2	3	4	5	6
Число фраз, задающих СЯУ	56	28	29	30	6	10
из них представляют эталон	8	9	7	9	1	2
Исходное число объектов СЯУ	18	17	15	13	12	14
Исходное число признаков СЯУ	177	186	173	162	94	81
Число объектов эталона	9	12	12	11	8	12
Число признаков эталона	82	90	80	69	35	53

i Ситуация языкового употребления

- 1 Связь переобучения с эмпирическим риском
- 2 Связь переусложнения модели с заниженностью средней ошибки на тренировочной выборке
- 3 Влияние переподгонки на частоту ошибок дерева принятия решений
- 4 Причина заниженности оценки обобщающей способности алгоритма
- 5 Зависимость оценки ошибки распознавания от выбора решающего правила
- 6 Зависимость обобщающей способности логического алгоритма классификации от числа закономерностей алгоритмической композиции

пример согласование знаний о синонимии оценка требуемого объёма памяти

i	1	2	3	4	5	6
n	12	15	16	17	10	14
$vol(n)$	$4.790 \cdot 10^8$	$1.308 \cdot 10^{12}$	$2.092 \cdot 10^{13}$	$3.557 \cdot 10^{14}$	$3.629 \cdot 10^6$	$8.718 \cdot 10^{10}$
$vol_1(n)$	648	795	416	442	20	42
$vol_2(n)$	168	225	80	187	20	42

Здесь:

i — порядковый номер СЯУ;

n — максимальное число слов во фразе;

$vol(n) = n!$ есть традиционно используемая оценка;

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

При этом:

$vol_1(n) = l_1 \cdot n$ есть оценка сверху, l_1 — число фраз, определяющих СЯУ;

$vol_2(n) = l_2 \cdot n$ есть оценка снизу, l_2 — число фраз, определяющих эталон СЯУ.

Рассмотрим представление тезауруса формальным контекстом

$$Kth = (Gth, Mth, Ith), \quad (5)$$

где Gth состоит из **символьных пометок** отдельных СЯУ;

Mth содержит **признаки** ФК вида (2) всех $gth \in Gth$.

Кроме того, в составе Mth выделяются **подмножества**:

- M_6 — указаний на **объекты** формальных контекстов вида (2) отдельных $gth \in Gth$;
- M_7 — множество связей «**основа–флексия**» для синтаксически зависимого слова;
- M_8 — множество **сочетаний основ** зависимого и главного слова.

По аналогии с ФК (2) отдельной СЯУ имеем $Ith \subseteq Gth \times Mth$.

пример представления СЯУ в формальном контексте тезауруса

При этом численная оценка схожести СЯУ определяется числом признаков, которые разделяются объектами сравниваемых ситуаций относительно формального контекста тезауруса.

согласование знаний о синонимии относительно разных СЯУ

Пример представления СЯУ в формальном контексте тезауруса



[назад к определению ФК тезауруса](#)

[далее](#)

Согласование знаний о синонимии относительно разных ситуаций языкового употребления

Пусть

St — основа слова (его неизменная часть);

Fl — его флексия;

S_1 и S_2 — некоторые СЯУ.

Предположим, что некоторое слово Wrd относительно S_1 представляется как $St_1 \odot Fl_1$, а относительно S_2 — как $St_2 \odot Fl_2$, причём $St_1 = St_2 \odot Sf$, где Sf содержит минимум один символ, а \odot есть операция конкатенации символьных строк.

Тогда относительно S_1 основа St_1 будет заменена на St_2 , флексия Fl_1 заменяется на $Fl_3 = Sf \odot Fl_2$, но только в том случае, если частоты встречаемости флексий Fl_3 и Fl_2 во всех лексико-синтаксических связях, представляемых формальным контекстом вида (5) для заданной предметной области, не уменьшаются при выполнении указанных замен.

Пример.

СЯУ №3, $St_1 =$ «является», $Fl_1 =$ «»,

СЯУ №1, $St_2 =$ «явля», $Fl_2 =$ «ется», $Sf =$ «ется»,

а относительно СЯУ №3 производится замена: Fl_1 — на $Fl_3 =$ «ется».

назад

далее

Практическое приложение: система тестирования знаний

Тестирование знаний и подготовка к ЕГЭ

База знаний Тесты Первое знакомство Window Помощь

Численные оценки близости правильному ответу

Испытуемые	Иванов Е.А.	Петров М.Н.	Сидоров Д.Л.	Зайцев Е.А.	Волков А.В.
Вопрос 1	0.857	1.000	0.4	1.000	0.857
Вопрос 2	1.000	0.733	0.868	0.75	0.545
Вопрос 3	0.75	0.63	0.000	0.703	0.42
Вопрос 4	0.861	0.861	0.717	0.662	1.000
Вопрос 5	0.725	0.657	0.000	0.5	0.471

Messages

Демо-версия системы представлена
на www.machinelearning.ru, на персональной странице автора.

[подробнее](#)

[далее](#)

Результат по испытуемому	
Испытуемый:	Петров М.Н.
Вопрос теста (вопрос №3):	
Как влияет переподгонка на частоту ошибок дерева принятия решений ?	
Полученный ответ:	
Именно с переобучение связана увеличение частоты ошибок дерева принятия решений на контрольной (= тестовой) выборке.	
Наиболее близкий вариант правильного ответа:	
Увеличение частоты ошибок дерева принятия решений на контрольной выборке связано с переподгонкой.	
Численная оценка близости правильному ответу:	0.63
Оценка за ответ:	удовл.

[назад](#)[далее](#)

Результаты группового тестирования после автоматического согласования знаний о синонимии относительно различных СЯУ

Тестирование знаний и подготовка к ЕГЭ

База знаний Тесты Первое знакомство Window Помощь

Численные оценки близости правильному ответу

Испытуемые	Иванов Е.А.	Петров М.Н.	Сидоров Д.Л.	Зайцев Е.А.	Волков А.В.
Вопрос 1	0.857	1.000	0.4	1.000	0.857
Вопрос 2	1.000	0.733	0.868	0.75	0.545
Вопрос 3	0.75	0.652	0.000	0.703	0.42
Вопрос 4	0.913	0.913	0.717	0.595	0.89
Вопрос 5	0.725	0.657	0.000	0.5	0.471

Messages

- Случай 1. *Неполный ответ* — все слова и словосочетания из ответа испытуемого нашли прообразы в наиболее близком варианте правильного ответа, но *часть слов правильного ответа не нашла прообразов в ответе испытуемого.*
Нулевое значение оценки схожести с объектом из ФК СЯУ правильного ответа будет для упущенного слова, которое в «правильном» варианте синтаксически зависимо относительно некоторого другого слова, присутствующего в анализируемом ответе.
- Случай 2. *Орфографические ошибки (из допустимых)* — слово из ответа испытуемого и слово правильного ответа есть *формы одного и того же слова в рамках некоторой лексико-синтаксической связи из известных системе.*
- Случай 3. *«Лишние» слова* — в анализируемом ответе имеются слова, не нашедшие прообразов в правильном «варианте».
Ответ засчитывается как неверный, если «лишнее» слова фигурируют в известных системе лексико-синтаксических связях.

- В предложенной концепции СЯУ все связи между главным и зависимым словом предполагались одинаково значимыми. Для их применения в задачах оценки знаний по отраслям *схожестть СЯУ* необходимо *переформулировать с позиций нечёткой логики*.
- Для описания функций принадлежности нечётким множествам необходим *системный анализ структуры профессиональных знаний* в конкретной области.
- *Базис импликаций* формального контекста ситуации языкового употребления может послужить основой *разработки стратегий и правил синтаксического анализа*.
- Концепцию модели линейной структуры предложения можно сделать более гибкой, введя *вероятности совместной встречаемости слов* относительно текстов заданной предметной области и жанра.