

УДК 004.912:519.7

**МОДЕЛЬ И МЕТОД МНОГОАСПЕКТНОГО ПОИСКА ФАКТОГРАФИЧЕСКИХ ДАННЫХ
ДЛЯ ПОДДЕРЖКИ ПРИНЯТИЯ РЕШЕНИЙ****В. А. Тертышный, И. В. Шевченко**Кременчугский национальный университет имени Михаила Остроградского
ул. Первомайская, 20, г. Кременчуг, 39600, Украина. E-mail: mainhousepost@hotmail.com

Построена формальная модель многоаспектного семантического пространства. Выявлены типы логических связей, которые могут быть использованы для многоаспектного фактографического поиска. Выбраны признаки и метрики для кластеризации контекстов. Сформирована специализированная логико-лингвистическая модель для выявления фактографических данных. Предложен метод многоаспектного анализа текстов, выявления релевантных контекстов и формирования фактографического отчета. При анализе запроса каждому аспекту присваивается весовой коэффициент. На этапе поиска релевантных документов используется нечеткий частотный анализ. На этапе фактографического поиска формируются матрицы связей термов в абзацах документов и матрицы частот пар термов по абзацам. Полученные матрицы сравниваются с эталонными матрицами по аспектам. Критерий сравнения – расстояние Кемени. На этапе определения степени принадлежности каждого абзаца к каждому аспекту используется пороговое преобразование. На последнем этапе осуществляется ранжирование аспектов и найденных абзацев по степени релевантности определенным аспектам.

Ключевые слова: фактографический поиск, логико-лингвистическая модель, многоаспектный анализ текстов.

**МОДЕЛЬ І МЕТОД БАГАТОАСПЕКТНОГО ПОШУКУ ФАКТОГРАФІЧНИХ ДАНИХ
ДЛЯ ПІДТРИМКИ ПРИЙНЯТТЯ РІШЕНЬ****В. А. Тертишний, І. В. Шевченко**Кременчуцький національний університет імені Михайла Остроградського
вул. Першотравнева, 20, м. Кременчук, 39600, Україна. E-mail: mainhousepost@hotmail.com

Побудовано формальну модель багатоаспектного семантичного простору. Виявлено типи логічних зв'язків, які можуть бути використані для багатоаспектного фактографічного пошуку. Обрано ознаки і метрики для кластеризації контекстів. Сформоване спеціалізовану логіко-лінгвістичну модель для виявлення фактографічних даних. Запропоновано метод багатоаспектного аналізу текстів, виявлення релевантних контекстів і формування фактографічного звіту. При аналізі запиту кожному аспекту присвоюється ваговий коефіцієнт. На етапі пошуку релевантних документів використовується нечіткий частотний аналіз. На етапі фактографічного пошуку формуються матриці зв'язків термів в абзацах документів і матриці частот пар термів по абзацах. Отримані матриці порівнюються з еталонними матрицями по аспектах. Критерій порівняння – відстань Кемени. На етапі визначення ступеня приналежності кожного абзацу до кожного аспекту використовується граничне перетворення. На останньому етапі здійснюється ранжування аспектів і знайдених абзаців за ступенем релевантності певним аспектам.

Ключові слова: фактографічний пошук, логіко-лінгвістична модель, багатоаспектний аналіз текстів.

АКТУАЛЬНОСТЬ РАБОТЫ. Компьютерная обработка информации, отображенной на естественном языке, является базовой проблемой в области построения интеллектуальных систем поддержки принятия решений на основе информационного поиска. Огромное количество информации, способной повлиять на принимаемые решения, представлено в виде текстов. Поэтому создание и развитие информационных технологий обработки текстовой информации является одной из важнейших задач на протяжении многих лет. Получено большое количество научных и практических результатов в осуществлении морфологического, синтаксического и частично семантического анализа [1], однако сложной задачей остается извлечение необходимой для принятия решений фактографической информации, особенно, в тех случаях, когда эта информация затрагивает несколько аспектов – различных сторон решения конкретной задачи. Поэтому современные средства автоматической обработки естественных языковых данных не всегда удовлетворяют потребностям пользователя. Кроме того, полное математическое обеспечение процессов син-

таксического и семантического анализа является сложным и его реализация в СППР, предназначенных для решения управленческих задач на предприятиях затруднена.

Таким образом, сегодня одной из актуальных задач является разработка методов и моделей фактографического многоаспектного анализа текстовых массивов с использованием статистических подходов и частотно-семантических методов [2].

Анализ литературных данных и постановка проблемы. Процесс обработки естественно-языковых текстов лежит в основе современных информационных технологий извлечения сведений из текстовых документов. При этом используются методы морфологического, синтаксического и частично семантического анализа. Тем не менее, извлечение нужных сведений, т.е. понимание содержания текста остается задачей, которая не решена в полном объеме [3–6].

Частотно-семантический анализ текстовых массивов является одним из перспективных направлений современных информационных технологий. Составляющими такого анализа является алгоритмы класси-

фикации и кластеризации текстовых документов. В этих алгоритмах используют векторную модель текстовых документов, основанную на представлении документов как векторов в некотором фазовом пространстве. Базис такого пространства часто образуют с помощью частотно-дистрибутивных характеристик лексем текстового словаря. Одна из проблем такого подхода обусловлена большой размерностью рассматриваемого векторного пространства. Считается, что такое пространство не позволяет выделить заданные семантические составляющие в интеллектуальном анализе текстов, несмотря на развитую теорию представления знаний. Проблеме представления знаний посвящено много литературных источников. Начиная с истоков проблемы формализации естественного языка [7, 8] и до настоящего этапа обработки текстовой информации [9–11], известны пять формальных моделей представления знаний: продукционные, фреймовые, логические, семантические сети и онтологии. Однако ни одна из этих моделей не может в полной мере отразить содержание текста на естественном языке. Методы интеллектуального анализа текстовой информации Data Mining помогают решить многие задачи, с которыми сталкивается аналитик, в частности, разработаны и усовершенствованы методы классификации и кластеризации описаны многими учеными [12, 13]. Именно поэтому в настоящее время в большинстве случаев используют для обработки текстов статистические методы без привлечения процедур синтаксического анализа. Но проблемой остается эффективное извлечение фактографических данных, – прецедентов и их параметров, в случае многоаспектного поиска.

Целью работы является упрощение процесса многоаспектного поиска фактографических данных путем модификации метода латентно-семантического анализа и соответствующей логико-лингвистической модели.

МАТЕРИАЛ И РЕЗУЛЬТАТЫ ИССЛЕДОВАНИЙ. В продолжение исследований [14, 15] и в соответствии с поставленной целью необходимо решить следующие задачи:

- построить формальную модель многоаспектного семантического пространства;
- выявить возможные типы логических связей, которые могут быть использованы для многоаспектного фактографического поиска;
- выбрать признаки и метрики для кластеризации контекстов;
- сформировать специализированную логико-лингвистическую модель (ЛЛМ) для выявления фактографических данных;
- сформировать последовательность этапов многоаспектного анализа, выявления релевантных контекстов и формирования фактографического отчета для пользователя.

Формирование семантического пространства. Имеется множество релевантных документов, специально отобранных для данной онтологии и данного аспекта. Эти документы, в частности, содержат клю-

чевые слова, обозначающие сущности ПрО (главный тезаурус) и лексемы, связывающие сущности с сущностями и сущности с признаками, действиями и параметрами (тезаурусы аспектов решаемых задач). Необходимо сформировать онтологию ПрО и соответствующее семантическое пространство.

Представим онтологию ПрО с учетом наличия аспектов в виде

$$O = \langle E(AT), ER, F, AS, AR \rangle, \quad (1)$$

где E – набор сущностей; AT – множество атрибутов сущностей; ER – множество отношений сущностей; $F: E \times ER$ – функции интерпретации отношений и сущностей; AS – множество аспектов, определяющих подмножества сущностей и связей; AR – пересечение аспектов.

Тогда многоаспектное семантическое пространство:

$$S = \langle O, M \rangle, \quad (2)$$

где M – набор метрик для вычисления степени близости сущностей и релевантности результатов поиска. Пространство SS условно разделено на пересекающиеся подпространства, каждое из которых соответствует одному аспекту. На рис. 1 представлена схема связей между ПрО, аспектами и семантическими пространствами.

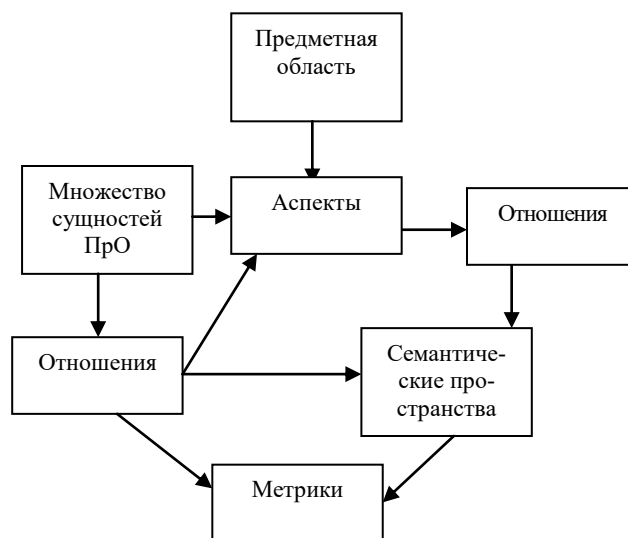


Рисунок 1 – Взаимосвязи между сущностями, аспектами и семантическими пространствами

Логические связи для многоаспектного фактографического поиска. Пусть T – входной текст, в котором имеются интерпретированные словоформы $X = \{x_1, \dots, x_m\}$ имеющие между собой скрытые отношения $R = \{r_1, r_2, \dots, r_n\}$. Проблема извлечения фактографических данных (знаний) состоит в поиске интерпретации $\varphi: X \rightarrow \Pi$, где Π – суперпозиция отношений R , выраженная на ЕЯ.

Особенность предлагаемой онтологии и семантического пространства – учет того факта, что между сущностями имеется несколько ассоциаций по нескольким аспектам.

Введем следующие допущения:

- а) семантическое пространство дискретно;
- б) набор элементов пространства конечен и обзорим;
- в) число комбинаций ключевых слов в контекстах велико, но конечно;
- г) каждый смысловой аспект можно представить кластером, центр которого определен на множестве релевантных документов.

Многоаспектность задачи предполагает, что одна и та же сущность находится в разных отношениях с другими сущностями. Если две сущности связаны несколькими отношениями в разных аспектах, тогда для каждого аспекта следует сформировать отдельную матрицу связей, построить свой тезаурус и своё пространство признаков.

В одной задаче могут быть важны одни аспекты, в другой – другие. Тогда между двумя узлами семантической сети есть несколько дуг, каждая из которых относится к одному аспекту. Каждой подсети соответствует своя матрица связей. Это дает возможность вычислять силу связей между сущностями с учётом аспектов, которые важны в данной задаче.

В частотно-семантическом методе в качестве критерия оценки «силы связи» обычно используются две характеристики: общность элементов и частота встречаемости.

Выбор признаков и метрики для кластеризации. Используем два пространства признаков. В первом пространстве признаков каждой оси координат соответствует относительная частота встречаемости ключевых слов, находящихся в тезаурусе ПрО. Во втором пространстве признаков каждой оси координат соответствует относительная частота встречаемости словосочетаний, находящихся в тезаурусах аспектов. Кластеры образуют документы, наиболее близкие по своим матрицам частот к эталонной матрице конкретного аспекта.

Учитывая, что в контексте решаемой задачи значимыми элементами языка являются документ, абзац (контексты) и слово, обозначим:

$N_A, f_A = N_A / N$ – количество и частота контекстов, где встретилось только слово А;

$N_B, f_B = N_B / N$ – количество и частота контекстов, где встретилось только слово В;

$N_{AB}, f_{AB} = N_{AB} / N$ – количество и частота контекстов, в которых наблюдалась совместная встречаемость слов А и В; N – общее количество контекстов.

В качестве основной метрики используем расстояние Кемени [16]. Как известно, бинарное отношение A на конечном множестве $Q = \{q_1, q_2, \dots, q_k\}$ – это подмножество декартова квадрата $Q^2 = \{(q_m, q_n), m, n = 1, 2, \dots, k\}$. При этом пара (q_m, q_n) входит в A тогда и только тогда, когда между q_m и q_n имеется рассматриваемое отношение.

Каждую кластеризованную ранжировку, как и любое бинарное отношение, можно задать квадратной матрицей $\|x(a, b)\|$ из 0 и 1 порядка $k \times k$. При этом $x(a, b) = 1$ тогда и только тогда, когда $a < b$ либо $a = b$. В первом случае $x(b, a) = 0$, а во втором $x(b, a) = 1$. При этом хотя бы одно из чисел $x(a, b)$ и $x(b, a)$ равно 1.

Расстоянием Кемени между бинарными отношениями A и B , описываемыми матрицами $\|a(i, j)\|$ и $\|b(i, j)\|$ соответственно, называется величина

$$KD(A, B) = \sum |a(i, j) - b(i, j)|, \quad (3)$$

где суммирование производится по всем i, j от 1 до k , т.е. расстояние Кемени между бинарными отношениями равно сумме модулей разностей элементов, стоящих на одних и тех же местах в соответствующих им матрицах.

При кластеризации с использованием расстояния Кемени центр кластера удобно определять как медиану Кемени – $Arg \min \sum D(A_i, A)$ [16].

Логико-лингвистическая модель для выявления фактографических данных. Формальной моделью представления знаний, учитывающей содержание предложений естественного языка, является логико-лингвистическая модель [2]. Поэтому, если построить формальную содержательную модель текста любой тематики и структуры, то можно будет анализировать электронные текстовые документы по содержанию, извлекать из них знания, сравнивать их.

Логико-лингвистическая модель текстового документа – это теоретико-множественная модель, которая включает основные свойства текста и его составных частей и отражает основные взаимосвязи между интересующими исследователя структурными компонентами текста.

С учетом перечисленных выше логических связей ЛЛМ текста для многоаспектного поиска выглядит следующим образом:

$$T = \langle CT, AS, DB, A, RA \rangle, \quad (4)$$

где CT – маркер принадлежности текста к определенному классу, маркер присваивается тексту на предварительном этапе поиска релевантных текстов; AS – множество аспектов в тексте; DB – база данных, представляющая собой набор тезаурусов, хранящих ключевые слова текста, лексемы и предложения; A – множество абзацев текста; RA – пересечение множества абзацев с множеством аспектов.

Таким образом, выражение (4) представляет усовершенствованную логико-лингвистическую модель, содержащую множество аспектов и проекцию аспектов на множество абзацев исследуемого текста, что позволяет реализовать многоаспектный фактографический поиск.

Метод многоаспектного анализа для фактографического поиска. Метод включает в себя несколько этапов, каждый из которых, в свою очередь, содержит определенные подэтапы.

На попередньому етапі відбувається формування онтології (1) і семантичного простору (2). Створюються таблиці зв'язків між аспектами. Кожна запис таблиці відповідає терму з тезауруса, а кожне поле – аспекту. Елементи таблиці – бінарні, сигналізуючі про наявність або відсутності даного терма в даному аспекті.

Формуються еталонні матриці частотних зв'язків аспекта на основі аналізу стопроцентно релевантних даним аспекту текстів. Ці тексти формують або підбирають експерти.

Етап 1. Обробка запиту.

1.1. З запиту вибираються ключові слова і лексеми.

1.2. Визначаються ведучі аспекти запиту шляхом аналізу таблиці зв'язків аспектів.

1.3. Уточнення у користувача аспектів і їх значущості (користувач може сам ранжувати аспекти і /або додавати нові). Повернення до п.1.2. Якщо користувач погодився, то перехід до п.1.4.

1.4. Групування ключових слів і лексем за аспектами. Генерація запитів за аспектами.

При аналізі запиту кожному аспекту присвоюється ваговий коефіцієнт

$$\alpha_j = \sum_{i=1}^n a_{ij}, \quad (5)$$

де a_{ij} – маркер участя i -го терма в j -м аспекті; n – кількість термів в тезаурусі.

Цей коефіцієнт буде використаний при формуванні результатів фактографічного пошуку.

Етап 2. Пошук і групування релевантних документів.

Спочатку з вихідного множини відбираються документи, для яких згідно векторного методу ступінь релевантності ПрО перевищує заданий поріг. Це класична процедура, використовувана в усіх пошукових системах.

Далі проводиться порівняння вектора ключових слів запиту q_k і вектора ключових слів документа x_k для первинного відбору документів в робочу підмножину аспектів A_j .

$$r_j = \sum_{k=1}^{L_q} x_k q_k, \quad (6)$$

де L_q – потужність множини термів в запиті з урахуванням розширення запиту за аспектами; k – номер терма.

Для кожного документа з робочої підмножини розраховується нормована частота спільної зустрічності ключових слів аспекта і документа:

$$r_j = \frac{\sum_{k=1}^t x_k^j q_k}{\sum_{k=1}^t (x_k^j)^2 + \sum_{k=1}^t (q_k)^2 - \sum_{k=1}^t x_k^j q_k}. \quad (7)$$

Отримані значення частот підлягають фазифікації з нечіткими оцінками «низька», «середня», «висока». Частоти представляються в нечіткій інтерпретації і ступені істинності вказаних термів порівнюються зі ступенями істинності цих термів в документах-образцях з бібліотеки СППР. Результат – сформовані підмножини релевантних документів за аспектами.

Етап 3. Формування робочих матриць частотних зв'язків подієльний документ. На даному етапі з вихідної еталонної матриць зв'язків аспекта видаляються рядки і стовпці цих термінів, які не виявлені в документі.

Етап 4. Обробка кожного знайденого документа для фактографічного пошуку.

4.1. Препозиції і абзаци в документі індексуються. Парі термів індексуються за порядком номерів.

4.2. Для i -тої кожної парі термів, зафіксованої в тезаурусах ПрО за аспектами, розраховуються відносні частоти fs_i зустрічності цих пар в препозиціях (sentences) документа

$$fs_i = \frac{n(t_1, t_2)}{n(t_1) + n(t_2)}. \quad (8)$$

Результат:

- квадратна матриця MLS_j зв'язків термів в препозиціях документа t за кожною аспектом з індексом j ;
- таблиця (матриця) індексів «препозиція–частота зустрічності парі термів».

4.3. Розраховуються відносні частоти fp_i зустрічності парі термів в абзацах (paragraphs) документа для кожної парі термів, зафіксованої раніше в тезаурусах ПрО за аспектами. Результат:

- квадратна матриця зв'язків MLP_j термів в абзацах документа;
- таблиця (матриця) індексів «абзац–частота зустрічності парі термів».

4.4. Отримані робочі матриці MLS_j і MLP_j порівнюються з еталонними матрицями за аспектами. Для кожної робочої матриць вираховується відстань Кемени (KD) за формулою (3) відносно еталонної матриць аспекта. З урахуванням того, розмір матриць може бути різним, для наступного відбору лексем і абзаци слід вираховувати нормовану відстань

$$F = n^2 / KD, \quad (9)$$

де n – фактичний розмір робочої матриць.

Використовуючи порогове перетворення, де T – заданий поріг,

$$H = \begin{cases} 1 & \text{if } F \geq T \\ 0 & \text{otherwise} \end{cases}, \quad (10)$$

визначаємо ступінь належності кожного абзаци за кожною аспектом. Абзац (прецедент) може належати до кількох аспектів, тому результат фактографічного пошуку може виявитися надлишковим. Але в даному випадку надлишковість переважить.

тельнее потери информации при рассмотрении реальной задачи в разных аспектах.

Этап 5. Ранжировка аспектов по весовым коэффициентам, ранжировка найденных абзацев (прецедентов), фраз и значений параметров по степени релевантности определенным аспектам и выдача пользователю подмножества прецедентов, фраз и параметров по каждому аспекту.

В качестве примера практического использования покажем фрагментарно решение задачи «Поиск новых покупателей (каналов сбыта) для предприятия сельскохозяйственного машиностроения». В табл. 1 показаны ключевые слова и лексемы для аспектов данной задачи. Согласно предлагаемому методу после введения запроса и извлечения из него ключевых слов и лексем определяются ключевые аспекты, и система пытается определить задачу. Пользователю предлагается уточнить тип задачи. После уточнения типа зада-

чи происходит формирование вторичных запросов, поиск в глобальной сети и отбор релевантных документов. Для каждого отобранного документа формируются рабочие эталонные матрицы частотных связей (этап 3 метода). Далее происходит формирование рабочих матриц документов, и рассчитываются частоты совместной встречаемости терминов и лексем в документах, абзацах и предложениях. Фрагмент матрицы частотных связей документа и аспекта «Функционал» показан в табл. 2.

Матрицы «абзац–частота встречаемости пары термов» и «предложение–частота встречаемости пары термов» не могут быть показаны из-за ограниченного объема статьи. Выбранные из документов фрагменты текста группируются согласно вычисленному значению расстояния (9).

Таблица 1 – Ключевые слова и лексемы для аспектов задачи «Поиск новых покупателей (каналов сбыта)»

Поиск новых покупателей (каналов сбыта)			
Аспекты			
Покупка–продажа	Продукция	Клиенты–покупатели–поставщики	Функционал
Ассортимент	Ассортимент	Адрес	Высокосортная мука
Дилер	Жатка	Ассортимент	Оборудование
Дистрибьютор	Жатка зерновая	Валюта	Отбор отрубей
Заказ	Жатка навесная	Выставка	Крупа
Импорт	Жатка валковая	Госзаказ	Мука
Клиент	Жатка валковая прицепная	Дилер	Переработка
Куплю	Жатка зерновая навесная	Дистрибьютор	Просеивание муки
Объем продаж	Жатка валковая прицепная	Е-mail	Производство
Партнерство	Мельница	Форум	Помол зерновых
Покупка	Мельница “Фермер”	Объем	Отопление
Покупатель	Мельница агрегатная	Телефоны	Точного высева
Потребитель	Мельница малогабаритная		Оптовая торговля
Продажа	Навесное оборудование		Розничная торговля
Покупка	Сеялка		Сельхозпродукты
Сбыт	Тракторы		
Сервис	Котлы отопительные		
Скидки	Крупополития		
Тендер			
Цена			
Экспорт			

Таблица 2 – Матрица частотных связей документа и аспекта

Наименование документа и аспекта	Оборудование	Крупа	Мука	Переработка	Производство	Помол зерновых	Оптовая торговля
Оборудование	–	0,34	0,76	0,96	0,85	0,73	0,0
Крупа		–	0,54	0,43	0,32	0,62	0,03
Мука			–	0,21	0,57	0,78	0,11
Переработка				–	0,23	0,15	0,12
Производство					–	0,16	0,13
Помол зерновых						–	0,11
Оптовая торговля							–

Решение указанной задачи с использованием СППР, реализующей предлагаемый метод фактографического поиска, позволило отделу сбыта и маркетинга выявить потенциальных клиентов и увеличить прибыли предприятия.

ВЫВОДЫ. Для упрощения процесса многоаспектного поиска фактографических данных с применением концепции латентно-семантического анализа построена формальная модель многоаспектного семантического пространства, выявлены типы логических связей, которые могут быть использованы для многоаспектного фактографического поиска. Выбраны признаки и метрики для кластеризации контекстов.

Сформирована специализированная логико-лингвистическая модель для выявления фактографических данных. Предложен метод многоаспектного анализа текстов, выявления релевантных контекстов и формирования фактографического отчета для пользователя.

ЛИТЕРАТУРА

1. Башмаков А.И., Башмаков И.А. Интеллектуальные информационные технологии: учебное пособие. – М.: МГТУ им. Баумана, 2005. – 304 с.
2. Вавиленкова А.И., Ланде Д.В., Литвиненко О.С. Теоретичні основи аналізу електронних текстів: монографія. – К.: НАУ, 2015. – 258 с.
3. Foltz W., Gilliam S., Kendall S. Supporting content based feedback in online writing evaluation with LSA // *Interactive Learning Environments*. – 2000. – No. 8. – PP. 111–128.
4. Gries S.Th. Corpus-based methods and cognitive semantics: the many meanings of to run // *Corpora in cognitive linguistics: corpus-based approaches to syntax and lexis*, 2006. – PP. 57–99.
5. Evans V. Lexical concepts, cognitive models and meaning-construction // *Journal of Cognitive semiotics*. – 2006. – PP. 73–107.
6. Кобозева И.М. Лингвистическая семантика. –

М.: Эдитореал УРСС, 2000. – 352 с.

7. Поспелов Д.А. Логико-лингвистические модели в системах управления. – М.: Энергоиздат, 1981. – 232 с.

8. Осуга С., Сазки Ю. Приобретение знаний. – М.: Мир, 1990. – 304 с.

9. Джарратано Д. Экспертные системы: принципы разработки и программирование – 4-е изд. / Пер. с англ. – М.: ООО «Вильямс», 2007. – 1152 с.

10. Корпусна лінгвістика / В.А. Широков, О.В. Булгаков, Т.О. Грязнухіна та ін. – К.: Довіра, 2005. – 471 с.

11. Design of a conceptual level programming environment based on task ontology / S. Kazuhisa, I. Mitsuru, K. Osamu, M. Riichiro // *Proc. of Successes and failures of knowledge based systems in real world applications*, 1996. – PP. 11–22.

12. Вавиленкова А.И. Извлечение смысла из предложений естественного языка // Программные продукты и системы. – 2012. – № 4 (100). – С. 87–90.

13. Палагин А.В., Крытый С.Л., Петренко Н.Г. Знание-ориентированные информационные системы с обработкой естественно-языковых объектов: основы методологии и архитектурно-структурная организация // *Управляющие системы и машины*. – 2009. – № 3. – С. 42–55.

14. Артамонов В.В., Тертышный В.А. Разработка модели информационного поиска с использованием связанных данных // *Системы обработки інформації*. – 2015. – № 10. – С. 69–75.

15. Тертышный В.А. Модель специализированной системы поиска сущностей на основе связанных данных // *Вісник Кременчуцький національний університет імені Михайла Остроградського*. – 2014. – Вип. 5/2014 (88). – С. 112–117.

16. Тоценко В.Г. Методы и системы поддержки принятия решений. Алгоритмический аспект. – К.: Наукова думка, 2002. – 381 с.

MODEL AND METHOD FOR MULTI-ASPECT SEARCH OF FACTOGRAPHIC FACTS FOR SUPPORT OF DECISION-MAKING

V. Tertishniy, I. Shevchenko

Kremenchuk Mykhailo Ostrohradskyi National University

vul. Pervomayskaya, 20, Kremenchuk, 39600, Ukraine. E-mail: mainhousepost@hotmail.com

Purpose. Formal model of aspect semantic space was built. Types of logical connections, which can be used for multi-aspect search, were identified. Features and metrics for clustering contexts were chosen. Specialized logical-linguistic model was formed for indication of factographic data. **Results.** Method of multi-aspect analysis indication of relevant contexts and formation of factographic reports was offered. During the analysis of request every aspect has got the weight coefficient. At the stage of relevant documents search an unclear frequent analysis is used. At the stage of factographic search matrixes of connections of terms in the documents and matrixes of terms by paragraphs have been built. **Practical value.** The Kemeny method has been applied in order to compare the aspects of these matrixes with those of the etalon matrixes. At the stage of defining rate of every paragraph belonging to each aspect threshold conversion is used. At the last stage the ranging of aspects and paragraphs which have been identified is conducted; the rate of their relevance to definite aspects is taken into account. References 16, tabl. 2, figure 1.

Key words: factographic search, logical-linguistic model; multi-aspect texts analysis.

REFERENCES

1. Bashmakov, A.I., Bashmakov, I.A. (2005), *Intellektualnyie informatsionnyie tehnologii: uchebnoe posobie* [Intelligent information technology: a training manual], MGTU im. Bauman, Moscow, Russia.
2. Vavilenkova, A.I., Lande, D.V., Litvinenko, O.E. (2015), *Teoretichni osnovi analizu elektronnih tekstiv: monografiya* [Basics of E-Teoretichni analizu tekstiv: monograph], NAU, Kyev, Ukraine.
3. Foltz, W., Gilliam, S., Kendall, S. (2000). "Supporting content based feedback in online writing evaluation with LSA", *Interactive Learning Environments*, no. 8, pp. 111–128.
4. Gries, S.Th. (2006). "Corpus-based methods and cognitive semantics: the many meanings of to run", *Corpora in cognitive linguistics: corpus-based approaches to syntax and lexis*, pp. 57–99.
5. Evans, V. (2006). "Lexical concepts, cognitive models and meaning-construction", *Journal of Cognitive Semiotics*, pp. 73–107.
6. Kobozeva, I.M. (2000), *Lingvisticheskaya semantika* [Linguistic semantics], Editoreal URSS, Moscow, Russia.
7. Pospelov, D.A. (1981), *Logiko-lingvisticheskie modeli v sistemah upravleniya* [Logico-linguistic models in control systems], Energoizdat, Moscow, USSR.
8. Osuga, S., Saeki, Yu. (1990), *Priobretenie znaniy* [Acquisition of knowledge] Mir, Moscow, USSR.
9. Dzharratano, D. (2007), *Ekspertnyie sistemy: printsipy razrabotki i programmirovaniya* [Expert Systems: Principles of design and programming], LLC «Vilyams», Moscow, Russia.
10. Shirokov, V.A., Bulgakov, O.V., Gryaznukhina, T.O. et al. *Korpusna lingvistika* [Corpus linguistics], Dovira, Kiev, Ukraine.
11. Kazuhisa, S., Mitsuru, I., Osamu, K., Riichiro, M. (1996). "Design of a conceptual level programming environment based on task ontology", *Proc. of Successes and failures of knowledge based systems in real world applications*, pp. 11–22.
12. Vavilenkova, A.I. (2012), "Removing the meaning of natural language sentences", *Software products and systems*, no. 4 (100), pp. 87–90.
13. Palagin, A.V., Kryvyyiy, S.L., Petrenko, N.G. (2009), "Knowledge-based information systems with the processing of natural language include the basics of methodology and architectural and structural organization", *Control systems and machines*, no. 3, pp. 42–55.
14. Artamonov, V.V., Tertyishnyiy, V.A. (2015), "Developing a model of information retrieval using linked data", *System on-timid Informácie*, no. 10, pp. 69–75.
15. Tertyishnyiy, V.A. (2014), "Model entities specialized search engine on the basis of related data", *Transaction of Kremenchuk Mykhailo Ostrohradskiy National University*, iss. 5, no. 88, pp. 112–117.
16. Totsenko, V.G. (2002), *Metody i sistemyi podderzhki prinyatiya resheniy. Algoritmicheskiy aspekt* [Methods and decision support systems. Algorithmic aspect], Naukova dumka, Kiev, Ukraine.

Стаття надійшла 24.10.2016.