

УДК 004.8

О.О. Арсірій, канд. техн. наук, доц.,  
А.А. Чугунов, канд. екон. наук, доц.,  
Ю.М. Ларченко, магістр,  
Одес. нац. політехн. ун-т

## ПОБУДОВА КОНТЕКСТНОЇ КАРТИ НА ОСНОВІ SOM ДЛЯ ВИДІЛЕННЯ КЛЮЧОВИХ СЛІВ ВЕБ-ДОКУМЕНТІВ ОСВІТНІХ ІНТЕРНЕТ-РЕСУРСІВ

*О.О. Арсірій, А.А. Чугунов, Ю.М. Ларченко. Побудова контекстної карти на основі SOM для виділення ключових слів веб-документів освітніх інтернет-ресурсів.* В результаті аналізу засобів просування освітніх інтернет-ресурсів запропоновано методику виділення ключових слів веб-документів на основі статистичних та морфологічних ознак шляхом побудови контекстної карти за допомогою SOM.

*Ключові слова:* веб-документ, ключові слова, СЯС, аналіз тексту, TF-IDF, самоорганізовані карти Кохонена, контекстні карти

*Е.А. Арсірій, А.А. Чугунов, Ю.Н. Ларченко. Построение контекстной карты на основе SOM для выделения ключевых слов веб-документов образовательных интернет-ресурсов.* В результате анализа средств продвижения образовательных интернет-ресурсов предложена методика выделения ключевых слов веб-документов на основе статистических и морфологических признаков путем построения контекстной карты с помощью SOM.

*Ключевые слова:* веб-документ, ключевые слова, СЯС, анализ текста, TF-IDF, самоорганизующиеся карты Кохонена, контекстные карты.

*Е.А. Arsiryi, A.A. Chugunov, Y.N. Larchenko. Creation contextual map based on SOM to searched keywords of web documents of educational internet-resources.* As a result of the analysis of tools of advancement internet resources offered the technique of searched keywords of web documents based on statistical and morphological characters by creating the contextual map based on SOM.

*Keywords:* web-document, the keywords, SCS, text analysis, TF-IDF, Kohonen self-organizing maps, contextual map.

**Актуальність проблеми. Постановка задачі.** Специфіка об'єкта управління, слабкий розвиток інформаційних систем та загострення проблеми профорієнтації в освітньої галузі роблять актуальним завдання розробки відвідуваних інтернет-ресурсів кафедр вищих навчальних закладів [1]. Розробка освітнього інтернет-ресурсу (ОІР) складається із етапів побудови, аналізу і просування в мережі Інтернет. Інформаційно-представницька спрямованість ОІР посилює важливість етапу просування, професійне виконання якого забезпечує оперативний доступ до довідкової, практичної і теоретичної інформації викладачів, студентів, абітурієнтів та інших зацікавлених осіб, підвищує кількість і якість відвідувачів сайту. Більшість потенційних відвідувачів потрапляють на ОІР при переході за посиланнями у результатах пошукових систем (ПС), робота яких заснована на визначенні так званих «ключових слів» (КС) веб-документу. З переліку таких слів складається семантичне ядро сайту (СЯС). СЯС являє собою список КС і їх комбінацій, записаних в метатеггах *keywords* і розподілених в контенті сайту, а саме, у тезі *title*, в *alt*-атрибутах, в тексті внутрішніх і зовнішніх посилань, у виділеннях жирним і курсивним шрифтом, на початку контенту сайту, в назві файлів, в *URL* та ін. Від повноти і точності розробки СЯС залежить положення сайту в списку видач ПС [2]. Розробка СЯС є основою так званих «білих» методів пошукової оптимізації, які, в свою чергу, є єдино можливими в умовах обмеженості грошового бюджету. Наявність обмежуючих умов ускладнює етап просування ОІР на відміну від просування конкуруючих з ними комерційних інформаційно-представницьких інтернет-ресурсів.

© О.О. Арсірій, А.А. Чугунов, Ю.М. Ларченко, 2013

З іншого боку тривалість розробки СЯС, а саме виділення КС веб-документів збільшується за рахунок наявності динамічного контенту в ОІР, а саме ведення тематичних форумів, та новинних блогів, періодичною зміною назв й переліків спеціальностей, появою нових наукових праць й видань і таке інше. При цьому тривалість виділення КС веб-документів може значно затримувати необхідну періодичність оновлення СЯС, що призводить до зниження повноти і точності СЯС, і ресурс втрачає свої позиції на сторінці видачі пошукових результатів. Тому метою дослідження є побудова методики автоматизованого виділення КС у веб-документі для скорочення часу розробки СЯС ОІР з динамічним контентом без втрати повноти і точності. Для розробки методики необхідно: проаналізувати існуючі методи аналізу тексту для виділення ознак КС веб-документів; визначити спосіб представлення виділених ознак наборів слів для автоматизації їх обробки; проаналізувати можливості SOM як засобу інтелектуальної візуалізації матриць даних при виділенні КС; запропонувати методику реалізації виділення КС за допомогою SOM на основі даних аналізу тексту.

**Методи аналізу тексту для виділення КС.** Всю сукупність представлених на сьогоднішній день методів аналізу тексту, відносно задачі виділення КС, можна умовно розділити на дві групи: лінгвістичний та статистичний аналіз.

Використання методів лінгвістичного аналізу дає можливість отримати сенс тексту за його семантичною структурою, що дозволяє точніше аналізувати текст, виділяючи його структурні особливості, але є більш трудомістким і складним у використанні, крім того передбачає орієнтацію на конкретну мову з її конкретними семантичними особливостями, що обумовлює погану міжмовну переносимість. Складність використання методів лінгвістичного аналізу для автоматизації виділення КС веб-документів пов'язана з багатством семантики та морфології природних мов. При цьому формальний опис правил природної мови та їх реалізація потребує залучення фахівців з області лінгвістики. Реалізація методів лінгвістичного аналізу у виділенні КС можлива за рахунок врахування морфологічних особливостей слова як елемента тексту (довжина слів, частини мови, рід, власна чи загальна назва, відміна, форма порівняння тощо).

Використання методів статистичного аналізу дає можливість виділити КС по частотному розподілу слів у тексті не залежно від мови, але не дозволяє повною мірою врахувати сенс тексту. При цьому автоматизація виділення КС веб-документів на основі використання статистичних методів більш проста та зручна у реалізації.

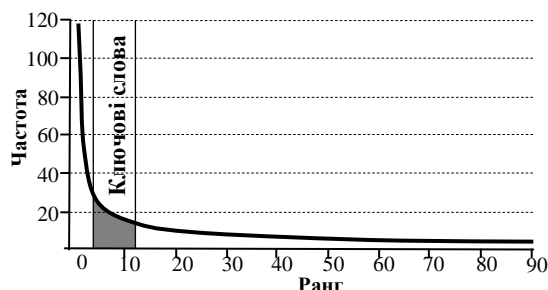
Методи статистичного аналізу для виділення КС базуються на емпіричних законах Ципфа [3]. Відповідно до першого закону Ципфа можна отримати графік залежності рангу  $r$  від частоти  $n$  (рис. 1, а). Кількість входження слова  $t$  у документ  $d$  буде його частотою  $n_t$ . Якщо розташувати частоти слів по мірі їх убавання і пронумерувати, тоді порядковий номер частоти буде рангом слова  $t$  —  $r_t$ . Дослідження показують, що найбільш значимі для документу слова лежать у середній частині графіка. Від вибору значень порогу зверху та знизу для обчислення рангу значимих слів залежить точність, якість і кількість виділення КС. Якщо встановити високе значення порогу зверху та низьке знизу, то в КС потраплять допоміжні слова, якщо навпаки — то можлива втрата смислових термінів. Тому використання тільки такого підходу для виділення КС викликає практичні складнощі.

Використання для визначення важливості слова веб-документу, який є частиною колекції веб-документів ОІР статистичної міри  $TF-IDF$  (від англ.  $TF$  — *term frequency*,  $IDF$  — *inverse document frequency*) — є продовженням робіт Ципфа і ліквідує складнощі їх практичного використання. При розрахунку  $TF-IDF$  вага слова  $t$  пропорційна кількості вживання цього слова в документі  $d$ , і обернено пропорційна частоті вживання слова в інших документах колекції  $D$ .

Таким чином, для визначення слів з високою вагою, які зустрічаються у досліджуваному документі необхідно визначити  $TF$  як відношення числа разів появи слова  $t$  до довжини документа  $d$

$$TF(t, d) = \frac{n_t}{N}, \quad (1)$$

де  $n_t$  — частота слова  $t$  у документі  $d$ ;  
 $N$  — кількість слів у документі  $d$ .



а

Слова	Показов в місяць
<a href="#">экономическая кибернетика</a>	1117
<a href="#">кафедра экономической кибернетики</a>	103
<a href="#">методы экономической кибернетики</a>	77
<a href="#">экономическая кибернетика специальность</a>	69
<a href="#">экономическая кибернетика скачать</a>	54
<a href="#">экономическая кибернетика учебник</a>	30
<a href="#">экономическая кибернетика работа</a>	25
<a href="#">факультет экономической кибернетики</a>	21
<a href="#">экономическая кибернетика лекции</a>	21
<a href="#">экономическая кибернетика реферат</a>	20
<a href="#">основные понятия экономической кибернетики</a>	14
<a href="#">экономическая кибернетика вакансии</a>	12

б

Рис. 1. Можливості виділення КС: за першим законом Ципфа (а), згідно сервісу Wordstat.yandex (б)

Для пониження ваги поширених слів, які зустрічаються майже у всіх документах колекції, вводять інверсну частоту IDF як відношення числа всіх документів колекції  $D$  до числа документів що містять слово  $t$ . Необхідно зазначити, що IDF буде близькою до нуля для поширених слів

$$IDF(t, D) = \frac{D}{d}, \quad (2)$$

де  $D$  — загальна кількість документів колекції;  
 $d$  — кількість документів, в яких зустрічається слово  $t$ .  
Статистична міра TF-IDF є добутком двох співмножників:

$$TFIDF(t, d, D) = TF(t, d) \times IDF(d, D). \quad (3)$$

Сучасні засоби статистики ПС (*wordstat.yandex*, *adstat.rambler*, *adwords.google*) за допомогою власних алгоритмів також дозволяють отримати список КС асоційованих з веб-документом на основі пошукових запитів користувачів мережі Інтернет (рис. 1, б).

**Представлення статистичних і морфологічних ознак наборів слів для автоматизації їх обробки.** Для автоматизації обробки веб-документу за статистичним та морфологічним аналізом набір слів можна представити у вигляді числових бінарних характеристичних векторів  $\{\vec{X}_1, \vec{X}_2, \dots, \vec{X}_n\}$ , де  $\vec{X}_t$  — характеристичний вектор слова  $t$ , а  $n$  — кількість слів в обраному для представлення фрагменті тексту. В свою чергу бінарний характеристичний вектор слова  $\vec{X}_t$  складається з статистичних і морфологічних ознак слова  $\vec{X}_t = \{x_1, x_2, \dots, x_k, y_1, y_2, \dots, y_m\}$ ,  $x = \{0, 1\}$ ,  $y = \{0, 1\}$ , де  $x, y$  — оцінки відповідно статистичної і морфологічної характеристики слова;  $k, m$  — кількість статистичних і морфологічних характеристик відповідно.

Значення оцінки характеристики може дорівнювати одиниці, якщо така властивість притаманна слову, або нулю, якщо слово такої властивості не має. Кількість характеристик визначається експертом. Набір характеристичних векторів можливо представити в вигляді

матриці даних в випадку якщо для всіх слів відомі значення всіх ознак, які їх характеризують. Кожен зі стовпців матриці відповідає одному слову із набору (табл. 1).

**SOM як засіб інтелектуальної візуалізації матриць даних при виділенні КС.** Для реалізації процедур інтелектуальної візуалізації матриць даних пропонується використовувати нейромережвий підхід, заснований на застосуванні самоорганізованих карт Кохонена (self-organizing map – SOM) [4]. Виділяють два методи візуалізації відображення вхідного простору в простір ознак за допомогою побудови самоорганізованих карт:

— у вигляді гнучкої решітки взаємозв'язаних нейронів, де вектори вагових коефіцієнтів є показниками, які спрямовані від відповідних нейронів карти до об'єктів із вхідного простору;

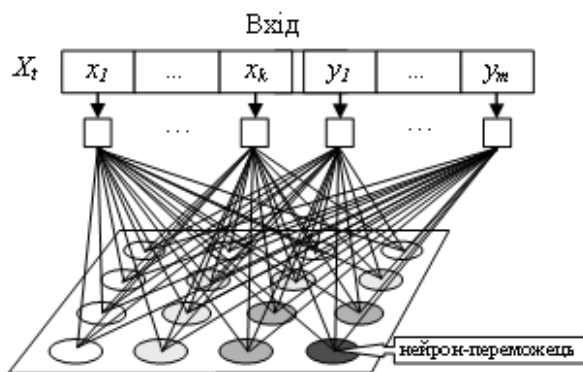
— у вигляді когерентних областей, які утворюють однаково помічені нейрони. При цьому мітки класів призначаються нейронам карти в залежності від їх збудження представниками вхідного простору, а під когерентністю розуміється те, що кожна із груп нейронів є відособленою множиною безперервних міток.

SOM які реалізовані за другим методом візуалізації називають контекстними або семантичними. Процес побудови контекстної SOM складається із процедур самоорганізації і градування. При виконанні процедури самоорганізації вектори  $X$  обираються із табл. 1 випадково і діють як набір значень на вході SOM (рис. 2). Вибірка і наступна за нею адаптація SOM продовжується ітераційно до поки отриманий асимптотичний стан SOM не буде визнаний стаціонарним. При виконанні процедури *градування* нейронам SOM присвоювалися мітки згідно найкращої відповідності векторів вагових коефіцієнтів нейронів вхідним векторам  $X$ .

Таблиця 1

Матриця вхідних даних

Ознака	Слова					
	$X_1$	$X_2$	...	$X_l$	...	$X_n$
$x_1$	1	0	...	0	...	1
$x_2$	0	0	...	1	...	0
...	...	...	...	...	...	...
$x_k$	0	1	...	0	...	0
$y_1$	0	0	...	1	...	1
$y_2$	1	0	...	1	...	0
...	...	...	...	...	...	...
$y_m$	1	0	...	0	...	1

Рис. 2. Відображення вхідного вектора  $X_l$  на SOM

**Методика автоматизованого виділення КС веб-документу.** Для виділення КС у веб-документах пропонується використати комплексний аналіз, який базується на поєднанні можливостей статистичного і лінгвістичного аналізу і використовує контекстну SOM як засіб інтелектуальної візуалізації. Запропонована методика складається з наступних етапів.

Етап 1. Розрахунок міри  $TF \times IDF$  та отримання статистики ПС.

Етап 2. Отримання морфологічних характеристик елементів тексту.

Етап 3. Представлення статистичних і морфологічних характеристик елементів фрагменту веб-документу у вигляді матриці вхідних даних (див. табл. 1).

Етап 4. Побудова контекстної SOM на основі матриці вхідних даних для отримання КС.

Покажемо реалізацію запропонованої методики на прикладі фрагменту тексту, що знаходиться на головній сторінці сайту кафедри ЕКІТ ОНПУ (<http://ecit.od.ua/index.php>).

Етап 1. Розрахунок міри  $TF \times IDF$  для елементів текстового фрагменту виконано за формулами (1)...(3). Для розрахунку використано наступні константи: довжина документу  $N$  дорівнює 677 слів; загальна кількість документів  $D$  у базі дорівнює 8 млрд. документів. При цьому за базу прийнято кількість веб-документів проіндексованих ПС Google у російському сегменті мережі Інтернет. За допомогою сервісів ПС (див. рис. 1, б) були отримані рекомендації



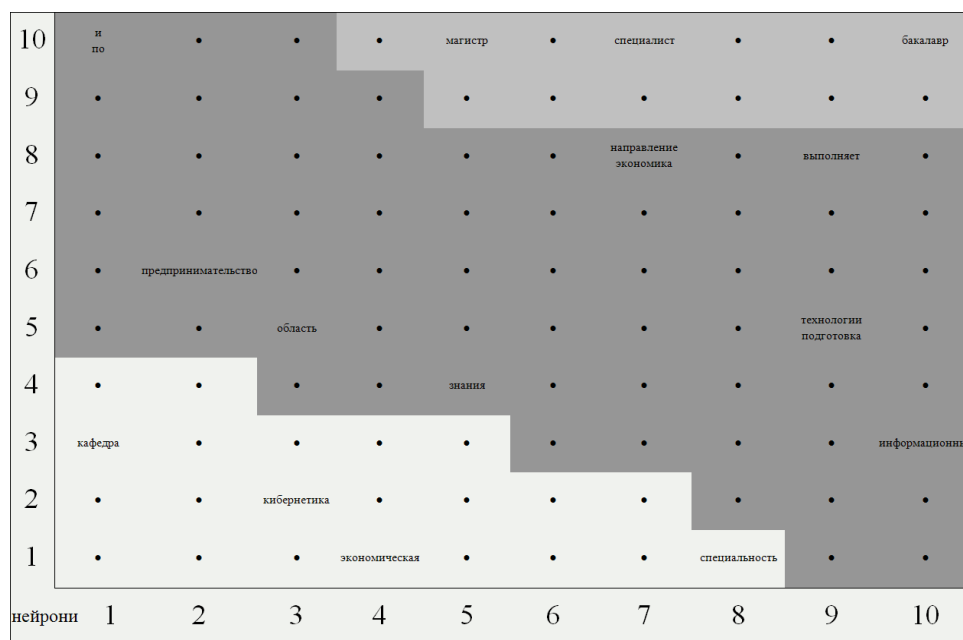


Рис. 3. Контекстна SOM із трьох когерентних областей

### Література

1. Солдатов, А.В. Информационная система как основа эффективного управления вузом / А.В. Солдатов // Университетское управление: практика и анализ. — 2004. — № 2(31). — С. 116 — 119.
2. Арсирій, Е.А. Автоматизация разработки и обновления семантического ядра сайта с динамическим контентом / Е.А. Арсирій, С.Г. Антошук, О.А. Игнатенко, Б.Ф. Трофимов // Искусственный интеллект. — 2012. — № 4. — С. 464 — 473.
3. Закон Ципфа — Вводная статья [Электронный ресурс] / Режим доступа: <http://webpavilion.ru/статья/закон-ципфа-вводная>. — 21.12.2012.
4. Кохонен, Т. Самоорганизующиеся карты / Т. Кохонен: Пер. 3-го англ. изд. — М.: БИНОМ, Лаборатория знаний, 2008. — 655 с.

### References

1. Soldatov, A.V. Informacionnaja sistema kak osnova jeffektivnogo upravlenija vuzom [Information system as the basis for effective university management]/ A.V.Soldatov // Universitetskoe upravlenie: praktika i analiz [University Management: Practice and Analysis]. — 2004. — # 2(31). — pp. 116 — 119.
2. Arsiriy, E.A. Automation of development and updating for semantic kernel of a site with dynamic content / E.A. Arsiriy, S.G. Antotshuk, O.A. Ignatenko, B.F. Trofimov // Artificial intelligence. — 2012. — # 4. — pp. 464 — 473.
3. Zakon Cipfa — Vvodnaja stat'ja [Zipf's Law — An introductory article] [Electronic resource] / Available at: <http://webpavilion.ru/stat'i/zakon-cipfa-vvodnaja>. — 21.12.2012.
4. Kohonen, T. Samoorganizujushiesja karty [Self-organizing map] / T. Kohonen: Per. 3-go angl.izd [Transl. 3-rd edition from Engl.] . — Moscow, 2008. — 655 p.

Рецензент д-р техн. наук, проф. Одес. нац. політехн. ун-ту Крисілов В.А.

Надійшла до редакції 21 грудня 2012 р.