

Київська міська адміністрація
Інститут прикладної Інформатики

На правах рукопису

ЖОГОВ Георгій Володимирович

АДАПТИВНИЙ ІНФОРМАЦІЙНИЙ ПОШУК. ПІДХІД І РЕАЛІЗАЦІЯ

05.13.17 - теоретичні основи Інформатики

Автореферат дисертації на здобуття наукового ступеня
кандидата фізико-математичних наук

Науковий керівник - старший науковий співробітник
кандидат фізико-математичних наук
ДРІЯНСЬКИЙ В.М.

Київ - 1994

АВ 31.003

Дисертація є рукопис

Робота виконана в Інституті прикладної інформатики
Науковий керівник кандидат фізико-математичних наук, старший нау-
ковий співробітник Дріанський Володимир Михайлович

Офіційні опоненти:

1. Доктор фізико-математичних наук, член-кореспондент АН України,
Ющенко Катерина Логвінівна
2. Кандидат фізико-математичних наук, старший науковий співробіт-
ник, Полумієнко Сергій Костянтинович

Провідна організація Інститут програмних систем
НАН України (м. Київ)

Захист відбудеться "25" травня 1994 р. об 10 год. на
засіданні Спеціалізованої вченої ради Д 166.01.01 в ін-
ституті прикладної інформатики (ІПрІн) за адресою:
252004, м. Київ, вул. Червоноармійська, 23-б.

З дисертацією можна ознайомитися у бібліотеці Інституту приклад-
ної інформатики.

Автореферат розісланий "23" вересня 1994 р.

/ Вчений секретар
Спеціалізованої вченої ради

Г.Б.

Мелент'єв Г.Б.

ЛНБ України ім.В.Стефаніка



ім. В. Стефаніка
АН України

00777063 (U)

ЗАГАЛЬНА ХАРАКТЕРИСТИКА РОБОТИ

Актуальність теми. В інформаційних системах, орієнтованих на ефективну обробку великих об'ємів текстографічної інформації: бібліотечних, науково-технічних, юридичних, офісних и т.п., особливе значення набувають можливості адаптації пошукових механізмів до інформаційних інтересів конкретних абонентів, тобто наявність засобів, які забезпечують, по-перше, початкове інформаційне позиціонування користувача і створення його особистого інформаційного середовища і, по-друге, підтримку цього середовища в актуальному стані.

У відомих пошукових схемах для рішення цієї проблеми в значній мірі використовуються інтелектуальні зусилля абонентів пошукової системи. Осв чому створення нових пошукових моделей з елементами штучного інтелекту, які значно зменшують інтелектуальне навантаження абонента, є актуальним.

Ціль роботи - розробка схеми адаптивного інформаційного пошуку і визначення в ній процесів, які дозволяють використання адаптаційних механізмів; моделювання цих процесів, розробка алгоритмів функціонування побудованих моделей і створення прототипу відповідної інформаційної системи.

Наукова новина. Запропонована концепція адаптивного інформаційного пошуку і схема інформаційної технології, яка базується на цій концепції. В рамках цієї схеми побудовані моделі локалізації інформаційної потреби абонентів і вивчені її властивості, запропоновані семантичні розширення цих моделей, розроблені алгоритми їх функціонування. Створено прототип системи адаптивного інформаційного пошуку і здійснено цикл експериментів, які підтвердили правильність запропонованої концепції.

Загальна методика досліджень. Для моделювання процесів локалізації інформаційної потреби виділяються два рівня опису: міжсеансовий та внутрішньосеансовий. Моделі міжсеансового рівня побудовані з використанням алгебраїчного підходу до теорії розпізнавання образів (по Ю.І.Журавльову). Для внутрішньосеансового моделювання використовується байєсовський підхід у статистичній теорії прийняття рішень, а для побудови семантичних розширень - теорія відношень.

Практична цінність. Реалізація концепції адаптивного пошуку

приводить до появи принципово нового класу інтелектуальних інформаційних систем з "прозорим" інтерфейсом користувача. Використання подібних систем дозволяє по-новому подивитись на цілий ряд галузей використання інформаційних технологій: перш за все це бібліографічні інформаційні системи, електронні каталоги і інформаційні мережі, оскільки в межах запропонованого підходу може бути здійснено автоматичне виявлення інформаційних потреб користувачів і подальше автоматичне вибіркове розповсюдження інформації в їх персональні бази даних (БД).

Реалізація результатів досліджень. Реалізовано і перевірено на експериментальній БД прототип системи адаптивного інформаційного пошуку. Система була запроваджена в Експо-центрі "Наука" АН України. Елементи цієї системи були використані при розробці експертної системи скринінгу хворих вторинним імунodefіцитом і системи обробки кореспонденції ДЖНТ України. За цикл робіт "Настраивающиеся документальные системы; математические модели, алгоритмы, программное обеспечение" автор спільно з Л.Г.Катериничем та О.Я.Колтуном був удостоєний у 1986 р. медалі АН України з премією для молодих вчених.

Апробація роботи. Матеріали дисертації доповідались на семінарах Наукової ради АН України з проблеми "Кібернетика"; "Проблеми проектування автоматизованих банків даних"; "Лингвистические проблемы проектирования информационных систем"; "Распознавание и оптимальное управление развитием систем" (Славская, 1989 р.); на Республіканській науково-технічній конференції "Пути дальнейшего совершенствования системы научно-технической информации и пропаганды в Украинской ССР в XII пятилетке" (Хмельницький, 1986 р.); на Всесоюзному науково-технічному семінарі "Мобильное программное обеспечение" (Калінін, 1988 р.); на 5 Всесоюзній конференції по БД та знань (Львів, 1991 р.).

Публікації. За результатами виконаних досліджень опубліковано 11 робіт, одна колективна монографія (глави 6, 10).

Обсяг і структура роботи. Дисертація складається з вступу, трьох глав, заключення, списку літератури (63 найменування) і додатків. Дисертація викладена на 130 сторінках, містить 4 малюнки та 2 таблиці. Додатки на 31 сторінках.

ЗМІСТ РОБОТИ

Дисертаційна робота присвячена питанням розробки і реалізації схеми адаптивного інформаційного пошуку:

- моделюванню процесів локалізації і ведення інформаційної потреби користувача, а також власне інформаційного пошуку і оцінки знайденої інформації;

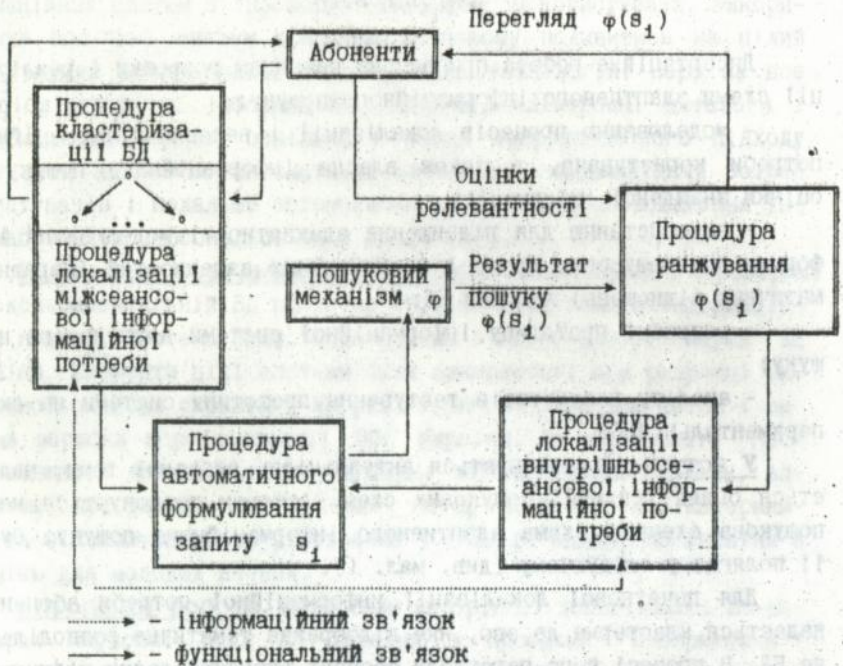
- використанню для підвищення адекватності моделювання інформаційно-пошукових процесів семантичних залежностей (парадигматичних відношень) лексики БД;

- розробці прототипу інформаційної системи адаптивного пошуку;

- аналізу результатів тестування прототипу системи на експериментальній БД.

У вступі обґрунтовується актуальність вибраної теми, надається огляд існуючих пошукових схем. Автором пропонується нова пошукова схема - схема адаптивного інформаційного пошуку. Суть її полягає у наступному (див. мал. 1).

Для початкової локалізації інформаційної потреби абоненту надається кластерне дерево, яке відображає тематичне розподілення БД. В процесі його перегляду абонент визначає деяку підмножину кластерів, які його зацікавили. На базі цієї інформації автоматично формується запит s_0 , по якому організується ітераційний пошук. Тобто, s_0 подається на вхід пошукового механізму ϕ . Результати пошуку $\phi(s_1)$ (в початковий момент $i=0$) проглядаються абонентом. Результати оцінок $\phi(s_1)$ подаються на вхід процедури внутрішньосеансової локалізації інформаційної потреби, результати роботи якої визначають подальшу роботу процедури автоматичного формування запиту s_{1+1} . Новий запит знову подається на вхід пошукового механізму. Ітераційний цикл зупиняється або системою, або користувачем. Результати спілкування з системою запам'ятовуються для подальшого використання у процедурі міжсеансової локалізації інформаційної потреби, яка дозволяє, по-перше, відслідкувати динаміку змінювання інформаційної потреби абонента від сеансу до сеансу i , по-друге, впливати на роботу процедури внутрішньосеансової локалізації.



Малюнок 1. Схема адаптивного інформаційного пошуку

У першій главі описуються розроблені автором математичні моделі локалізації інформаційної потреби користувача в схемі адаптивного інформаційного пошуку. Для вирішення задачі локалізації інформаційної потреби на базі схеми адаптивного пошуку пропонується виділити два види інформаційної потреби, — міжсеансову і внутрішньосеансову. Локалізація міжсеансової інформаційної потреби спрямована на визначення динаміки змінювання інформаційних інтересів користувачей і використовується для вирішення задачі класифікації абонентів. Локалізація внутрішньосеансової інформаційної потреби необхідна для ефективного навігації абонента по БД в межах одного сеансу спілкування його з системою.

Визначення 1. Міжсеансовою інформаційною потребою будемо називати підмножину кластерів класифікаційної структури БД, які були визнані абонентом пертинентними в процесі її перегляду.

Очевидно, що задачу локалізації міжсеансової інформаційної потреби можна розглядати як класичну задачу розпізнавання, а са-

ме, по опису множини документів $K = K_1 \cup K_2$, про яку відомо, що документи з K_1 визнані абонентом пертинентними, а документи з K_2 - непертинентними, і по опису центра кластера $d_s \in \mathcal{M}$ (\mathcal{M} - класифікаційна структура БД) вилічити предикати $d_s \in \mathcal{R}_1$ і $d_s \in \mathcal{R}_2$, де \mathcal{R}_1 і \mathcal{R}_2 відповідно класи пертинентних і непертинентних документів.

Нехай $\mathcal{M}_1 = H_K(\mathcal{M})$ - результат застосування розпізнаючого алгоритму (по типу вирахування оцінок) H до \mathcal{M} на базі навчальної вибірки K ; θ - дискретний час відносно деякої точки відліку, який пов'язаний з номером сеансу зв'язку абонента з системою (так, якщо T_0 -й сеанс прийняти за точку відліку, то в сеансі $T_0 + k$: $\theta = k$ ($k \geq 0$)); $\mathcal{M}_{\theta,z}^{T_0}$ - термінальні вершини структури \mathcal{M}_1 , знайденої в кінці сеансу θ абонента z ; $[0, T^*]$ - деякий цілочисловий інтервал (T^* - деяка постійна, яка характеризує час стаціонарності інформаційної потреби абонента z) і нехай кожній вершині $a_1 \in \mathcal{M}$ певним чином поставлено у відповідність цілочисловий параметр $T_{a_1} \in [0, T^*]$, причому $\forall a_1 \in \mathcal{M}_{\theta,z}^{T_0}$: $T_{a_1} = T^*$. Для користувача з наступним чином побудуємо множину $A_{\theta,z}^{T_0}$:

$$1) \text{ при } \theta = 1 \quad A_{\theta,z}^{T_0} = \begin{cases} \mathcal{M}, & T_0 = 0, \\ \emptyset, & T_0 > 0; \end{cases}$$

$$2) \text{ при } \theta = 2 \quad A_{\theta,z}^{T_0} = \mathcal{M}_{\theta-1,z}^{T_0};$$

$$3) \text{ при } \theta > 2 \quad A_{\theta,z}^{T_0} = \mathcal{M}_{\theta-1,z}^{T_0} \cup (X_\theta - Y_\theta - W_\theta),$$

де $X_\theta = \{a_1 \mid a_1 \in A_{\theta-1,z}^{T_0} \setminus \mathcal{M}_{\theta-1,z}^{T_0}\}$, причому

$$\forall a_1 \in X_\theta: T_{a_1} = T_{a_1} - 1;$$

$$Y_\theta = \{a_1 \mid a_1 \in X_\theta \ \& \ T_{a_1} = 0\};$$

$$W_\theta = \{a_1 \mid a_1 \in X_\theta \ \& \ T_{a_1} \neq 0 \ \& \ \exists a_j \in \mathcal{M}_{\theta-1,z}^{T_0}: W_{a_1}^{a_j}\},$$

$W_{a_1}^{a_j}$ - шлях з a_1 в a_j .

Визначення 2. Множину $A_{\theta,z}^{T_0}$ будемо називати множиною активних вершин.

Змістовно, $A_{\theta,z}^{\circ}$ формується на базі зостосування алгоритма Н до \mathcal{M} з урахуванням досвіду, накопиченого за сеанс $\theta-1$. При цьому в $A_{\theta,z}^{\circ}$ включаються всі вершини з $\mathcal{M}_{\theta-1,z}^{\circ}$ і ті вершини з $A_{\theta-1,z}^{\circ}$ (що не увійшли в $\mathcal{M}_{\theta-1,z}^{\circ}$), у яких значення параметрів $T_{a_1} > 0$ за умовою, що не існує шлях з a_1 в $\mathcal{M}_{\theta-1,z}^{\circ}$. Ті вершини, котрі за час T^* не розпізнаються Н як релевантні, вилучаються з множини активних вершин.

Визначимо на множині $A_{\theta,z}^{\circ}$ лінійний порядок ρ так, щоб в першу чергу виділити найбільш актуальні вершини (на базі значень T_a), а серед них найбільш пертинентні (з використанням Н), які становлять кластери нижніх рівней ієрархічної структури \mathcal{M} .

Визначення 3. Структуру $G_{T^{\circ},z}^{\circ} = \langle A_{T^{\circ},z}^{\circ}, \rho \rangle$ будемо називати глобальною інформаційною структурою абонента z в інтервалі часу $[T_0, T_0 + T^*]$, а T^* - часом стаціонарності.

Зафіксуємо деякий інтервал часу $[T_1, T_2] \subset [T_0, T_0 + T^*]$, і нехай T_1 - нова точка відліку.

Визначення 4. Структуру $L_{T_2,z}^1 = \langle A_{T_2-T_1,z}^1, \rho \rangle$ назвемо локальною інформаційною структурою абонента z в $[T_1, T_2]$, якщо

$$\text{card} \left\{ \theta = 1, \dots, T_2 - T_1 \mid \mathcal{M}_{\theta,z}^1 \right\} / \max_{\theta = 1, \dots, T_2 - T_1} \left\{ \text{card}(\mathcal{M}_{\theta,z}^1) \right\} > \Phi, \quad (1)$$
 де Φ - деяка порогова величина.

Визначення 4. Ядром структури $L_{T_2,z}^1$ назвемо множину

$$\text{Ker}(L_{T_2,z}^1) = \bigcap_{\theta = 1, \dots, T_2 - T_1} \mathcal{M}_{\theta,z}^1.$$

Зауваження. $\forall a_1 \in \text{Ker}(L_{T_2,z}^1) : T_{a_1} = T^*$.

Визначення 5.

1. Якщо $\forall \theta \in [1, T^*] : A_{\theta-1,z}^{\circ} \cap A_{\theta,z}^{\circ} = A_{\theta,z}^{\circ}$, то будемо вважати, що абонент z має постійну інформаційну потребу. Клас таких абонентів позначимо G_1 .

2. Якщо $\forall t \in [0, T^*] : A_{T^{\circ},z}^{\circ} = A_{T^{\circ}+t,z}^{\circ} \& A_{T^{\circ},z}^{\circ} \neq \mathcal{M}$, (2) то будемо вважати, що інформаційна потреба абонента обмежена деяким постійним колом інтересів. Клас таких абонентів позначимо

через C_2 .

3. Якщо не виконується (2), то будемо вважати, що коло інтересів абонента z не має фіксованих границь в заданому інтервалі стаціонарності. Клас таких абонентів позначимо C_3 .

Вивимо співвідношення між глобальною і локальною структурами для різних класів абонентів.

Теорема 2. Для $z \in C_1$: $L_{T_2, z}^{T_1} = G_{T, z}^{T_1}$ з точністю до ρ .

Теорема 3. Для $z \in C_2$: $G_{T, z}^{T_1 + T^*} = G_{T, z}^{T_1}$ з точністю до ρ .

Результат, отриманий в теоремі 3, дозволяє для абонента $z \in C_2$ не вести з момента $T_0 + T^*$ на протязі кванту T^* його глобальну інформаційну структуру, що значно зменшує кількість обчислень.

Нехай, як і раніше, $z \in C_2$; \mathfrak{A} - деякий поріг. Розглянемо послідовність $\Omega_{T, z}^{T_0} = (a_{1, z}^{T_0}, a_{1, z}^{T_0+1}, \dots, a_{1, z}^{T_0+T^*})$.

Визначення 6. Інформаційна потреба абонента z має марківський характер, якщо $\forall \theta \in (0, T^*]: \exists L_{T_0+1+\theta, z}^{T+\theta} \& \exists L_{T_0+\theta+2, z}^{T+\theta}$.

Очевидно, що для абонента z з марківським характером інформаційної потреби $\Omega_{T, z}^{T_0}$ має властивість:

$$\text{card} \left\{ \text{Ker} \left[L_{T_0+\theta+1, z}^{T+\theta} \right] \right\} > \mathfrak{A} \&$$

$$\max \left\{ \text{card} \left\{ a_{1, z}^{T_0+\theta} \right\}, \text{card} \left\{ a_{1, z}^{T_0+\theta+1} \right\} \right\}$$

$$\& \frac{\text{card} \left\{ \text{Ker} \left[L_{T_0+\theta+2, z}^{T+\theta} \right] \right\}}{\max_{\theta=1, 2, 3} \left\{ \text{card} \left\{ a_{1, z}^{T_0+\theta} \right\} \right\}} < \mathfrak{A}.$$

Побудуємо ймовірнісний простір $\langle \{t_1\}, 2^{\{t_1\}}, P \rangle$, де t_0 - елементарна подія, яка полягає у тому, що перетин двох сусідніх елементів $\Omega_{T, z}^{T_0}$ є пустим; t_1 ($i = 1, \dots, T^* - 1$) - елементарна подія, яка полягає у тому, що перетин тільки $i + 1$ безпосередньо слідуючих один за одним елементів $\Omega_{T, z}^{T_0}$ не є пустим і задоволь-

няє (1); p_1 - ймовірність події ξ_1 ; $P : \{\xi_1\} \rightarrow \{p_1\}$ - функція така, що $P(\xi_1) = p_1$.

Теорема 4. Якщо $z \in C_1$, то $P(\xi_{T^*+1}) = 1$. (3)

Визначення 7. Абонент $z \in C_1$ ("Е" - слабо належить), якщо $\text{card}\{\text{Ker}(L_{T^*}^{\varphi, z})\} / \max_{\theta=1, \dots, T^*} \{\text{card}\{\mathcal{U}_{\theta, z}^{\varphi}\}\} > \varphi$ ($\varphi \in (0, 1)$).

Теорема 5. Справедливо, що $z \in C_1 \implies z \in C_0$.

Теорема 6. Справедливо, що $P(\xi_{T^*+1}) = 1 \implies z \in C_1$.

Нехай ζ - ймовірна величина така, що $\zeta(\xi_1) = 1 + i$ з ймовірністю p_1 , а $M\zeta$ - її математичне сподівання. Тоді:

якщо $P(\xi_{T^*+1}) = 1$, то з урахуванням теореми 6 $z \in C_1$ і для z на протязі часу $[T_0 + T^*, T_0 + 2T^*]$ достатньо вести в кожному момент $\theta \in [1, T^*]$ тільки структуру $L_{T_0 + T^* + \theta}^{T^* + \theta, z}$;

якщо $\forall [T_1, T_2] \subseteq [T_0, T_0 + T^*]: L_{T_2}^{T_1, z} = G_{T_2}^{\varphi, z}$, то $z \in C_1$ і в цьому випадку в інтервалі часу $[T_0 + T^*, T_0 + 2T^*]$ для обслуговування достатньо зберігати тільки $L_{T_0 + T^*}^{\varphi, z}$;

якщо $z \in C_2$ (у відповідності з (2)), то з урахуванням теореми 3 в інтервалі $[T_0 + T^*, T_0 + 2T^*]$ необхідно зберігати $G_{T_0 + T^*}^{\varphi, z}$; крім цього для z з марківським характером інформаційної потреби в кожному момент $\theta \in [1, T^*]$ необхідно зберігати тільки $L_{T_0 + T^* + \theta}^{T_0 + T^* + \theta - 1, z}$, а для z з немарківським характером інформаційної потреби в інтервалі $[T_0 + T^*, T_0 + 2T^*]$ необхідно зберігати $L_{T_0 + T^* + \theta + M\zeta}^{T_0 + T^* + \theta, z}$, де $\theta = 1, 1 + M\zeta, \dots, 1 + kM\zeta$; $k = [(T^* - 1) / M\zeta]$.

Випадок, коли $z \in C_3$, можна інтерпретувати як ситуацію $z \in C_2$ з марківським характером інформаційної потреби при $T^* \rightarrow \infty$. Для таких z вести $G_{T_0 + T^*}^{\varphi, z}$ не має сенсу, а його положення в \mathcal{U} в кожному момент θ слід визначати на базі $L_{T_0 + T^* + \theta}^{T_0 + T^* + \theta - 1, z}$.

Таким чином, при допущенні існування H , вирішено задачу локалізації міжсеансової інформаційної потреби абонентів. (Для побудови H можуть бути застосовані алгоритмічні схеми, які базуються на алгебраїчному підході до теорії розпізнавання образів

по Ю.І.Журавльову.)

Далі описується модель внутрішньосеансової інформаційної потреби, яка базується на гіпотезі зв'язності (ван Рійсберген, С.Джоунс): якщо термін x_1 добре розділяє класи релевантних і нерелевантних документів, то будь-який термін, що часто зустрічається разом з x_1 , також володіє цією властивістю. Тому задача полягає в знаходженні множини термінів UD, які розділяють документи БД на пертинентні і непертинентні (для конкретного абонента). Крім цього, так як терміни з UD використовуються у процедурі автоматичного формулювання запиту, то бажано в UD включати як терміни, що підвищують точність пошуку, так і терміни, що зменшують пошуковий шум.

Визначення 8А. Внутрішньосеансовою інформаційною потребою користувача UD^+ , яка забезпечує підвищення точності пошуку, будемо називати множину термінів БД, ймовірність яких знаходиться в пертинентних документах вище деякого порогу M^+ .

Визначення 8В. Внутрішньосеансовою інформаційною потребою користувача UD^- , яка забезпечує мінімізацію пошукового шуму, будемо називати множину термінів БД, ймовірність яких знаходиться в непертинентних документах вище деякого порогу M^- .

Нехай $D_1 = \{d_1\}$ - множина документів, переглянутих абонентом на протязі першого циклу роботи процедури ранжування, і нехай $T_1 = \{t_1\}$ - множина термінів, які знаходяться у всіх цих документах. Розглянемо трійку $\langle t_j, a_j, b_j \rangle$, де $t_j \in T_1$, a_j - кількість документів, які відмічені абонентом як пертинентні і в яких знаходиться термін t_j , b_j - кількість документів в БД, в яких знаходиться термін t_j .

На множині $X_1 = \{\langle t_j, a_j, b_j \rangle\}$ визначимо функцію $I: X_1 \rightarrow \mathbb{R}$ (де \mathbb{R} - множина дійсних чисел), яку будемо називати інформативністю терміна. Як приклад може виступати

$$I(t_j) = \log \left[\frac{p_j \cdot (1 - q_j)}{q_j \cdot (1 - p_j)} \right],$$

де $p_j = \frac{a_j + 0.5}{N + 1}$, $q_j = \frac{(b_j - a_j) + 0.5}{M + 1}$, N - загальна кількість пертинентних документів, M - кількість документів БД за виключенням N .

Змістовно, інформативність відображає здатність терміну розділяти документи БД на пертинентні і непертинентні.

Видно, що задача локалізації внутрішньосеансової інформа-

ційної потреби може бути зведена до задачі виділення з множини термінів, які входять в уже переглянуті користувачем документи, таких термінів, у яких інформативність більше (для позитивних значень інформативності) або менше (для негативних значень інформативності) деякого порогу.

Розглянемо задачу визначення UD^+ . Нехай $I^+ = \{I_k^+\}$ - множина різних позитивних значень інформативності термінів T_1 . Побудуємо ймовірнісний простір $\langle \{\xi_1\}, 2^{\{\xi_1\}}, P \rangle$, де ξ_1 ($i=1, \dots, \text{card}(I^+)$) - елементарна подія, яка полягає у тому, що інформативність терміна прийме значення I_k^+ ; p_1 - ймовірність події ξ_1 ; $P: \{\xi_1\} \rightarrow \{p_1\}$ - функція така, що $P(\xi_1) = p_1$.

Нехай ζ - випадкова величина така, що $\zeta(\xi_1) = I_k^+$ з ймовірністю p_1 , а $M\zeta$ - її математичне сподівання.

Тоді з T_1 виділимо підмножину термінів \mathfrak{E}_1^+ таку, що для будь-якого $t_j \in \mathfrak{E}_1^+$ інформативність I_j терміна t_j більша або дорівнює $M\zeta$. З X_1 виділимо підмножину \mathfrak{X}_1^+ , яка відповідає \mathfrak{E}_1^+ . Побудована таким чином \mathfrak{X}_1^+ (якщо покласти $M^+ = M\zeta$) якраз і є внутрішньосеансовою інформаційною потребою користувача, яка забезпечує підвищення точності інформаційного пошуку перед другим циклом роботи процедури ранжування.

Для 1-ої ($i \geq 2$) пошукової транзакції T_1 будеться наступним чином: $T_1 = \mathfrak{X}_{1-1}^+ \cup T(D_1)$, де $T(D_1)$ - множина термінів документів D_1 , які були переглянуті користувачем за час 1-ої транзакції. Для T_1 будеться відповідна їй множина \mathfrak{X}_1^+ , яка в свою чергу є внутрішньосеансовою інформаційною потребою користувача, що забезпечує підвищення точності інформаційного пошуку перед наступною пошуковою транзакцією.

Ймовірність p_1 випадкової події ξ_1 будемо відновлювати таким чином, щоб більші значення інформативності мали більшу ймовірність. Тоді p_1 і $M\zeta$ можна обчислювати за формулами:

$$p_1 = I_k^+ / \sum_{I_k^+ \in I^+} I_k^+ ; M\zeta = \left[\sum_{I_k^+ \in I^+} (I_k^+)^2 \right] / \sum_{I_k^+ \in I^+} I_k^+.$$

Множина \mathfrak{X}_1^+ може бути використана для підвищення точності пошуку процедурою формування запиту на (i+1)-ому циклі роботи у схемі адаптивного інформаційного пошуку.

Для визначення UD^- слід урахувувати терміни з негативними

значеннями інформативності. У цьому випадку їх можна було б використовувати в запитах зі зв'язкою "не". Для здобуття таких термінів, аналогічно X_1^+ , будується множина X_1^- , яка містить терміни t_{1j} , у яких значення інформативності негативне і менше M_1^- , де ζ^- - випадкова величина, визначена на множині різних негативних значень інформативностей.

В залежності від того, до якого класу буде віднесено конкретного користувача, застосовуються наступні дії:

1. Для користувачів з постійною інформаційною потребою термінологічні словники не змінюються і використовуються для автоматичного забезпечення режиму вибіркового розповсюдження інформації.

2. Для користувачів, інформаційна потреба у яких змінюється в межах деякого постійного кола інтересів, процедура локалізації внутрішньосеансової інформаційної потреби постійно провадить "рафінування" термінологічного словника перед кожною новою пошуковою транзакцією.

3. Для користувачів, у яких інформаційна потреба постійно змінюється, термінологічні словники не зберігаються і такі абоненти повинні переглядати кластерне дерево (для початкового позиціонування в БД) перед черговим сеансом зв'язку з системою.

Друга глава присвячена опису семантичної моделі та її використанню у схемі адаптивного інформаційного пошуку. Ключові моменти пропонованої моделі полягають у наступному.

Нехай $X = \{x_1\}$ - множина термінів, які позначають поняття $P = \{p_1\}$ вихідної предметної галузі, а $\hat{\omega}$ - множина парадигматичних відношень на множині P (всі парадигматичні відношення є толерантними або транзитивно-антисиметричними).

Відношення $\hat{\omega} \in \hat{\omega}$ індукує відношення ω на множині X :

$$\forall p_1, p_j \in P: p_1 \hat{\omega} p_j \implies x_1 \omega x_j.$$

Транзитивне відношення $\omega \in \Omega$ ϵ - (ϵ^{-1}) -орієнтоване, якщо $\forall (x_1, x_j) \in \omega: M_1 \rho M_j (M_1 \rho^{-1} M_j)$, де M_1, M_j - множина об'єктів, які відповідають поняттям p_1 і p_j , ρ - деяке відношення уточнення властивостей, дроблення об'єктів і т.п. Додамо до Ω відношення зворотне до транзитивних, розглянемо транзитивне замкнення Ω^+ множини Ω і введемо на Ω^+ асоціативну операцію "*" добутку відношень: $x_1 \omega_1 * \omega_2 x_j \iff \exists x_s: x_1 \omega_1 x_s \& x_s \omega_2 x_j (\omega_1, \omega_2 \in \Omega^+)$. Замикаючи Ω^+ відносно "*", утворимо універсальну алгебру складних від-

ношень $U_{\Omega^c} = \langle \Omega^c, (*) \rangle$.

Представлення $\omega \in \Omega^c$ у вигляді $K(\omega) = \omega_{1_1} * \dots * \omega_{1_s}$ назовемо канонічним, якщо $\forall k = 1, \dots, s-1 : \omega_{1_k} * \omega_{1_{k+1}} \in \Omega^+$.

Теорема 1. Для будь-якого $\omega \in \Omega^c$ $K(\omega)$ єдине.

Добуток $\omega_1 * \omega_2$ назовемо добутком $\omega_{1_k} \omega_{2_1}$ -типу, якщо

$$K(\omega_1) = \omega_{1_1} * \dots * \omega_{1_k} \quad \text{і} \quad K(\omega_2) = \omega_{2_1} * \dots * \omega_{2_s}.$$

Визначення 1. Канонічне представлення називається направленим, якщо в ньому не зустрічаються добутки $\omega_{1_j} \omega_{1_j}^{-1}$ -типу.

Теорема 2. Якщо $\omega_1 * \omega_2$ добуток $\omega_{1_k} \omega_{2_1}$ -типу і $\omega_{1_k} \neq \omega_{2_1}$, то $K(\omega_1 * \omega_2) = K(\omega_1) * K(\omega_2)$.

Наслідок. В канонічному представленні не можуть зустрічатися добутки $\omega_{1_j} \omega_{1_j}^{-1}$ -типів.

Можливість представлення складних відношень у канонічному вигляді дозволяє виконувати кількісну оцінку складних відношень на базі простих.

Використовуючи властивості алгебри U_{Ω^c} , будується функція семантичного зв'язку $f: \Omega^c \rightarrow (0,1)$. При цьому враховуються як топологічні характеристики графа залежності (довжина шляху між термінами, ширина основи поділу, тип відношень відповідних компонентів зв'язності), так і деякий набір семантичних передумов (характерних для задач кластеризації БД, локалізації інформаційної потреби, інформаційного пошуку), які обумовлюють властивості функції f (симетричність; значення f для транзитивних відношень більше її значень для толерантних відношень; значення f для відношень ω , канонічні представлення яких не містить толерантних відношень, перевершують значення f для відношень, які містять толерантні відношення; значення f узгоджені з порядком τ , який задається на Ω^+).

Задамо на Ω лінійний порядок τ і перенумеруємо відношення у відповідності з цим порядком: $\omega_1 \tau \omega_j \iff 1 < j$. Тоді для складного відношення ω , з урахуванням вигляду графа залежності і введених семантичних передумов, функція семантичного зв'язку має вигляд:

$$f_{\omega_{j_1} \dots \omega_{j_m} \omega_{1_1} \dots \omega_{1_n}}(x, y) = \prod_{l=1}^m (A_{j_l})^{k_l} \prod_{k=1}^n \frac{1}{B_{1_k}} \left[1 - \sum_{k=1}^n \frac{B_{1_k}}{B_{1_k-1}} \right]^k \sum_{s=1}^p \gamma_{\omega_s}^{OT},$$

де B_{1_k} - параметр, який відповідає транзитивному відношенню ω_{1_k} ;

A_{j_1} - параметр, який відповідає толерантному відношенню ω_{j_1} ; p - довжина шляху i_k -ої транзитивної компоненти складного відношення; k_1 - кількість ланок толерантного відношення ω_{j_1} , які входять у шлях, що з'єднує терміни x і y ; n - кількість транзитивних відношень в $K(\omega)$; m - кількість толерантних відношень в $K(\omega)$; $\gamma_{\omega_s}^{OT}$ - відносний коефіцієнт уточнення, який є функцією від основи поділу; параметри B_{i_k} і A_{j_1} задовольняють умовам

$$B_{i_0} \equiv B_{i_1}; \quad 0 < A_{j_1} < A_{j_{1+1}} < 1; \quad B_{i_k} > B_{i_{k+1}} > 1.$$

Для транзитивних (простих і складних) відношень значення функції f_{ω} повинні задовольняти умові

$$L_1 < f_{\omega_{i_1} \dots \omega_{i_n}} < L_2, \quad (1)$$

де $L_1 = \begin{cases} \max f_{\omega_{i_1} \dots \omega_{i_{j-1}}}, & \text{якщо } i_n = N, \text{ де } j \in \{1, \dots, n\} \text{ таке;} \\ \max f_{\omega_{i_1} \dots \omega_{i_n} \omega_N}, & \text{якщо } \forall k = \overline{0, n-j} \quad i_{j+k} = i_j + k; \\ & \text{у іншому випадку;} \end{cases}$

$$L_2 = \begin{cases} \min f_{\omega_{i_1} \dots \omega_{i_{n-1}}}, & \text{якщо } i_n = N; \\ \min f_{\omega_{i_1} \dots \omega_{i_{n-1}} \omega_{i_n+1} \omega_{i_n+2} \dots \omega_N}, & \text{у іншому випадку;} \end{cases}$$

N - кількість транзитивних відношень; $i_0 = N+1$; $f_{\omega_{N+1}} \equiv 1$ і $f_{\omega_0} \equiv 0$.

Для транзитивно-толерантних відношень, які містять тільки по одному простому толерантному відношенню, значення функції f_{ω} повинно задовольняти умові ($k=1, 2$):

$$\left[(A_{j_m})^{k_m, L_1}, (A_{j_m})^{k_m, L_2} \right], \quad (2)$$

Теорема 3. Для того, щоб $0 < f_{\omega} < 1$, достатньо

$$B_i/B_j \leq 2 / \left(1 + \sqrt{1 + 4\gamma_{\omega_s}^{\max}} \right), \quad i > j; \quad 0 < A_k < 1;$$

де $\gamma_{\omega_s}^{\max} = \max(\gamma_{\omega_s}) \cdot \max(\max(\gamma_{i_k}))$, γ_{i_k} - довжина k -го шляху в i -ій компоненті зв'язності, γ_{ω_s} - ширина основи поділу, ω_s - клас основи поділу.

Теорема 4. Рекурентні співвідношення

$$B_{N-1} \geq \gamma_{\omega_s}^{\max} \cdot B_N/b; \quad B_1 \geq \gamma_{\omega_s}^{\max} \cdot B_{i+1} \dots B_N/b;$$

де $0 < b < 1$, є достатніми для виконання (1).

Наслідок 1. Значення параметрів $B_1, 1=\overline{1,N}$, можна обчислювати за формулою $B_1 = \left(\gamma_{\Omega}^{\max} / b \right)^{2^{N-1}}$.

Наслідок 2. Умови

$A_1 < A_j; A_1 > A_M^2; A_M \leq \left[1 - \gamma_{\Omega}^{\max} \cdot \frac{c^{N-1} - 1}{c^N (c-1)} \right] / \left[(\gamma_{\Omega}^{\max} + 1) \cdot c^{2^{N-1}} \right] = L, A_j < A_{j+1} \cdot L;$
де $1 < j, M$ - кількість толерантних відношень в $\Omega, c = \gamma_{\Omega}^{\max} / b, \epsilon$ достатніми для виконання (2).

Наслідок 3. Лінійний порядок, індукований по f_{ω} на Ω природним порядком на $(0,1)$, співпадає з τ .

Побудуємо скалярно-семантичну і семантичну функції кореляції документів d_1, d_j .

1. Скалярно-семантична функція кореляції документів f^* вибирається у вигляді $f^* \stackrel{\text{def}}{=} S \left(\sum_{t \in d_j} \varphi_t(d_j) \right)$, де функція φ_t визначається наступним чином (t - термін):

- 1) $\forall t: \varphi_t(d_1) \cdot \varphi_t(d_j) = 1 \Rightarrow \varphi_t(d_j) = \varphi_t(d_1);$
- 2) $\forall t: \varphi_t(d_1) \cdot \varphi_t(d_j) = 0 \Rightarrow \varphi_t(d_j) = \max_{r \in d_j; \omega \in \Omega^0} (f_{\omega}(t,r))$

і $\varphi_t(d_1) = 1$, якщо $t \in d_1$, і $\varphi_t(d_1) = 0$ - у протилежному випадку.

2. Скалярно-семантична функція f^* дозволяє оцінити лише найбільш суттєві семантичні зв'язки між термінами порівнюваних документів. Для врахування всіх семантичних зв'язків визначаються поняття семантичної сили терміну $\beta_{1,K}$ у деякій множині термінів $K = (t_1, \dots, t_s)$ і узагальнює його поняття семантичної сили множини термінів K_1 у деякій множині термінів $K_2 - \beta_{K_1, K_2}$:

$$1) \beta_{1,K} = \sum_{t_1 \in T_1 \cap K} \sqrt[n]{f_{\omega}(1,1_j)}, T_1 = \{t_{1_j} \mid \exists \omega \in \Omega^0: t_1 \omega t_{1_j}\},$$

причому якщо існує декілька $\tilde{\omega} \in \Omega^0: t_1 \tilde{\omega} t_{1_j}$, то для визначення f_{ω} вибирається відношення, яке дає найбільшу силу зв'язку, тобто $f_{\omega}(1,1_j) = \max_{\tilde{\omega} \in \Omega^0} (f_{\tilde{\omega}}(1,1_j));$

$$2) \beta_{K_1, K_2} = \sum_{t_1 \in K_1} \beta_{1, K_2}$$

Тоді семантичну функцію f^{**} кореляції документів, яка вра-

ховує всі взаємозв'язки досліджуваних документів, можна визначити наступним чином: $f^*(d_i, d_j) = \beta_{d_i, d_j} / \min(\beta_{d_i, T}, \beta_{d_j, T})$.

Відмітимо, що функції f^* і f^{**} узагальнюють кореляційні функції, у яких враховується тільки синонімія.

Третя глава присвячена опису реалізованого автором програмного прототипу системи адаптивного інформаційного пошуку та аналізу проведених експериментів. Даний прототип реалізовано в оточенні IPIC, усі програмні компоненти написані на мові Сі. Для роботи системи необхідно мати персональний комп'ютер, сумісний з IBM PC/AT у стандартній конфігурації, та операційну систему MS DOS версії 3.0 та вище. Для експериментів використовувалась БД в галузі обчислювальної техніки і програмного забезпечення (300 документів). При реалізації системи велика увага приділялася інтерфейсу користувача, який був зведений до того, що від абонента вимагалось лише здійснювати оцінку пред'явлених йому системою документів. Вся решта роботи по визначенню інформаційної потреби, формулюванню запиту, здійсненню пошукових ітерацій, упорядкуванню результатів пошуку перед переглядом з урахуванням інформаційних інтересів користувача здійснювалась автоматично самою системою. Для кожного абонента система створює і веде його персональну віртуальну БД.

Експерименти, проведені на даній системі, показали, що вона досить добре адаптується до інформаційної потреби користувача. Середній показник точності при майже 100 % повноті становив біля 50 % (у звичайних пошукових системах показник точності не перевершує 20%), характер оцінок документів у пред'являемих системою порціях мав коливальний вигляд, тобто. погано - краще - добре - погано - краще - добре - ... Причому перехід до наступного періоду означає, що система вже вичерпала розділяючі властивості найбільш інформативних термінів і намагається "зачепитися" за іншу лексику, семантично близьку до вибраної тематики.

У таблиці приведено оцінки документів, які давались у кожній порції, пред'являемої системою, та Стоп-правила, які використовувала система у випадку адаптації її на інформаційні інтереси у галузі мережевого забезпечення.

Після всіх шагів "рафінування" словника користувача він містив 124 терміна, які повністю покривали мережеву термінологію

БД (відмітимо, що на початковій стадії, після проглядання кластерного дерева, словник містив 11 термінів).

	Пор-ція	Оцінки документів										Стоп-правило
		d_0	d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9	
$\varphi(S_0) \neq \emptyset$	1	+	+	-	-	-	-	-				переупорядкувати
	2	+	-	+	+	+	+	-	-	+	-	переупорядкувати
	3	+	+	+	-	+	+	-	-	-	+	переупорядкувати
	4	+	+	+	+	+	+	+	-	+	+	продовжити
	5	+	+	+	+	+	+	+	-	-	-	продовжити
	6	-	-									новий запит
$\varphi(S_1) \neq \emptyset$	1	-	+									новий запит
$\varphi(S_2) \neq \emptyset$	1	+	+	-								новий запит
$\varphi(S_3) = \emptyset$												КІНЕЦЬ

"+" - документ відмічений як релевантний,

"-" - документ відмічений як нерелевантний.

Таблиця.

Основні висновки і результати, отримані автором по темі дисертації, полягають у наступному:

1. Автором запропонована нова схема інформаційного пошуку - адаптивний інформаційний пошук, - яка в розвитку ітераційних пошукових схем.

2. Для реалізації цієї схеми автором запропонована модель локалізації міжсезонної інформаційної потреби абонента. Запропоновані стратегії локалізації міжсезонної інформаційної потреби, вивчені види її мінливості. З цією метою введені поняття глобальної і локальної інформаційних структур, виділено три класи абонентів і знайдено умови віднесення абонента до того чи іншого класу на підставі співвідношень між цими структурами.

3. Для локалізації інформаційної потреби в межах одного сеансу зв'язку абонента з системою введено поняття внутрішньосезонної інформаційної потреби і запропоновано ймовірнісний підхід до побудови відповідних термінологічних словників.

4. Запропонована семантична модель інформаційно-пошукових процесів для задач кластеризації БД, локалізації інформаційної потреби користувача, пошуку і ранжування знайдених документів. Побудована всюди визначена функція сили зв'язку між двома термінами, на базі якої побудовані скалярно-семантична і семантична

функції кореляції документів.

5. Для апробації схеми адаптивного інформаційного пошуку розроблено і реалізовано прототип системи адаптивного пошуку. При реалізації системи велика увага приділялась інтерфейсу користувача, який був зведений до того, що від абонента вимагалось лише здійснювати оцінку пред'явлених йому системою документів. Вся решта роботи по визначенню інформаційної потреби, формулюванню запиту, здійсненню пошукових ітерацій, упорядкуванню результатів пошуку перед переглядом з урахуванням інформаційних інтересів користувача здійснювалась автоматично самою системою.

6. З метою перевірки адекватності моделей була проведена серія експериментів з використанням експериментальної БД в галузі обчислювальної техніки і програмування. Експерименти показали, що дана схема і використані автором для її реалізації моделі дають, по-перше, добрі показники повноти і точності інформаційного пошуку і, по-друге, термінологічний словник користувача, який будується системою, фактично покриває всю представницьку (природно, у статистичному розумінні) термінологію по вибраній темі в конкретній БД.

Основні результати дисертації опубліковані в роботах:

1. Дрянский В.М., Жогов Г.В. Модель локализации информационной потребности абонентов настраивающихся документальных информационно-поисковых систем // Кибернетика. - 1986. - № 4. - С. 81-84.

2. Дрянский В.М., Жогов Г.В. Система программных интерфейсов // Базы данных и знаний в автоматизированных региональных системах. - К.: Наукова думка, 1990. - С. 269 - 278.

3. Дрянский В.М., Жогов Г.В., Алешкина С.М. и др. Персональная электронная библиотека // Проблемы информатики (Научно-практическая конф. с межд. участием). - Самара-Астрахань-Самара, 11 - 18 мая 1991. - С. 35 - 37.

4. Дрянский В.М., Жогов Г.В., Катеринич Л.Г. Алгоритм принятия решения в моделях локализации входной информационной потребности // Проблемы создания ретроспективных поисковых массивов в автомат. центрах НТИ. Тез. док. XV Всес. научн. сем. "Систем. исследов. ГАСНТИ" (г.Рига, 20 - 22 мая 1985) Ч.1 - С. 134-135.

5. Дрянский В.М., Жогов Г.В., Катеринич Л.Г. Алгоритмы классификации в моделях локализации инф. потребности пользовате-

лей документальной информационно-поисковой системы // Математические методы в автоматизированных информационных системах и банках данных. - К.: ИК АН Украины, 1985. - 49 - 55.

6. Дриянский В.М., Жогов Г.В., Катеринич Л.Г. Концепция настраиваемого информационного поиска в АСИО по ВТ // Создание автоматизированных систем информационного обеспечения научных исследований. - К.: ИК АН Украины, 1986. - С. 27 - 38.

7. Дриянский В.М., Жогов Г.В., Катеринич Л.Г. Математические модели настраиваемого документального поиска // В кн. Автоматизация информ. поиска. - К.: Наукова думка, 1990. - С. 121 - 181.

8. Дриянский В.М., Жогов Г.В., Катеринич Л.Г. Нужен ли язык запросов в информационно-поисковых системах // УСИМ. - 1991. - N 7. - С. 158 - 160.

9. Дриянский В.М., Жогов Г.В., Капустин В. А., Катеринич Л.Г. Информационно-поисковая система с открытой архитектурой // Проблемы информатики города: Сб. науч. тр. - К.: Наукова думка, 1990. - С. 218 - 228.

10. Дриянский В.М., Жогов Г.В., Колтун А.Я. Количественные оценки парадигматических отношений и их использование в документальных ИПС // НТИ. Сер. 2. - 1985. - N 10. - С. 6 - 16.

11. Дриянский В.М., Жогов Г.В., Чалый И.А. Концепция операционной оболочки для создания информационных технологий // Проблемы разработки и внедрения программного обеспечения ЭВМ и систем. - К.: ИК АН Украины, 1988. - С. 24 - 30.

12. Дриянский В.М., Жогов Г.В., Чалый И.А. Мобильные инструментальные средства поддержки протоколов коммуникации прикладных процессов // Программирование. - 1990. - N 6. - С. 103-108.

Підписано до друку: 20 . 09 . 94г . Формат 60x84 1/16 .
Об'єм 1,05 д. а. Зак. № 4034 . Тираж 100 привітників .

Державне конунальне поліграфічне підприємство "Тираж"
м. Київ

454110

AB 31.003

AB 31.003