

VẤN ĐỀ TỐI THIỂU HÓA HỆ ĐIỀU CHỈNH THÔNG TIN

ĐỖ ĐỨC GIÁO, VŨ NGỌC LOÃN

Trong [7] tác giả đã đưa ra định nghĩa tập các cây thông tin TREE. Với mỗi cây $T \in \text{TREE}$, ta lập bộ ba $S(T) = [R, D_T, \sigma_T]$ và gọi là hệ điều chỉnh thông tin ở đây R — tập các công thức logic, D_T — tập tất cả các thông tin khác nhau từng cặp có mặt trong cây T , còn σ_T là ánh xạ từ tập R vào tập $2D_T$. Như vậy tập các hệ điều chỉnh thông tin $S(\text{TREE}) = [R, D_T, \sigma_T] / R \subseteq \text{FORM}, T \in \text{TREE}$ là một lớp các retrieval Systems mà SALTON đã định nghĩa trong [4]. Trong bài này, ta sẽ nghiên cứu vấn đề tối thiểu hóa cấu trúc hệ điều chỉnh thông tin ứng với một cây thông tin cho trước.

1. Cây thông tin với các tập công thức logic

Z — tập các số nguyên. D — tập các Documents. $D^+ = : DU[r] \tau$ — kí hiệu cây rỗng. Tập các cây thông tin định nghĩa trong [7] kí hiệu là TREE.

RESULT : $\text{TREE} \times Z \rightarrow D^+$ định nghĩa như sau :

1. RESULT $(d, 1) = d, \forall d \in D^+$ và $\forall 1 \in Z$.

2. RESULT $(K < T_1, T_2, T_3 >, 1) = \begin{cases} \text{RESULT}(T_1, 1) & \text{nếu } K < 1 \\ \text{RESULT}(T_2, 1) & \text{nếu } K = 1 \\ \text{RESULT}(T_3, 1) & \text{nếu } K > 1 \end{cases}$

Hai cây T_1 và T_2 là tương đương nhau ($T_1 \sim T_2$) $\Leftrightarrow \forall 1 \in Z$
đn

ta có RESULT $(T_1, 1) = \text{RESULT}(T_2, 1)$

Giả sử $p, q \in Z, x$ là đối số nhận giá trị trong Z .

INEQU = : $[x \leq p, x > q/p, q \in Z]$.

t là kí hiệu đúng (nhận giá trị 1). f là kí hiệu sai (nhận giá trị 0)

Ta định nghĩa các công thức :

1. Mỗi phần tử $u(x) \in \text{INEQU}$ gọi là một công thức (công thức cơ bản). Kí hiệu t, f cũng là những công thức.

2. Giả sử $H(x), H_1(x), H_2(x)$ là các công thức. Khi đó dãy kí hiệu

$(\neg H(x)), (H_1(x) \square H_2(x)),$ ở đây $\square \in [\vee, \wedge, \rightarrow, \leftrightarrow]$.

Cũng là những công thức. Các kí hiệu \neg (non), \vee (vel), \wedge (et), \rightarrow (Seq) và \leftrightarrow (Seq) và \leftrightarrow (Saq) là các kí hiệu hàm logic. Định nghĩa các hàm này xem trong [1].

Kí hiệu tập tất cả các công thức logic định nghĩa như trên là FORM. b là ánh xạ từ tập $\text{INEQU} \times Z$ vào tập $[0, 1]$ và được định nghĩa như sau :

$b(u(x), K) = \begin{cases} 1 & \text{nếu } u(K) \text{ đúng} \\ 0 & \text{nếu } u(K) \text{ sai,} \end{cases}$

ở đây $u(x) \in \text{INEQU}, k \in \mathbb{Z}$

Tất nhiên ta có thể mở rộng ánh xạ trên thành ánh xạ $\text{Val} : \text{FORM}_x \mathbb{Z} \rightarrow [0, 1]$ như sau:

$$1) \text{Val}(u(x), k) = b(u(x), k) \quad \forall u(x) \in \text{INEQU}, \forall k \in \mathbb{Z}.$$

$$\text{Val}(f, k) = 0, \forall k \in \mathbb{Z}.$$

$$\text{Val}(t, k) = 1, \forall k \in \mathbb{Z}.$$

$$2) \text{Val}(\neg H(x), k) = \text{non Val}(H_1(x), k)$$

$$\text{Val}(H_1(x) \square H_2(x), k) = 0 (\text{Val}(H_1(x), k) \text{ val } H_2(x), k) \text{ Val}(H_2(x), k); k$$

ở đây $(\square, 0) \in \{(\vee, \text{vel}), (\wedge, \text{et}), (\rightarrow, \text{Seq}), (-), \text{Seq}), (\leftrightarrow, \text{Saq})\}$

Giả sử $H_1(x), H_2(x) \in \text{FORM}$. Ta nói $H_1(x)$ là tương đương (tương đương yếu) với $H_2(x)$, kí hiệu $H_1(x) \sim H_2(x)$ ($H_1(x) \sim H_2(x)$) khi và chỉ khi

$\text{Val}(H_1(x), k) = \text{Val}(H_2(x), k)$ ($\text{val}(H_1(x), k) \leq (\text{val}(H_2(x), k)) \quad \forall k \in \mathbb{Z}$. Bây giờ ta định nghĩa ánh xạ $\text{RETR} : \text{TREE}_x \text{FORM} \rightarrow 2^{\mathbb{D}^+}$ như sau:

$$\text{RETR}(T, H(x)) = [\text{RESULT}(T, K) / k \in \mathbb{Z}, T \in \text{TREE}, \text{Val}(H(x), k) = 1].$$

Đặt $\text{RETR}(T, H(x)) =: \sigma_T(H)$. Khi đó dễ dàng thấy:

$$1. \text{Nếu } H_1(x) \sim H_2(x) \text{ thì } \sigma_T(H_1) = \sigma_T(H_2)$$

$$2. \text{Nếu } H_1(x) \sim H_2(x) \text{ thì } \sigma_T(H_1) \subseteq \sigma_T(H_2)$$

$$3. \sigma_T(f) = \phi \text{ (tập rỗng)}$$

$$4. \sigma_T(\neg H) = \sigma_T(x) \setminus \sigma_T(H)$$

$$5. \sigma_T(H_1 \vee H_2) = \sigma_T(H_1) \cup \sigma_T(H_2)$$

$$6. \sigma_T(H_1 \wedge H) = \sigma_T(H_1) \cap \sigma_T(H_2)$$

Giả sử $T, T' \in \text{TREE}$ và $R \subseteq \text{FORM}$.

Ta nói T là R -tương đương với T' ($T \approx_R T'$) $\Leftrightarrow \forall H \in R$ ta có

$$\text{RETR}(T, H) = \text{RETR}(T', H).$$

Hiển nhiên nếu $T \approx_R T'$ thì $T \approx_R T'$

2. Hệ điều chỉnh thông tin và vấn đề tối thiểu hóa cấu trúc của nó.

Với mỗi $T \in \text{TREE}$ và $R \subseteq \text{FORM}$, kí hiệu D_T là tập tất cả các thông tin có mặt trong T và khác nhau từng cặp.

$$\sigma_T : R \rightarrow 2^{D_T} \text{ ở đây } \sigma_T =: \text{TREE}_x R \rightarrow 2^{D_T}.$$

Khi đó bộ ba $S(T) = [R, D_T, \sigma_T]$ với R là tập các ngôn ngữ vào, D_T là tập các thông tin ra và σ_T là ánh xạ điều chỉnh thông tin) được gọi là hệ điều chỉnh thông tin ứng với cây T .

$$S(\text{TREE}) =: \{ [R, D_T, \sigma_T] / R \subseteq \text{FORM}, T \in \text{TREE} \}$$

là tập tất cả các hệ điều chỉnh thông tin ứng với TREE .

Định nghĩa 1.

$$\text{Giả sử } S(T) = [R, D_T, \sigma_T] \in S(\text{TREE})$$

1. Tập $A \subseteq D_T$ gọi là tập mô tả được trong $S(T) \Leftrightarrow \exists H \in R$ sao cho $\sigma_T(H) \stackrel{dn}{=} A$. Tập $D(S(T)) = [A \mid A \subseteq D_T \wedge \exists H \in R \wedge \sigma_T(H) = A]$ gọi là tập tất cả các tập mô tả được trong $S(T)$.

2. Hệ $S(T) = [R, D_T, \sigma_T]$ gọi là hệ mô tả được $\Leftrightarrow \forall d \in D_T \exists H \in R \wedge \sigma_T(H) \stackrel{dn}{=} [d]$.

3. Phần tử $d \in D_T$ là dư thừa trong $S(T) \Leftrightarrow \forall H \in R$ ta luôn luôn có $\sigma_T(H) \neq [d]$.

Định nghĩa 2.

1. Ta nói $S(T_1)$ tương đương với $S(T_2)$ ($S(T_1) \stackrel{dn}{\approx} S(T_2)$) $\Leftrightarrow T_1 \approx T_2$.

2. Ta nói $S(T_1)$ là R -tương đương với $S(T_2)$ ($S(T_1) \stackrel{R}{\approx} S(T_2)$) $\Leftrightarrow T_1 \stackrel{R}{\approx} T_2$.

3. Ta nói $S(T_1)$ đồng nhất với $S(T_2)$ ($S(T_1) \equiv S(T_2)$) $\Leftrightarrow T_1 \equiv T_2$ và $R_1 \equiv R_2$.

Định lý 1

1. $\forall S(T) \in S(TREE) \exists S(T') \in S(TREE)$ sao cho $S(T) \approx S(T')$ và $S(T')$ là mô tả được.

2. Hệ $S(T) = [R, D_T, \sigma_T]$ là mô tả được khi và chỉ khi mỗi tập con của D_T là mô tả được trong $S(T)$.

Chứng minh

1. Để chứng minh phần này, ta đưa vào khái niệm cây chuẩn N . N gọi là cây chuẩn nếu nó có dạng là d ($d \in D^+$) hoặc $K \langle d_0, d_1, d_2 \rangle$ với $d_0 \neq d_1$, hoặc $k \langle d_0, d_1, k+1 \langle \tau, d_2, \dots, k+s \langle \tau, d_{s+1}, d_{s+2} \rangle \dots \rangle$ với $d_0 \neq d_1$ và $d_{s+1} \neq d_{s+2}$ ($s \geq 1$). Đối với T trong $S(T)$ có tồn tại duy nhất cây chuẩn N sao cho $N \approx T$. Xem trong [7]. Kiểm tra lại $S(N) = [R, D_N, \sigma_N]$ là hệ mô tả được, và $S(T) \approx S(N)$.

2. Dựa vào định nghĩa hệ mô tả được và sử dụng các tính chất đã đưa ra của ánh xạ σ_T trong phần I.

Định nghĩa 3

Giả sử $S(T_0) \in S(TREE)$. Hệ $S(T_0) = [R, D_{T_0}, \sigma_{T_0}]$ được gọi là D -tối giản $\Leftrightarrow dn[D_{T_0}] \leq \min \{ [D_T] \mid S(T) = [R, D_T, \sigma_T] \in S(TREE) \}$
 $S(T_0) \approx S(\tau)$

Định lý 2

1. $\forall S(T) \in S(TREE) \exists S(T_0)$ sao cho $S(T) \approx S(T_0)$ và $S(T_0)$ là D -tối giản.

2. Hệ $S(T)$ là D -tối giản $\Leftrightarrow S(T)$ là mô tả được.

Chứng minh.

Phương pháp chứng minh như định lý 1

Định nghĩa 4

Giả sử $S(T) = [R, D_T, \sigma_T] \in S(TREE)$

Tập $L(S(T)) = [H(x) \mid H(x) \in R \wedge \exists A \subseteq D_T \wedge \sigma_T(H) = A]$ gọi là tập ngôn ngữ vào mô tả được trong $S(T)$.

Tập $L_0(S(T)) = \{H(x) \mid H(x) \in R \vee \exists d \in D_T \wedge \sigma_T(H) = [d]\}$ gọi là tập nhân của tập $L(S(T))$.

Chú ý $L_0(S(T)) \subseteq L(S(T))$

Định nghĩa 5

Giả sử $H_1(x), H_2(x) \in L_0(S(T))$ với $S(T) = [R, D_T, \sigma_T]$.

$H_1(x)$ là đồng dạng với $H_2(x)$ trên $L_0(S(T))$ (kí hiệu $H_1(x) \equiv H_2(x)$)

$\Rightarrow \exists d \in D_T \wedge \sigma_T(H_1) = \sigma_T(H_2) = [d]$.

Đn

Quan hệ \equiv trên $L_0(S(T))$ là quan hệ tương đương và nó tạo ra một phân hoạch tương đương. Mỗi lớp trong phân hoạch tương đương đó ta giữ lại một phần tử đại diện và tập các phần tử đó ta kí hiệu là $L_0(S(T))/\equiv$. Rõ ràng là $L_0(S(T))/\equiv \subseteq L(S(T) \subseteq L(S(T))$.

Định lý 3

Trên mỗi $(S(T))$ đều có tồn tại tập $L_0(S(T))/\equiv$ mô tả được trong $S(T)$

Chứng minh. Đối với T trong $S(T)$ có tồn tại cây chuẩn N (xem [7]). Dễ dàng kiểm tra lại $D_N = \bigcup \sigma_N(H)$, trong đó $L_0(S(N))/\equiv$ được xây dựng từ quá trình

$$H \in L_0(S(N))/\equiv$$

ây dựng dạng chuẩn N . Như vậy $L_0(S(N))/\equiv$ là mô tả được trong $S(N)$ mà $(T) \sim S(N)$ (Xem định lý 1) nên $L_0(S(N))/\equiv$ mô tả được trong $S(T)$.

Định nghĩa 6

Giả sử $S(N) = [R, D_N, \sigma_N] \in S(\text{TREE})$, N là cây chuẩn. Tập ngôn ngữ vào R được gọi là tập cốt yếu trong $S(N)$ nếu nó thỏa mãn các tính chất sau đây:

- 1) R chỉ gồm các công thức cơ bản trong INEQU (tính cơ bản).
- 2) $H \in R \Leftrightarrow \exists d \in D_N \wedge \sigma_N(H) = [d]$ (tính mô tả được).
- 3) Nếu bỏ đi khỏi R một công thức bất kỳ H thì tập $R \setminus \{H\}$ sẽ không mô tả được trong $S(N)$ tính không thừa).
- 4) Nếu thêm vào R một công thức $H \in \text{FORM}$ thì hoặc $H \in R$, hoặc $\exists H' \in R$ sao cho $H \equiv H'$ hoặc $R \cup \{H\}$ sẽ không thỏa mãn một trong ba tính chất trên. (Tính không thiếu).

Định lý 4

Nếu $N \in \text{TREE}$ là cây chuẩn thì tập $L_0(S(N))/\equiv$ là tập cốt yếu duy nhất trong $S(N)$

Chứng minh

Kiểm tra lại 4 tính chất của tập cốt yếu. Tính duy nhất suy từ tính duy nhất của dạng chuẩn N của T (xem trong [7]).

Định nghĩa 7

Hệ $S(T) = [R, D_T, \sigma_T]$ được gọi là tối giản nếu nó thỏa mãn

- 1) $S(T)$ là hệ D - tối giản
- 2) R là tập cốt yếu của $S(T)$.

Định lý 5

$\forall S(T) \in S(\text{TREE}) \exists$ hệ tối giản $S(T_0)$ sao cho $S(T) \sim S(T_0)$.

Chứng minh

Sử dụng các kết quả ở trên.

Chú ý: Dễ dàng chứng minh được nếu $S(T_1) \approx S(T_2)$ thì $S(T_1) \approx S(\Pi_2)$. Điều ngược lại nói chung không đúng với mọi $R \subseteq \text{FORM}$. Nhưng nếu cho $R_0 \in [L_0(S(N_1)) \equiv, L_0(S(N_2))/\equiv]$ thì ta có $S(T_1) \approx S(T_2) \leftrightarrow S(T_1) \approx S(T_2)$. Ở đây N_i là cây chuẩn của T_i ($i = 1, 2$).

TÀI LIỆU THAM KHẢO

1. G. Asser; Einführung in die mathematische logik. Teil 1, Leipzig, 1959.
2. W. Marek, Z. Pawlac, Bll. Acad. Polon., Sci., Sei. Math. Astronom. Phys. 22, 447. (1974)
3. W. Marek, Z. Pawlak, Z. Pawlak, Mathematical Foundations I, CC PA Reports, No 149, 1974
4. Salton G. - Automatic Information orgacisation and Retrieval New York 1968 (Mc Grow - Hill Book Company).
5. H. Thiele: Colloques Internationaux du C. N. R. S. No 296 (THEORIE DE L' INFORMATION 1977)
6. H. Thiele: Pr. IPI. Pan. 1980. No 411, 87 - 88.
7. Do Duc Giao. Diss. A, Humboldt - Uni. Berlin, 1986.
8. Do Duc Giao: Pr. IPI. PAN. 1980. No 411. 33 - 35.
9. Đỗ Đức Giáo: «Thông báo khoa học 1983», Khoa Toán Cơ trường Đại học Tổng hợp Hà nội, 1984, 64-70.
10. Đỗ Đức Giáo «Tập chí khoa học» Toán Lý, số 2/1985, 22-26.
11. Đỗ Đức Giáo «Tóm tắt báo cáo Hội nghị toán học toàn quốc lần thứ 12 (22-25-1985), trang 97.

До Дык Зао, Ву Нгок Loan

ОБ ПРОБЛЕМЕ МИНИМИЗАЦИИ СИСТЕМЫ РЕГУЛИРУЕМОЙ ИНФОРМАЦИИ

Трио $S(T) = [R, D_T, \sigma_T]$ (T -TREE) называется система регулируемой информации тогда и только тогда, когда R и D - непустые множества, а σ_T преобразование из R в множество всех подмножеств множества D_T . Каждый элемент из R является формулой и множество R называется множеством языков системы $S(T)$. D_T есть множество documents присутствующих в дереве T и называется множеством documents системы $S(T)$. σ_T - преобразование регулирования системы $S(T)$. Для кажзого элемента $H \in R$, $\sigma_T(H)$ - подмножество documents в D_T . В этой работе рассматривается проблема минимизации системы регулируемой информации.

Do Duc Giao, Vu Ngoc Loan

MINIMIZATION PROPLEM FOR INFORMATION RETRIEVAL SYSTEM

A triple $S(T) = [R, D_T, \sigma_T]$ is called a information retrieval system if and only if R and D_T are non - empty sets and a mapping $\sigma_T : R \rightarrow 2^{D_T}$. The elements of R are the formulas and the Set R is said to be the language of $S(T)$. D_T is the set of the documents in the tree T and the set D_T is interpreted as the document set of $S(T)$. σ_T is calles the retrieval function of $S(T)$ which give for each $H \in R$ the subset $\sigma_T(H) \subseteq D_T$ of documents retrieved by σ_T . In this paper we shall discuss the problem of minimization for the information retrieval systems consireded.

Nhận bài ngày 20-4-1986