

## MỘT SỐ KẾT QUẢ VỀ HIỆU QUẢ CỦA THUẬT TOÁN NHẬN DẠNG TỔ HỢP

NGÔ QUỐC TẠO, HOÀNG KIẾM

### I - ĐẶT VẤN ĐỀ

Phương pháp tổ hợp các thuật toán nhằm nhận được thuật toán có hiệu quả hơn hiệu quả của mỗi thuật toán ban đầu, cụ thể là, nếu xác suất sai lầm của mỗi thuật toán là  $\epsilon_i$  thì xác suất sai lầm của thuật toán tổ hợp có thể nhỏ hơn  $\epsilon_i$ . Tính cho đến nay đã xuất hiện nhiều công trình để cập đến thuật toán nhận dạng tổ hợp với các công thức tính  $R_n$  [1], [2]. Song đó là những công thức đánh giá chưa được chứng minh chặt chẽ. Trong bài này trình bày một số đánh giá cho  $R_n$  chặt hơn một số kết quả đã công bố.

Phương pháp tổ hợp các thuật toán được dựa trên quy tắc quyết định theo đa số được phát biểu như sau:

Giả sử cho trước  $N$  thuật toán nhận dạng  $\mathcal{A}_1, \dots, \mathcal{A}_N$  phân loại các dạng  $\{S_i\}_{i=1, \dots, n}$  thành 2 lớp  $K_1$  và  $K_2$ . Người ta xác định các ma trận  $\|a_{ij}^k\|_{n \times 2}, k=1, \dots, N$  như sau:

$$a_{ij}^k = \begin{cases} 1, & \text{nếu thuật toán } \mathcal{A}_k \text{ xếp } S_i \text{ vào lớp } K_j \\ 0, & \text{ngược lại.} \end{cases}$$

Quy tắc quyết định của thuật toán nhận dạng tổ hợp như sau:

$$D_j(S_i) = \begin{cases} 1, & \text{nếu } \sum_{i=1}^n a_{ij}^k > E \\ 0, & \text{ngược lại.} \end{cases}$$

Ở đây ký hiệu  $E = \lfloor LN/2 \rfloor$  (số nguyên tố lớn nhất không vượt quá  $N/2$ ).

$D_j(S_i) = 1$  có nghĩa là thuật toán nhận dạng tổ hợp xếp  $S_i$  vào lớp  $K_j$ ,  $i=1, \dots, n$ ;  $j = 1, 2$ . Khi các thuật toán  $\mathcal{A}_1, \dots, \mathcal{A}_n$  độc lập và mỗi thuật toán có xác suất sai lầm  $0 < \epsilon < 1/2$  thì xác suất sai lầm của thuật toán nhận dạng tổ hợp được xác định theo công thức sau:

$$R_N = \sum_{k>E} C_N^k \epsilon^k (1-\epsilon)^{N-k}$$

Trong lì sẽ chỉ ra rằng  $\forall \epsilon, 0 < \epsilon < 1/2$  thì  $R_N < \epsilon$ . (Bất đẳng thức này chứng tỏ rằng thuật toán tổ hợp có hiệu quả hơn mỗi thuật toán đã cho). Và  $R_N \leq \exp\{-N \times C(\epsilon)\}$  (tốc độ tiến tới 0 của  $R_N$  có cấp hàm mũ theo  $N$ ).

Ta ký hiệu  $\vec{\lambda}$  là vector có  $N$  tọa độ nhận giá trị 0,1 tức là

$$\vec{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_N)$$

$$\lambda_k \in \{0, 1\}, k = 1, \dots, N$$

và

$$TA = \left\{ \vec{\lambda} \mid \sum_{k=1}^N \frac{\lambda_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} = \frac{1}{2} \sum_{k=1}^N \frac{1}{\sqrt{\epsilon_k(1-\epsilon_k)}} \right\}$$

Trong trường hợp các thuật toán độc lập và có xác suất sai lầm  $\epsilon$  tương ứng, thì xác suất sai lầm  $R_N$  có thể xác định theo công thức:

$$R_N = \sum_{\lambda \in TA} \prod_{k=1}^N \lambda_k^k (1 - \lambda_k)^{1-\lambda_k}$$

Hơn nữa

$$R_N \leq \frac{4}{N^3} \left( \sum_{k=1}^N \sqrt{\epsilon_k(1-\epsilon_k)} \right)^2 / \left( 1 - 2 \sum_{k=1}^N \epsilon_k/N \right)^2.$$

Các khẳng định này được chỉ ra trong III.

## II - ĐÁNH GIÁ XÁC SUẤT SAI LẦM CỦA THUẬT TOÁN NHẬN DẠNG TỔ HỢP ĐỐI VỚI N THUẬT TOÁN ĐỘC LẬP CÓ CÙNG XÁC SUẤT SAI LẦM GIỐNG NHAU $\epsilon$ .

Bài đề 1 [1]: Nếu cho trước  $N$  thuật toán nhận dạng độc lập có cùng xác suất sai lầm  $\epsilon$ , thì xác suất sai lầm của thuật toán nhận dạng tổ hợp thỏa mãn:

$$R_N = \sum_{k>E} C_N^k \epsilon^k (1-\epsilon)^{N-k}.$$

Mệnh đề 1: Với  $\epsilon$  thỏa mãn  $0 < \epsilon < 1/2$  thì xác suất sai lầm của thuật toán nhận dạng tổ hợp thỏa mãn:

$$R_N < \epsilon.$$

**Chứng minh:** Trước hết ta nhận thấy hàm  $f_{a, b}(\epsilon) = \epsilon^a(1-\epsilon)^b$ , với  $0 < b \leq a$ , tăng chéo trên  $[0, 1/2]$ , tức là  $\forall \epsilon_1, \epsilon_2$  thỏa mãn  $0 \leq \epsilon_1 < \epsilon_2 \leq 1/2$

thì  $f_{a, b}(\epsilon_1) < f_{a, b}(\epsilon_2)$ .

$$\text{Đặt } T_N = \sum_{k>E} C_N^k \epsilon^{k-1} (1-\epsilon)^{N-k}, \text{ do bài đe 1 ta có } R_N = T_N \cdot \epsilon.$$

Do nhận xét trên, với  $k > E$  hàm  $f_{k-1, N-k}(\epsilon)$  tăng chéo trên  $[0, 1/2]$ . Do đó  $f_{k-1, N-k}(\epsilon) \leq f_{k-1, N-k}(1/2)$ .

$$\text{Hay } \epsilon^{k-1} (1-\epsilon)^{N-k} \leq \left(\frac{1}{2}\right)^{k-1} \left(1 - \frac{1}{2}\right)^{N-k} = 2^{N-1}.$$

Vì thế

$$T_N \leq \sum_{k>E} C_N^k \left(\frac{1}{2}\right)^{N-1} = \left(\frac{1}{2}\right)^{N-1} \sum_{k>E} C_N^k.$$

Dễ dàng nhận thấy rằng:

$$\sum_{k>E} C_N^k \leq 2^{N-1},$$

Do đó  $T_N < 1$  hay  $R_N < \epsilon$ . Mệnh đề được chứng minh.

Mệnh đề 2: Với các điều kiện tương tự như bài đe 1, ta nhận được

$$R_N < \exp\{-N \times c(\epsilon)\}, \text{ với } c(\epsilon) = \ln \frac{1}{2\sqrt{\epsilon(1-\epsilon)}}.$$

**Chứng minh:** Trước hết ta thấy rằng với  $0 < \epsilon < 1/2$

$$\text{thì } \ln \frac{1-\epsilon}{\epsilon} > 0 \quad \text{và} \quad \ln \frac{1}{2\sqrt{\epsilon(1-\epsilon)}} > 0.$$

$\forall \alpha > 0, \forall k > E$  ta đều có:

$\exp(\alpha k) \geq \exp(\alpha \cdot N/2)$  hay  $\exp(-\alpha \cdot N/2 + \alpha \cdot k) \geq 1$ . Kết hợp với bài đe 1 ta suy ra:

$$R_N < \epsilon \exp(-\alpha \cdot N/2) \sum_{k>E} C_N^k \epsilon^k (1-\epsilon)^{N-k} \exp(\alpha \cdot k) =$$

$$= \exp(-\alpha \cdot N/2) \sum_{k=1}^N C_N^k \epsilon^k (1-\epsilon)^{N-k} \exp(\alpha \cdot k) = \left(\epsilon \cdot \exp\left(\frac{\alpha}{2}\right) + (1-\epsilon) \exp\left(-\frac{\alpha}{2}\right)\right)^N.$$

Chọn  $\alpha = \ln \frac{1-\varepsilon}{\varepsilon}$  suy ra  $R_N \leq (2\sqrt{\varepsilon(1-\varepsilon)})^N$ . Đặt  $c(\varepsilon) = \ln \frac{1}{2\sqrt{\varepsilon(1-\varepsilon)}}$  ta được

$R_N < \exp \{-N \cdot c(\varepsilon)\}$ . Mệnh đề được chứng minh. Như vậy, nếu như ta tò hợp càng nhiều thuật toán thành toán tò hợp thì hiệu quả của thuật toán tò hợp càng cao. Phương pháp này được ứng dụng rộng rãi trong trí tuệ nhân tạo và chẩn đoán.

### III - ĐÁNH GIÁ XÁC SUẤT SAI LẦM CỦA THUẬT TOÁN NHẬN DẠNG TÒ HỢP TRONG TRƯỜNG HỢP CÁC XÁC SUẤT SAI LẦM CỦA MỖI THUẬT TOÁN KHÁC NHAU

Quy tắc quyết định của thuật toán nhận dạng tò hợp như sau:

$$D_j(S_i) = \begin{cases} 1, & \text{nếu } \sum_{k=1}^N \frac{a_{ij}^k}{\sqrt{\varepsilon_k(1-\varepsilon_k)}} > \frac{1}{2} \sum_{k=1}^N \frac{1}{\sqrt{\varepsilon_k(1-\varepsilon_k)}} \\ 0, & \text{ngược lại} \end{cases}$$

Bđ đ# 2: Nếu mỗi thuật toán ban đầu có xác suất sai lầm  $\varepsilon_i$  tương ứng, thì xác suất sai lầm của thuật toán tò hợp được xác định như sau:

$$R_N = \sum_{\lambda \in TA} \prod_{k=1}^N \varepsilon_k^{\lambda_k} (1-\varepsilon_k)^{1-\lambda_k}$$

$\rightarrow$   $\lambda$  và TA được xác định như trong I.

Chứng minh: Dễ dàng suy ra từ tính độc lập của các thuật toán đã cho. Bđ đ# 1 là trường hợp đặc biệt của bđ đ# 2 khi  $\varepsilon_i = \varepsilon$ ,  $\forall i$ .

Mệnh đ# 3: Với các điều kiện tương tự như bđ đ# 2 thỏa mãn thì

$$R_N \leq \frac{4}{N^3} \left( \sum_{k=1}^N \sqrt{\varepsilon_k(1-\varepsilon_k)} \right)^2 \times \left( 1 + 2 \sum_{k=1}^N \tau_k / N \right)^2$$

Chứng minh:

Trước hết ta chứng minh bđ đ# sau:

Bđ đ# 3: Cho hai dãy số thực  $a_1, a_2, \dots, a_N$  và  $b_1, b_2, \dots, b_N$ .

1) Nếu  $\forall i, j : i > j, a_i \leq a_j, b_i \geq b_j$  thì

$$\sum_{i=1}^N a_i \sum_{j=1}^N b_j \geq N \sum_{i=1}^N a_i b_i$$

2) Nếu  $\forall i, j : (a_i - a_j)(b_i - b_j) \leq 0$  thì

$$\sum_{i=1}^N a_i \sum_{j=1}^N b_j \geq N \sum_{i=1}^N a_i b_i$$

Chứng minh: 1) chính là bất đẳng thức Trèbusep.

2) Chỉ cần sắp xếp dãy  $\{b_i\}$  theo thứ tự tăng dần thì ta cũng nhận được  $\{a_i\}$  có thứ tự hoán vị tương ứng với  $\{b_i\}$  sẽ giảm dần. Áp dụng 1) ta được 2) - điều phải chứng minh.

*Chứng minh mệnh đề:*

Từ bđ đe 2 ta suy ra được

$$\begin{aligned} R_N &\leq \left( \frac{1}{2} \sum_{k=1}^N \frac{1-2\epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} \right)^{-2} \sum_{\lambda \in T \Lambda} \prod_{k=1}^N \epsilon_k^{\lambda_k} (1-\epsilon_k)^{1-\lambda_k} \\ &\quad \times \left( \sum_{k=1}^N \frac{\lambda_k - \epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} \right)^2 \\ &\rightarrow R_N \leq \left( \frac{1}{2} \sum_{k=1}^N \frac{1-2\epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} \right)^{-2} \sum_{\lambda \in \{0,1\}^N} \prod_{k=1}^N \epsilon_k^{\lambda_k} (1-\epsilon_k)^{1-\lambda_k} \\ &\quad \times \left( \sum_{k=1}^N \frac{1-\epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} \right)^2. \end{aligned}$$

Dễ dàng thấy rằng :

$$\sum_{\lambda} \prod_{k=1}^N \epsilon_k^{\lambda_k} (1-\epsilon_k)^{1-\lambda_k} \left( \sum_{k=1}^N \frac{1-\epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} \right)^2 = N.$$

Từ đó ta có bất đẳng thức :

$$R_N \leq \left( \frac{1}{2} \sum_{k=1}^N \frac{1-2\epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} \right)^{-2}$$

Đặt  $a_k = \frac{1-2\epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}}, b_k = \sqrt{\epsilon_k(1-\epsilon_k)}$  và áp dụng bđ đe 3 ta có :

$$\begin{aligned} \sum_{k=1}^N \frac{1-2\epsilon_k}{\sqrt{\epsilon_k(1-\epsilon_k)}} \sum_{k=1}^N \sqrt{\epsilon_k(1-\epsilon_k)} &\geq N \sum_{k=1}^N (1-2\epsilon_k) \\ &= N^2 (1-2 \sum_{k=1}^N \epsilon_k/N). \end{aligned} \tag{2}$$

Áp dụng (1) và (2) ta được

$$R_N \leq \frac{4}{N^3} \left( \sum_{k=1}^N \sqrt{\epsilon_k(1-\epsilon_k)} \right)^2 / (1-2 \sum_{k=1}^N \epsilon_k/N)^2,$$

Nhận xét : -- Nếu áp dụng bất đẳng thức Svac cho  $2N$  số thực ta có :

$$\sum_{k=1}^N \sqrt{\epsilon_k(1-\epsilon_k)} \leq \sqrt{\sum_{k=1}^N \epsilon_k} \sqrt{\sum_{k=1}^N (1-2\epsilon_k)}$$

$$\text{Ta thu được } R_N \leq \frac{4}{N} \frac{\bar{\epsilon}(1-\bar{\epsilon})}{(1-2\bar{\epsilon})^2}; \text{ trong đó } \bar{\epsilon} = \sum_{k=1}^N \epsilon_k,$$

$$\text{Trong trường hợp } \epsilon_1 = \epsilon, \forall i \text{ ta thu được } R_N \leq \frac{4}{N} \frac{\epsilon(1-\epsilon)}{(1-2\epsilon)^2}.$$

Như vậy trong trường hợp các  $\epsilon_i = \epsilon, \forall i$ , thì kết quả II mạnh hơn III.

Một số vấn đề cần mở :

- Phát biểu định nghĩa cho thuật toán đệ quy lập.

— Trong trường hợp có định nghĩa về thuật toán độc lập thì hiệu quả của thuật toán tề hợp từ những thuật toán không độc lập sẽ ra sao? Chúng tôi chân thành cảm ơn Tiến sĩ Bạch Hưng Khang và phòng nhận dạng Viện KHTTDK về sự giúp đỡ nhiệt tình và có những gợi ý quý báu cho bài báo này. Chúng tôi xin cảm ơn đồng nghiệp Nguyễn Thành Thủy trường đại học Bách Khoa Hà Nội, đã có những ý kiến bđ ích cho bài báo.

Nhận ngày 8.3.1986

### TÀI LIỆU THAM KHẢO

1. Hoàng Kiếm, Some methods improving the efficiency of pattern recognition algorithms. Computer and artificial intelligence, 3(1984), №4, pp. 347–359, Czechoslovakia.
2. Gaillat G., borne Superleuse de la probabilité d'erreur avec la règle de décision majoritaire. Congrès AFCET-IRIA, Paris 1978.
3. ЖУРАВЛЁВ Ю. И., Об алгебраическом походе к решению задач распознавания или классификации. ВКН: Побл. Кибернетики, Выч. зз. М.; Наука, 1978 с. 5–68.
4. РЯЗАНОВ В.В., Комитетный синтез алгоритмов распознавания и классификации. Ж. вычисл. Матем. и матем. физ., 1981, Т. 21, №6. с. 1538–1549.
5. РЯЗАНОВ В.В., О синтезе классифицирующих алгоритмов на конечных множествах алгоритмов классификации. Ж. вычисл. Матем. и Матем. физ., Т.32, №2. с. 429–440.

### ABSTRACT

#### Some results of the efficiencies of the combined pattern recognition algorithm

The paper presents the evalution of the error probability of the combined pattern recognition algorithm based on N. independent algorithms. These main results are follows: If every given algorithm has error probability  $0 < \varepsilon < 1/2$  then  $R_N < \exp \{ -N \cdot c(\varepsilon) \}$ , here

$$c(\varepsilon) = \ln \frac{1}{2\sqrt{\varepsilon(1-\varepsilon)}}.$$

When all of the given algorithms have error probability  $0 < \varepsilon_i < 1/2$  respectively, then the error probability of combined pattern recognition algorithms satisfies

$$R_N \leq \frac{4}{N^3} \left( \sum_{i=1}^N \sqrt{\varepsilon_i(1-\varepsilon_i)} \right)^2 / (1 - 2 \sum_{i=1}^N \varepsilon_i/N)^2,$$