



Thuật toán ILA

Thuật toán nâng cao (Đại học Thủy lợi)



Scan to open on Studocu

Thuật toán ILA(Inductive Learning Algorithm) được dùng để xác định các luật phân loại cho tập hợp các mẫu học. Thuật giải này thực hiện theo cơ chế lặp, để tìm luật riêng đại diện cho tập mẫu của từng lớp. Sau khi xác định được luật, ILA loại bỏ các mẫu liên quan khỏi tập mẫu, đồng thời thêm luật mới vào tập luật. Kết quả có được là một danh sách có thứ tự các luật chứ không là một cây quyết. Các ưu điểm của thuật giải này có thể được trình bày như sau:

-Dạng các luật sẽ phù hợp cho việc khảo sát dữ liệu, mô tả mỗi lớp một cách đơn giản để dễ phân biệt với các lớp khác.

- Tập luật được sắp thứ tự, riêng biệt cho phép quan tâm đến một luật tại thời điểm bất kỳ. Khác với việc xử lý luật theo phương pháp cây quyết định, vốn rất phức tạp trong trường hợp các nút cây trở nên khá lớn.

Bước 1: Chia mẫu ban đầu thành n bảng con. Mỗi bảng con ứng với một giá trị của thuộc tính quyết định của tập mẫu

Thực hiện lần lượt các bước từ 2 đến 8 cho mỗi bảng con có được

Bước 2: $j=1$ (j là số thuộc tính của tổ hợp T)

Bước 3: Trên mỗi bảng con đang khảo sát, chia danh sách các thuộc tính thành các tổ hợp khác nhau, mỗi tổ hợp bao gồm j thuộc tính

Bước 4: Với mỗi tổ hợp thuộc tính có được tính số lần giá trị thuộc tính xuất hiện theo cùng tổ hợp thuộc tính trong các dòng còn lại của bảng con đang xét (mà đồng thời không xuất hiện tổ hợp giá trị này trên tất cả các bảng còn lại).

Tổ hợp T ^{Thuộc tính}
Giá trị của thuộc tính

Gọi tổ hợp đầu tiên(trong bảng con) có số lần xuất hiện nhiều nhất là tổ hợp lớn nhất.

Bước 5: Nếu tổ hợp lớn nhất có giá trị bằng 0, tăng j lên 1 và quay lại bước 3.

Bước 6:Loại bỏ các dòng thỏa mãn tổ hợp lớn nhất ra khỏi bảng con đang xử lý

Bước 7: Thêm luật mới vào tập luật R , với vế trái là tập các thuộc tính của tổ hợp lớn nhất(Kết hợp các thuộc tính bằng toán tử AND) và vế phải là giá trị thuộc tính quyết định tương ứng.

Bước 8: Nếu tất cả các dòng đều đã được loại bỏ, tiếp tục thực hiện bước 2 cho các bảng còn lại. Ngược lại(nếu còn dòng chưa bị loại bỏ) thì quay lại bước 4. Nếu tất cả các dòng con đã được xét thì kết thúc. Tập R chính là tập luật cần tìm.

2. Minh họa thuật toán:

Minh họa giải thuật ILA cho bảng dữ liệu sau đây:

Tên	Nước tiểu	Giảm cân	Đau lưng	Sốt	Kết quả
An	Vàng	Không	Không	Có	Không bệnh sỏi thận
Cường	Vàng	Không	Không	Không	Không bệnh sỏi thận
Châu	Có vôi	Không	Không	Có	Bệnh sỏi thận
Dung	Có máu	ít	Không	Có	Bệnh sỏi thận
Diễm	Có máu	Nhanh	Có	Có	Bệnh sỏi thận
Hương	Có máu	Nhanh	Có	Không	Không bệnh sỏi thận
Hoa	Có vôi	Nhanh	Có	Không	Bệnh sỏi thận
Phương	Vàng	ít	Không	Có	Không bệnh sỏi thận
My	Vàng	Nhanh	Có	Có	Bệnh sỏi thận
Nhung	Có máu	ít	Có	Có	Bệnh sỏi thận
Thu	Vàng	ít	Có	Không	Bệnh sỏi thận
Thương	Có vôi	ít	Không	Không	Bệnh sỏi thận
Tuấn	Có vôi	Không	Có	Có	Bệnh sỏi thận
Tùng	Có máu	ít	Không	Không	Không bệnh sỏi thận

Chia bảng mẫu thành 2 bảng con bởi 2 loại quyết định: “Bệnh sỏi thận” và “Không bệnh sỏi thận” như sau:

Bảng 1: “Bệnh sỏi thận”

Tên	Nước Tiểu	Giảm Cân	Đau Lưng	Sốt	Kết Quả
Châu	Có vôi	Không	Không	Có	Bệnh Sỏi Thận
Dung	Có máu	Ít	Không	Có	Bệnh Sỏi Thận
Diễm	Có máu	Nhanh	Có	Có	Bệnh Sỏi Thận
Hoa	Có vôi	Nhanh	Có	Không	Bệnh Sỏi Thận
My	Vàng	Nhanh	Có	Có	Bệnh Sỏi Thận
Nhung	Có máu	Ít	Có	Có	Bệnh Sỏi Thận
Thu	Vàng	Ít	Có	Không	Bệnh Sỏi Thận
Thương	Có vôi	Ít	Không	Không	Bệnh Sỏi Thận
Tuấn	Có vôi	Không	Có	Có	Bệnh Sỏi Thận

Bảng 2: “Không bệnh sỏi thận”

Tên	Nước Tiểu	Giảm Cân	Đau Lưng	Sốt	Kết Quả
An	Vàng	Không	Không	Có	Không Bệnh Sởi Thận
Cường	Vàng	Không	Không	Không	Không Bệnh Sởi Thận
Hương	Có máu	Nhanh	Có	Không	Không Bệnh Sởi Thận
Phương	Vàng	Ít	Không	Có	Không Bệnh Sởi Thận
Tùng	Có máu	Ít	Không	Không	Không Bệnh Sởi Thận

Bảng 1: “Bệnh sởi thận”

Tổ hợp: T

Với $j=1$, có 4 tổ hợp:

- {Nước tiểu}
 - {Giảm cân}
 - {Đau lưng}
 - {Sốt}
- Với tổ hợp {Nước tiểu}: Thuộc tính “Có vôi” xuất hiện 4 lần trong bảng 1 và không xuất hiện trong bảng 2. Thuộc tính “Có máu” và “Vàng” xuất hiện trên cả hai bảng.
- $$T_{\text{Nước tiểu}}^{\text{Có vôi}} = 4$$
- $$T_{\text{Nước tiểu}}^{\text{Có máu}} = 0$$
- $$T_{\text{Nước tiểu}}^{\text{Vàng}} = 0$$
- Với tổ hợp {Giảm cân}: Thuộc tính “Không”, “Nhanh”, “Ít” xuất hiện trên cả hai bảng.
- $$T_{\text{Giảm cân}}^{\text{Không}} = 0$$
- $$T_{\text{Giảm cân}}^{\text{Nhanh}} = 0$$
- $$T_{\text{Giảm cân}}^{\text{Ít}} = 0$$
- Với tổ hợp {Đau lưng}: Thuộc tính “Không”, “Có” xuất hiện trên cả hai bảng.
- $$T_{\text{Đau lưng}}^{\text{Không}} = 0$$
- $$T_{\text{Đau lưng}}^{\text{Có}} = 0$$
- Với tổ hợp {Sốt}: Thuộc tính “Không”, “Có” xuất hiện trên cả hai bảng.

$$T_{\text{Không}}^{\text{Sốt}} = 0$$

$$T_{\text{Có}}^{\text{Sốt}} = 0$$

⇒ Ta có $T_{\text{Có vôi}}^{\text{Nước tiểu}} = 4$ là lớn nhất. Ta chọn $T_{\text{Có vôi}}^{\text{Nước tiểu}}$

RULE 1: IF Nước tiểu = có vôi THEN Kết quả = Bệnh sỏi thận

Tiếp theo ta loại bỏ những dòng thỏa mãn tổ hợp lớn nhất tương ứng với Nước tiểu = Có vôi ra khỏi bảng 1 ta có bảng sau:

Tên	Nước tiểu	Giảm cân	Đau lưng	Sốt	Kết quả
Dung	Có máu	Ít	Không	Có	Bệnh sỏi thận
Diễm	Có máu	Nhanh	Có	Có	Bệnh sỏi thận
My	Vàng	Nhanh	Có	Có	Bệnh sỏi thận
Nhung	Có máu	Ít	Có	Có	Bệnh sỏi thận
Thu	Vàng	Ít	Có	Không	Bệnh sỏi thận

Các dòng trong bảng trên mọi giá trị của thuộc tính đều xuất hiện trong cả hai bảng(mọi giá trị T đều bằng 0) nên ta sẽ tăng j lên 1.

Với j=2, có 6 tổ hợp:

- {Nước tiểu, Giảm cân }
- {Nước tiểu, Đau lưng}
- {Nước tiểu, Sốt}
- {Giảm cân, Đau lưng}
- {Giảm cân, Sốt}
- {Đau lưng, Sốt}

$$T_{\text{Có máu}, \text{Ít}}^{\text{Nước tiểu, Giảm cân}} = 0$$

$$T_{\text{Có máu}, \text{Nhanh}}^{\text{Nước tiểu, Giảm cân}} = 0$$

$$T_{\text{Vàng}, \text{Ít}}^{\text{Nước tiểu, Giảm cân}} = 0$$

$$T_{\text{Vàng}, \text{Nhanh}}^{\text{Nước tiểu, Giảm cân}} =$$

$$T_{\text{Có máu}, \text{Không}}^{\text{Nước tiểu, Đau lưng}} = 0$$

$$T_{\text{Có máu}, \text{Có}}^{\text{Nước tiểu, Đau lưng}} = 0$$

T Nước tiểu Vàng , Đau lưng Có = 1

T Nước tiểu Vàng , Sốt Không = 0

T Nước tiểu Có máu , Sốt Có = 3

T Giảm cân Ít , Đau lưng Không = 0

T Giảm cân Nhanh , Đau lưng Có = 0

T Giảm cân Ít , Đau lưng Có = 2

T Giảm cân Ít , Sốt Có = 0

T Giảm cân Ít , Sốt Không = 0

T Giảm cân Nhanh , Sốt Có = 1

T Đau lưng Không , Sốt Có = 0

T Đau lưng Có , Sốt Có = 2

T Đau lưng Có , Sốt Không = 0

⇒ Ta có T Nước tiểu Có máu , Sốt Có = 3 là lớn nhất, ta chọn T Nước tiểu Có máu , Sốt Có

RULE 2: IF Nước tiểu= Có máu AND Sốt = Có THEN Kết Quả = Bệnh sỏi thận

Kể tiếp, loại bỏ những dòng ứng với Nước tiểu = có máu và sốt = có ra khỏi bảng ta được:

Tên	Nước tiểu	Giảm cân	Đau lưng	Sốt	Kết quả
My	Vàng	Nhanh	Có	Có	Bệnh sỏi thận
Thu	Vàng	Ít	Có	Không	Bệnh sỏi thận

Với j=2, ta có 6 tổ hợp mỗi tổ hợp gồm 2 thuộc tính:

- {Nước tiểu, Giảm cân }
- {Nước tiểu, Đau lưng}
- {Nước tiểu, Sốt}
- {Giảm cân, Đau lưng}
- {Giảm cân, Sốt}
- {Đau lưng, Sốt}
-

T Nước tiểu vàng, giảm cân ít=0
T Nước tiểu vàng, giảm cân Nhanh=1

T Nước tiểu vàng, đau lưng có=2

T Nước tiểu vàng, sốt không=0
T Nước tiểu vàng, sốt Có=0

T Giảm cân ít, đau lưng có=1

T Giảm cân ít, sốt Không=0
T Giảm cân Nhanh, sốt Có=1

T Đau lưng có, sốt không=0
T Đau lưng có, sốt có=1

⇒ Ta có T Nước tiểu vàng, đau lưng có=2 là lớn nhất. Ta chọn T Nước tiểu vàng, đau lưng có và ta có luật:

RULE 3: IF Nước tiểu = Vàng AND Đau lưng=Có THEN Kết quả=Bệnh sỏi thận

Loại bỏ các dòng tương ứng với Nước tiểu = vàng, đau lưng=có, như vậy tất cả các dòng trong bảng 1 bị loại bỏ

Bảng 2: “Không bệnh sỏi thận”

Tên	Nước Tiểu	Giảm Cân	Đau Lưng	Sốt	Kết Quả
An	Vàng	Không	Không	Có	Không Bệnh Sỏi Thận
Cường	Vàng	Không	Không	Không	Không Bệnh Sỏi Thận
Hương	Có máu	Nhanh	Có	Không	Không Bệnh Sỏi Thận
Phương	Vàng	Ít	Không	Có	Không Bệnh Sỏi Thận
Tùng	Có máu	Ít	Không	Không	Không Bệnh Sỏi Thận

Trong bảng 2, mọi giá trị của thuộc tính đều xuất hiện trong cả hai bảng(mọi giá trị T đều bằng 0) nên ta sẽ tăng j lên 1.

Với j=2, có 6 tổ hợp mỗi tổ hợp có 2 thuộc tính:

- {Nước tiểu, Giảm cân }
- {Nước tiểu, Đau lưng}
- {Nước tiểu, Sốt}
- {Giảm cân, Đau lưng}
- {Giảm cân, Sốt}
- {Đau lưng, Sốt}

$$T(\text{nước tiểu có máu, giảm cân ít})=0$$

$$T(\text{nước tiểu có máu, giảm cân nhanh})=0$$

$$T(\text{nước tiểu vàng, giảm cân ít})=0$$

$$T(\text{nước tiểu vàng, giảm cân không})=2$$

$$T(\text{nước tiểu có máu, đau lưng không})=0$$

$$T(\text{nước tiểu có máu, đau lưng có})=0$$

$$T(\text{nước tiểu vàng, đau lưng không})=3$$

$$T(\text{nước tiểu vàng, sốt không})=0$$

$$T(\text{nước tiểu vàng, sốt có})=0$$

$$T(\text{nước tiểu có máu, sốt không})=2$$

$$T(\text{giảm cân ít, đau lưng không})=0$$

$$T(\text{Giảm cân nhanh, đau lưng có})=0$$

$$T(\text{giảm cân không, đau lưng không})=0$$

$$T(\text{giảm cân không, sốt có})=0$$

$$T(\text{giảm cân không, sốt không})=1$$

$$T(\text{giảm cân nhanh, sốt không})=0$$

$$T(\text{giảm cân ít, sốt có})=0$$

T giảm cân ít, sốt không = 0

T đau lưng không, sốt có = 0

T đau lưng không, sốt không = 0

T đau lưng có, sốt không = 0

⇒ Ta có T nước tiểu vàng, đau lưng không = 3 là lớn nhất. ta chọn T nước tiểu vàng, đau lưng không và ta có luật

RULE 4: IF Nước tiểu = Vàng AND Đau lưng = Không THEN Kết Quả = Không bệnh sỏi thận

Kết tiếp ta loại bỏ những dòng ứng với Nước tiểu=vàng và Đau lưng = Không ra khỏi Bảng 2:

Tên	Nước tiểu	Giảm cân	Đau lưng	Sốt	Kết quả
Hương	Có máu	Nhanh	Có	Không	Không bệnh sỏi thận
Tùng	Có máu	Ít	Không	Không	Không bệnh sỏi thận

Với j=2, có 6 tổ hợp gồm 2 thuộc tính:

- {Nước tiểu, Giảm cân}
- {Nước tiểu, Đau lưng}
- {Nước tiểu, Sốt}
- {Giảm cân, Đau lưng}
- {Giảm cân, Sốt}
- {Đau lưng, Sốt}

T nước tiểu có máu, giảm cân ít = 0

T nước tiểu có máu, giảm cân nhanh = 0

T nước tiểu có máu, đau lưng không = 0

T nước tiểu có máu, đau lưng có = 0

T nước tiểu có máu, sốt không = 2

T giảm cân ít, đau lưng không = 0

T giảm cân nhanh, đau lưng có = 0

T giảm cân nhanh, sốt không = 0

T giảm cân ít, sốt không =0

T đau lưng không, sốt không =0

T đau lưng có, sốt không =0

⇒ Ta có T nước tiểu có máu, sốt không =2 là lớn nhất. ta chọn T nước tiểu có máu, sốt không và ta sẽ có luật:

RULE 5: IF Nước tiểu=Có máu AND Sốt = Không THEN Kết Quả=Không bệnh sỏi thận

Kết tiếp, loại bỏ những dòng ứng với Nước tiểu= có máu và sốt = không ra khỏi bảng

⇒ Như vậy ta đã loại bỏ tất cả các dòng trong bảng 2

⇒ Thuật toán kết thúc vì tất cả các bảng đã được xét đến và các dòng trong các bảng đã được loại bỏ.

Tổng hợp các luật:

RULE 1: IF Nước tiểu = có vôi THEN Kết quả = Bệnh sỏi thận

RULE 2: IF Nước tiểu= Có máu AND Sốt = Có THEN Kết Quả = Bệnh sỏi thận

RULE 3: IF Nước tiểu = Vàng AND Đau lưng=Có THEN Kết quả=Bệnh sỏi thận

RULE 4: IF Nước tiểu = Vàng AND Đau lưng = Không THEN Kết Quả = Không bệnh sỏi thận

RULE 5: IF Nước tiểu=Có máu AND Sốt = Không THEN Kết Quả=Không bệnh sỏi thận

3. Cài đặt ứng dụng minh họa

Hai chuyên đề nổi bật là giải thuật để xây dựng cây định danh và tìm ra tri thức cho

mẫu dữ liệu thực tế là giải thuật

Quinlan và giải thuật ILA. Trong phần này ứng dụng

chỉ minh họa cho giải thuật ILA để tìm ra tri thức cho bảng dữ liệu.