



KHAI PHÁ DỮ LIỆU



BUỔI THỰC HÀNH 5

Trần Mạnh Tuấn

Bộ môn Hệ thống thông tin, Khoa CNTT

Trường đại học Thủy Lợi

NỘI DUNG

- Giới thiệu về luật kết hợp
- Luật kết hợp trên weka

Giới thiệu về luật kết hợp

Khai phá luật kết hợp:

- Tìm tần số mẫu, mối kết hợp, sự tương quan, hay các cấu trúc nhân quả giữa các tập đối tượng trong các cơ sở dữ liệu giao tác, cơ sở dữ liệu quan hệ, và những kho thông tin khác.

Tính hiểu được: dễ hiểu

Tính sử dụng được: Cung cấp thông tin thiết thực

Tính hiệu quả: Đã có những thuật toán khai thác hiệu quả

Các ứng dụng:

- Phân tích bán hàng trong siêu thị, cross-marketing, thiết kế catalog, loss-leader analysis, gom cụm, phân lớp, ...

Giới thiệu về luật kết hợp

Các khái niệm

Cho $I = \{I_1, I_2, \dots, I_m\}$ là tập các đơn vị dữ liệu. Cho D là tập các giao tác, mỗi giao tác T là tập các đơn vị dữ liệu sao cho $T \subseteq I$

Định nghĩa 1: Ta gọi giao tác T chứa X , với X là tập các đơn vị dữ liệu của I , nếu $X \subseteq T$

Định nghĩa 2: Một luật kết hợp là một phép suy diễn có dạng $X \rightarrow Y$, trong đó $X \subset I$, $Y \subset I$ và $X \cap Y = \emptyset$

Định nghĩa 3: Ta gọi luật $X \rightarrow Y$ có mức xác nhận(support) là s trong tập giao tác D , nếu có $s\%$ giao tác trong D chứa $X \cup Y$.
Ký hiệu: $\text{Supp}(X \rightarrow Y) = s$

Giới thiệu về luật kết hợp

Định nghĩa 4: Ta gọi luật $X \rightarrow Y$ là có độ tin cậy c (Confidence) trên tập giao tác D ,

$$\text{Ký hiệu: } c = \text{Conf}(X \rightarrow Y) = \text{Supp}(X \rightarrow Y) / \text{Supp}(X)$$

Nhận xét: Các xác nhận và độ tin cậy chính là các xác suất sau:

$$\text{Supp}(X \rightarrow Y) = P(X \cup Y) : \text{Xác suất của } X \cup Y \text{ trong } D$$

$$\text{Conf}(X \rightarrow Y) = P(Y/X) : \text{Xác suất có điều kiện}$$

Định nghĩa 5: Cho trước $\text{Min_Supp} = s_0$ và $\text{Min_Conf} = c_0$

Ta gọi luật $X \rightarrow Y$ là xảy ra nếu thỏa:

$$\text{Supp}(X \rightarrow Y) > s_0 \text{ và } \text{Conf}(X \rightarrow Y) > c_0$$

Giới thiệu về luật kết hợp

- Thuật toán Apriori
- Thuật toán FP-growth

Luật kết hợp trên weka

Luật kết hợp trên weka

- Là một chức năng của Explorer
- Hỗ trợ người dùng huấn luyện và kiểm chứng các thuật toán luật kết hợp cơ bản

Các bước thực hiện luật kết hợp

- Bước 1: tại tab Preprocess, chọn tập dữ liệu và tiền xử lý dữ liệu: các trường dữ liệu dạng Nominal. Nếu ở dạng khác thì dùng bộ lọc để chuyển về: NumericToNominal
- Bước 2: Chọn thuật toán luật kết hợp và tham số
- Bước 3: Tiến hành thực hiện thuật toán
- Bước 4: Ghi nhận và phân tích kết quả

Giao diện chính của luật kết hợp

The screenshot displays the Weka Explorer application window. The 'Associate' tab is selected in the top menu. The 'Associator' section shows the 'Apriori' algorithm with its specific command-line parameters: `-N 50 -T 0 -C 0.7 -D 0.05 -U 1.0 -M 0.4 -S -1.0 -c -1`. Below this, 'Start' and 'Stop' buttons are visible. On the left, a 'Result list (right-click...)' pane shows a single entry: '21:28:00 - Apriori'. The main 'Associator output' pane contains the following text:

```
=== Run information ===  
  
Scheme:      weka.associations.Apriori -N 50 -T 0 -C 0.7 -D 0.05 -U 1.0 -M 0.4 -S -1.0 -c -1  
Relation:    supermarket-weka.filters.unsupervised.attribute.Remove-R1-9,11,15-weka.filters.unsupervised.attribute.F  
Instances:   4627  
Attributes:  107  
              [list of attributes omitted]  
=== Associator model (full training set) ===  
  
Apriori  
=====
```

Minimum support: 0.4 (1851 instances)
Minimum metric <confidence>: 0.7
Number of cycles performed: 12

Generated sets of large itemsets:

Size of set of large itemsets L(1): 17
Size of set of large itemsets L(2): 16

Best rules found:

1. biscuits=t 2605 ==> bread and cake=t 2083 <conf:(0.8)> lift:(1.11) lev:(0.04) [208] conv:(1.4)
2. milk-cream=t 2939 ==> bread and cake=t 2337 <conf:(0.8)> lift:(1.1) lev:(0.05) [221] conv:(1.37)
3. fruit=t 2962 ==> bread and cake=t 2325 <conf:(0.78)> lift:(1.09) lev:(0.04) [193] conv:(1.3)
4. baking needs=t 2795 ==> bread and cake=t 2191 <conf:(0.78)> lift:(1.09) lev:(0.04) [179] conv:(1.29)

Tổng hợp so sánh luật kết hợp

- **Chạy 1 bộ dữ liệu với các phương pháp thuật toán khác nhau**
- **Chạy thuật toán Apriori với các bộ dữ liệu khác nhau**



THỰC HÀNH