

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

NGUYỄN HỒNG SÂM

KHAI PHÁ TẬP MỤC LỢI ÍCH CAO
SỬ DỤNG CẤU TRÚC CÂY TIỀN TỔ

Chuyên ngành: Khoa học máy tính

Mã số: 60.48.01

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

NGƯỜI HƯỚNG DẪN KHOA HỌC

TS. NGUYỄN HUY ĐỨC

Thái Nguyên - 2014

LỜI CẢM ƠN

Lời đầu tiên tôi xin gửi lời cảm ơn chân thành và biết ơn sâu sắc tới TS. Nguyễn Huy Đức – Trường Cao đẳng Sư phạm Trung ương, người đã chỉ bảo và hướng dẫn tận tình cho tôi trong suốt quá trình nghiên cứu khoa học và thực hiện luận văn này.

Tôi xin chân thành cảm ơn sự dạy bảo, giúp đỡ, tạo điều kiện và khuyến khích tôi trong quá trình học tập và nghiên cứu của các thầy cô giáo của Viện Công nghệ Thông tin, Trường Đại học Công nghệ Thông tin và Truyền thông – Đại học Thái Nguyên.

Và cuối cùng, tôi xin gửi lời cảm ơn tới gia đình, người thân và bạn bè – những người luôn ở bên tôi những lúc khó khăn nhất, luôn động viên tôi, khuyến khích tôi trong cuộc sống và trong công việc.

Tôi xin chân thành cảm ơn!

Thái Nguyên, ngày 12 tháng 03 năm 2014

Tác giả

Nguyễn Hồng Sâm

LỜI CAM ĐOAN

Tôi xin cam đoan Luận văn *“Khai phá tập mục lợi ích cao sử dụng cấu trúc cây tiền tố”* được thực hiện theo đúng mục tiêu đề ra dưới sự hướng dẫn của TS. Nguyễn Huy Đức. Trong toàn bộ luận văn, những điều được trình bày là của cá nhân hoặc là được tổng hợp từ nhiều nguồn tài liệu. Tất cả các loại tài liệu đều có xuất xứ rõ ràng và được trích dẫn hợp pháp. Tôi xin chịu trách nhiệm về luận văn của mình.

Thái Nguyên, ngày 12 tháng 03 năm 2014

Tác giả

Nguyễn Hồng Sâm

MỤC LỤC

Trang phụ bì	Trang
LỜI CẢM ƠN	i
LỜI CAM ĐOAN	ii
MỤC LỤC	iii
DANH MỤC CÁC KÝ HIỆU VÀ CÁC CHỮ VIẾT TẮT	v
DANH MỤC CÁC BẢNG	vi
DANH MỤC HÌNH VẼ	vii
LỜI MỞ ĐẦU	1
CHƯƠNG 1: TỔNG QUAN VỀ KHAI PHÁ DỮ LIỆU VÀ KHAI PHÁ TẬP MỤC THƯỜNG XUYÊN	3
1.1. Khái niệm về khai phá tri thức và khai phá dữ liệu	3
1.2. Kiến trúc của hệ thống khai phá dữ liệu	4
1.3. Quá trình khai phá dữ liệu	5
1.4. Một số kỹ thuật khai phá dữ liệu	6
1.5. Các cơ sở dữ liệu phục vụ cho khai phá dữ liệu	9
1.6. Một số ứng dụng của khai phá dữ liệu	11
1.7. Khai phá tập mục thường xuyên	12
1.7.1. Các khái niệm cơ bản	12
1.7.1.1. Cơ sở dữ liệu giao tác	12
1.7.1.2. Tập mục thường xuyên và luật kết hợp	14
1.7.1.3. Bài toán khai phá luật kết hợp	15
1.7.2. Cách tiếp cận khai phá tập mục thường xuyên	16
1.7.3. Một số thuật toán điển hình tìm tập mục thường xuyên	17
1.7.3.1. Thuật toán Apriori	17
1.7.3.2. Thuật toán COFI-tree	21
1.7.4. Mở rộng bài toán khai phá tập mục thường xuyên	26
1.8. Kết luận chương 1	27
CHƯƠNG 2: KHAI PHÁ TẬP MỤC LỢI ÍCH CAO SỬ DỤNG CẤU TRÚC CÂY TIỀN TỔ	28
2.1. Mở đầu	28

2.2. Bài toán khai phá tập mục lợi ích cao.....	29
2.3. Khai phá tập mục lợi ích cao sử dụng cấu trúc cây tiền tố.....	33
2.3.1. Thuật toán COUI-Mine	33
2.3.1.1. Xây dựng cây TWUI-tree.....	34
2.3.1.2. Khai phá cây TWUI-tree	39
2.3.1.3. Nhận xét và đánh giá thuật toán COUI-Mine	45
2.3.1.4. Khai phá tương tác với cây TWUI-tree.....	47
2.3.2. Các cấu trúc cây tiền tố cho khai phá lợi ích cao	48
2.3.3. Thuật toán UP-Growth	51
2.4. Kết luận chương 2.....	54
CHƯƠNG 3: THỰC NGHIỆM KHAI PHÁ TẬP MỤC LỢI ÍCH CAO	56
3.1. Bài toán phát hiện nhóm các mặt hàng có lợi nhuận cao	56
3.2. Mô tả dữ liệu.....	56
3.3. Xây dựng chương trình.....	60
3.4. Thực nghiệm khai phá tìm tập mục lợi ích cao	60
3.5. Kết quả thực nghiệm.....	61
KẾT LUẬN	62
TÀI LIỆU THAM KHẢO	63
Tiếng Việt	63
Tiếng Anh	63
PHỤ LỤC.....	65

DANH MỤC CÁC KÝ HIỆU VÀ CÁC CHỮ VIẾT TẮT

Trong luận văn này, dùng thống nhất các ký hiệu và chữ viết tắt sau:

Các ký hiệu:

$I = \{i_1, i_2, \dots, i_n\}$: Tập n mục dữ liệu.

$DB = \{T_1, T_2, \dots, T_m\}$: Cơ sở dữ liệu có m giao tác.

Db : cơ sở dữ liệu giao tác con của DB , $db \subseteq DB$.

I_p : Mục dữ liệu thứ p .

T_q : Giao tác thứ q .

n : Số mục dữ liệu một cơ sở dữ liệu giao tác.

m : Số giao tác một cơ sở dữ liệu giao tác.

A, B, C, \dots : Tên các mục dữ liệu trong cơ sở dữ liệu giao tác ví dụ.

X, Y, \dots : Tập con của tập mục dữ liệu I , $X, Y \subseteq I$.

$X = ABC$ thay cho $X = \{A, B, C\}$ trong các cơ sở dữ liệu giao tác ví dụ.

Nếu $X \subseteq Y$ thì X gọi là tập con của tập Y , Y gọi là tập cha của tập X .

minsup: Ngưỡng độ hỗ trợ tối thiểu.

minutil: Giá trị lợi ích tối thiểu.

$|X|$: Số phần tử của tập hợp X .

Viết tắt:

KPDL: Khai phá dữ liệu

CSDL: Cơ sở dữ liệu.

CNTT: Công nghệ thông tin.

CNTT và TT: Công nghệ Thông tin và Truyền thông.

DANH MỤC CÁC BẢNG

Bảng 1.1: Biểu diễn ngang của cơ sở dữ liệu giao tác.	13
Bảng 1.2: Biểu diễn dọc của cơ sở dữ liệu giao tác.	13
Bảng 1.3: Ma trận giao tác của cơ sở dữ liệu cho ở bảng 1.1.	14
Bảng 1.4: Cơ sở dữ liệu giao tác minh họa thực hiện thuật toán Apriori.	20
Bảng 1.5: Cơ sở dữ liệu giao tác minh họa thực hiện thuật toán COFI-tree.	22
Bảng 1.6: Các mục dữ liệu và độ hỗ trợ.	23
Bảng 1.7: Các mục dữ liệu thường xuyên đã sắp thứ tự.	23
Bảng 1.8: Các mục dữ liệu trong giao tác sắp xếp giảm dần theo độ hỗ trợ.	23
Bảng 2.1: Cơ sở dữ liệu giao tác.	30
Bảng 2.2: Bảng lợi ích.	30
Bảng 2.3: Lợi ích các giao tác của cơ sở dữ liệu bảng 2.1 và bảng 2.2.	36
Bảng 2.4: Lợi ích TWU của các mục dữ liệu.	36
Bảng 2.5: Các mục dữ liệu có lợi ích TWU cao sắp giảm dần theo twu.	36
Bảng 2.6: Các mục dữ liệu trong giao tác sắp giảm dần theo lợi ích TWU.	37
Bảng 2.7: Lợi ích các tập mục ứng viên.	43
Bảng 2.8: Cơ sở dữ liệu ví dụ cho thuật toán UP-Growth.	52
Bảng 2.9: Bảng lợi ích của CSDL bảng 2.8.	53
Bảng 2.10: Các giao tác được sắp lại các mục dữ liệu theo TWU giảm dần.	53
Bảng 3.1: Dữ liệu đã trích chọn để khai phá.	57
Bảng 3.2: Bảng lợi ích các mặt hàng.	58
Bảng 3.3: Mã hóa các mặt hàng.	58

DANH MỤC HÌNH VẼ

Hình 1.1: Quá trình phát hiện tri thức	4
Hình 1.2: Kiến trúc của một hệ thống khai phá dữ liệu	5
Hình 1.3: Quá trình KPDL	6
Hình 1.4: Cây quyết định	7
Hình 1.5: Mẫu kết quả của nhiệm vụ phân cụm dữ liệu	8
Hình 1.6: Mẫu kết quả của nhiệm vụ hồi quy	8
Hình 1.7: Phân loại các thuật toán khai phá tập mục thường xuyên.....	17
Hình 1.8: Cây FP-tree của CSDL bảng 1.5.....	24
Hình 1.9: Cây COFI-tree của mục D.	24
Hình 1.10: Các bước khai phá cây D-COFI-tree.	25
Hình 2.1: Cây TWUI-tree sau khi lưu giao tác $T1$	37
Hình 2.2: Cây TWUI-tree sau khi lưu giao tác $T1$ và $T2$	38
Hình 2.3: Cây TWUI-tree của cơ sở dữ liệu bảng 2.1 và bảng 2.2.....	38
Hình 2.4: Cây C-COUI-tree sau khi lưu mẫu CBE.....	40
Hình 2.5: Cây C-COUI-tree sau khi lưu mẫu CBE và CE.....	40
Hình 2.6: Cây C-COUI-tree sau khi xây dựng xong.....	40
Hình 2.7: Cây D-COUI-tree.....	41
Hình 2.8: Cây B-COUI-tree.	41
Hình 2.9: Các bước khai phá cây D-COUI-Tree.	42
Hình 2.10: Cây TWUI-tree có các mục dữ liệu sắp tăng dần theo trật tự từ điển của cơ sở dữ liệu bảng 2.1 và bảng 2.2.....	49
Hình 2.11: Cây TWUI-tree có các mục dữ liệu sắp giảm dần theo số lần xuất hiện của chúng trong cơ sở dữ liệu bảng 2.1 và bảng 2.2.	49
Hình 2.12: Cây TWUI-tree có các mục dữ liệu sắp giảm dần theo TWU của chúng trong cơ sở dữ liệu bảng 2.1 và bảng 2.2.	50
Hình 2.13: Cây TWUI-tree của CSDL bảng 2.8 với $minutil = 40$	54
Hình 2.14: Cây UP-tree của CSDL bảng 2.8 với $minutil = 40$	54
Hình 3.1: Dữ liệu đã mã hóa chuẩn bị cho khai phá.....	59
Hình 3.2: Dữ liệu mã hóa của bảng 3.2.....	59
Hình 3.3: Giao diện chương trình	60
Hình 3.4: Giao diện kết quả khai phá.....	61

LỜI MỞ ĐẦU

Khai phá tập mục thường xuyên đóng vai trò quan trọng trong nhiều nhiệm vụ khai phá dữ liệu. Khai phá tập mục thường xuyên xuất hiện như là bài toán con của nhiều lĩnh vực khai phá dữ liệu như khám phá luật kết hợp, khám phá mẫu tuần tự,... Bài toán khai phá luật kết hợp do Agrawal, T.Imielinski và A. N. Swami đề xuất và nghiên cứu lần đầu vào năm 1993 với mục tiêu là phát hiện các tập mục thường xuyên, từ đó tạo các luật kết hợp. Trong mô hình của bài toán này, giá trị của mỗi mục dữ liệu trong một giao tác là 0 hoặc 1, tức là chỉ quan tâm mục dữ liệu có xuất hiện trong giao tác hay không. Bài toán cơ bản này có nhiều ứng dụng, tuy vậy, do tập mục thường xuyên chỉ mang ngữ nghĩa thông kê nên nó chỉ đáp ứng được phần nào nhu cầu của thực tiễn.

Nhằm khắc phục hạn chế của bài toán cơ bản khai phá luật kết hợp, nhiều nhà nghiên cứu đã mở rộng bài toán theo nhiều hướng khác nhau. Năm 1997, Hilderman và các cộng sự đề xuất bài toán khai phá tập mục cổ phần cao. Trong mô hình này, giá trị của mục dữ liệu trong giao tác là một số. Năm 2004, nhóm các nhà nghiên cứu H. Yao, Hamilton và Butz, mở rộng tiếp bài toán, đề xuất mô hình khai phá tập mục lợi ích cao.

Trong mô hình khai phá tập mục lợi ích cao, giá trị của mục dữ liệu trong giao tác là một số (như số lượng đã bán của mặt hàng, gọi là giá trị khách quan), ngoài ra còn có bảng lợi ích cho biết lợi ích mang lại khi bán một đơn vị hàng đó (gọi là giá trị chủ quan). Lợi ích của tập mục là số đo lợi nhuận mà tập mục đó mang lại. Khai phá tập mục lợi ích cao là khám phá tất cả các tập mục có lợi ích không nhỏ hơn ngưỡng lợi ích tối thiểu của người sử dụng.

Trong những năm gần đây, bài toán này đã và đang thu hút sự quan tâm của nhiều nhà nghiên cứu trong và ngoài nước. Đã có một số thuật toán phát hiện tập mục lợi ích cao được đề xuất. Các thuật toán này có thể phân thành hai loại:

- Thuật toán kiểu Apriori (Apriori-like), sinh ra các tập mục ứng viên, duyệt theo chiều rộng.

- Thuật toán không sinh ứng viên, chuyển đổi cơ sở dữ liệu thành cấu trúc cây, duyệt theo chiều sâu để phát hiện các tập mục lợi ích cao. Các thuật toán loại này hiệu quả hơn và tránh được khối lượng tính toán lớn.

Với ý nghĩa đó và mục đích tìm hiểu bài toán tìm tập mục lợi ích cao và các thuật toán khai phá sử dụng cấu trúc cây tiền tố, em đã quyết định lựa chọn đề tài luận văn: “ **KHAI PHÁ TẬP MỤC LỢI ÍCH CAO SỬ DỤNG CẤU TRÚC CÂY TIỀN TỐ**”

Nội dung luận văn gồm 3 chương:

Chương 1: Tổng quan về khai phá dữ liệu và khai phá tập mục thường xuyên.

Chương 2: Khai phá tập mục lợi ích cao sử dụng cấu trúc cây tiền tố.

Chương 3: Chương trình thực nghiệm và ứng dụng.