

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CNTT&TT



ĐỖ THỊ HẢI YẾN

**KHAI PHÁ TẬP MỤC LỢI ÍCH CAO
TRONG CƠ SỞ DỮ LIỆU LỚN**

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

THÁI NGUYÊN - 2011

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CNTT&TT



ĐỖ THỊ HẢI YẾN

**KHAI PHÁ TẬP MỤC LỢI ÍCH CAO
TRONG CƠ SỞ DỮ LIỆU LỚN**

CHUYÊN NGÀNH: KHOA HỌC MÁY TÍNH

MÃ SỐ: 60 48 01

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

HƯỚNG DẪN KHOA HỌC: PGS.TS. NGUYỄN THANH TÙNG

THÁI NGUYÊN - 2011

LỜI CAM ĐOAN

Tôi xin cam đoan Luận văn "Khai phá tập mục lợi ích cao trong cơ sở dữ liệu lớn" là công trình nghiên cứu của riêng tôi dưới sự hướng dẫn của PGS.TS Nguyễn Thanh Tùng. Kết quả đạt được trong luận văn là sản phẩm của riêng cá nhân tôi, không sao chép lại của người khác. 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 tài liệu tham khảo đều có xuất xứ rõ ràng và được trích dẫn hợp pháp.

Tôi xin chịu hoàn toàn trách nhiệm và chịu mọi hình thức kỷ luật theo quy định cho lời cam đoan của mình.

Thái Nguyên, ngày 30 tháng 9 năm 2011

Người cam đoan

Đỗ Thị Hải Yến

LỜI CẢM ƠN

Lời đầu tiên tôi xin gửi lời cảm ơn chân thành và biết ơn sâu sắc tới PGS.TS. Nguyễn Thanh Tùng - Viện Công nghệ thông tin, người thầy 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 30 tháng 9 năm 2011

Tác giả

Đỗ Thị Hải Yến

MỤC LỤC

Trang

Trang bìa phụ	
Lời cảm ơn	
Lời cam đoan	
Mục lục.....	i
Danh mục các từ, các ký hiệu viết tắt	iii
Danh mục các bảng	iv
Danh mục các hình	v
LỜI MỞ ĐẦU	1
Chương 1. KHÁI QUÁT VỀ KHAI PHÁ DỮ LIỆU VÀ BÀI TOÁN KHAI PHÁ TẬP MỤC THƯỜNG XUYÊN	4
1.1. Khai phá dữ liệu	4
1.2. Khai phá tập mục thường xuyên.....	8
1.2.1. Cơ sở dữ liệu giao tác	8
1.2.2. Tập mục thường xuyên và luật kết hợp	10
1.2.3 Bài toán khai phá luật kết hợp	11
1.3. Các cách tiếp cận khai phá tập mục thường xuyên	12
1.3.1 Thuật toán Apriori	13
1.3.2 Thuật toán FP-growth	17
1.4. Mở rộng bài toán khai phá tập mục thường xuyên.....	23
1.5. Kết luận chương 1	24
Chương 2. KHAI PHÁ TẬP MỤC LỢI ÍCH CAO: BÀI TOÁN VÀ BA THUẬT GIẢI KIỂU APRIORI	25
2.1. Mở đầu.....	25
2.2. Bài toán khai phá tập mục lợi ích cao	26
2.3. Ba thuật toán khai phá tập mục lợi ích cao kiểu Apriori.....	30
2.3.1. Thuật toán UMining	30
2.3.2. Thuật toán UMining-H	32
2.3.3. Thuật toán hai pha HUMining	34

2.4. Kết luận chương 2	41
Chương 3. THUẬT TOÁN HIỆU QUẢ KHAI PHÁ TẬP MỤC LỢI ÍCH	
CAO KIỂU FP-GROWTH.....	42
3.1 Mở đầu.....	42
3.2. Thuật toán COUI-Mine	42
3.2.1. Xây dựng cây TWUI-tree	44
3.2.2. Khai phá cây TWUI-tree	48
3.2.3. Đánh giá độ phức tạp của thuật toán COUI-Mine.....	55
3.2.4. Nhận xét thuật toán COUI-Mine	58
3.2.5. Khai phá tương tác với cây TWU-tree	59
3.3. Kết luận chương 3	60
KẾT LUẬN	62
TÀI LIỆU THAM KHẢO	64

DANH MỤC CÁC TỪ VIẾT TẮT

STT	Cụm từ viết tắt	Nghĩa của cụm từ viết tắt
1	CNTT	Công nghệ thông tin
2	CSDL	Cơ sở dữ liệu
3	KDD	Khám phá tri thức trong cơ sở dữ liệu

DANH MỤC CÁC BẢNG

	<i>Trang</i>
Bảng 1.1: Biểu diễn ngang của cơ sở dữ liệu giao tác	9
Bảng 1.2: Biểu diễn dọc của cơ sở dữ liệu giao tác	9
Bảng 1.3: Ma trận giao tác của cơ sở dữ liệu bảng 1.1	10
Bảng 1.4: Cơ sở dữ liệu giao tác minh họa thực hiện thuật toán Apriori.	16
Bảng 1.5: Cơ sở dữ liệu giao tác minh họa thực hiện thuật toán COFI-tree	19
Bảng 1.6: Các mục dữ liệu và độ hỗ trợ.....	20
Bảng 1.7: Các mục dữ liệu thường xuyên đã sắp thứ tự	20
Bảng 1.8: Các mục dữ liệu trong giao tác sắp giảm dần theo độ hỗ trợ.	21
Hình 1.4: Các bước khai phá cây D-COFI-tree.	23
Bảng 2.1. Cơ sở dữ liệu giao tác	27
Bảng 2.2. Giá trị lợi ích chủ quan của các mục trong bảng 1.	27
Hình 2.1. Dàn tập mục trong cơ sở dữ liệu bảng 1	29
Bảng 3.1: Lợi ích các giao tác của cơ sở dữ liệu	45
Bảng 3.2: Lợi ích TWU của các mục dữ liệu.....	45
Bảng 3.3: Các mục dữ liệu có lợi ích TWU cao sắp giảm dần theo twu	46
Bảng 3.4: Các mục dữ liệu trong giao tác sắp giảm dần theo lợi ích TWU.	46
Hình 3.7: Cây D-COUI-tree.....	50
Bảng 3.5: Lợi ích các tập mục ứng viên	53

DANH MỤC CÁC HÌNH

	<i>Trang</i>
Hình 1.1. Các bước thực hiện của quá trình khai phá dữ liệu.....	6
Hình 1.2: Cây FP-tree của CSDL bảng 1.5.....	21
Hình 1.3: Cây COFI-tree của mục D	21
Hình 2.2. Không gian tìm kiếm tập mục lợi ích cao theo thuật toán UMining	32
Hình 2.3. Không gian tìm kiếm tập mục lợi ích cao theo thuật toán UMining-H	33
Hình 2.4. Không gian tìm kiếm tập mục lợi ích cao theo thuật toán HUMining	39
Hình 3.1: Cây TWUI-tree sau khi lưu thao tác T1.....	47
Hình 3.2: Cây TWUI-tree sau khi lưu thao tác T1 và T2.	47
Hình 3.3: Cây TWUI-tree của cơ sở dữ liệu bảng 3.1 và 3.2.	47
Hình 3.4: Cây C-COUI-tree sau khi lưu mẫu CBE.....	49
Hình 3.5: Cây C-COUI-tree sau khi lưu mẫu CBE và CE.....	50
Hình 3.6: Cây C-COUI-tree sau khi xây dựng xong.....	50
Hình 3.8: Cây B-COUI-tree	51
Hình 3.9: Các bước khai phá cây D-COUI-tree.....	52

LỜI MỞ ĐẦU

Trong những năm gần đây, cùng với sự phát triển vượt bậc của công nghệ thông tin, truyền thông, khả năng thu thập và lưu trữ thông tin của các hệ thống thông tin không ngừng được nâng cao. Với lượng dữ liệu khổng lồ và luôn gia tăng theo thời gian, rõ ràng các phương pháp phân tích dữ liệu truyền thống sẽ không còn hiệu quả, gây tốn kém và dễ dẫn đến những kết quả sai lệch. Để có thể khai thác hiệu quả các cơ sở dữ liệu lớn, một lĩnh vực khoa học mới đã ra đời: Khám phá tri thức trong cơ sở dữ liệu (Knowledge Discovery in Databases – KDD). Khai phá dữ liệu (Data Mining) là một công đoạn chính trong quá trình khám phá tri thức, nhằm *tìm kiếm, phát hiện các tri thức mới, hữu ích tiềm ẩn trong các cơ sở dữ liệu lớn.*

Khai phá luật kết hợp là một nhiệm vụ quan trọng của khai phá dữ liệu. Bài toán truyền thống (hay còn gọi bài toán nhị phân) khai phá luật kết hợp do R. Agrawal, T. Imielinski và A. N. Swami đề xuất và nghiên cứu lần đầu tiên vào năm 1993 khi phân tích các cơ sở dữ liệu của các siêu thị. Mục tiêu của nó là phát hiện các *tập mục thường xuyên*, từ đó tạo các *luật kết hợp* phản ánh hành vi mua hàng của khách hàng. Những thông tin như vậy giúp nhà quản lý có thể lựa chọn phương án tiếp thị, kinh doanh hiệu quả hơn.

Cho đến nay, bài toán khai phá luật kết hợp truyền thống có nhiều ứng dụng, tuy vậy do tập mục thường xuyên chỉ mang ngữ nghĩa thống kê nên mô hình bài toán truyền thống chỉ đáp ứng được phần nào nhu cầu ứng dụng thực tiễn. Thật ra, trong kinh doanh, điều mà người quản lý quan tâm hơn là phát hiện những khách VIP, đem lại lợi nhuận cao. Trong thực hành, có những tập mục thường xuyên nhưng chỉ đóng góp phần rất nhỏ, ngược lại có những tập mục không thường xuyên lại đóng góp phần đáng kể vào lợi nhuận chung của công ty.

Gần đây, nhằm khắc phục hạn chế của bài toán truyền thống khai phá luật kết hợp, các nhà nghiên cứu đã mở rộng nó theo nhiều hướng khác nhau, trong đó có vấn đề khai phá tập mục lợi ích cao. Lợi ích của một tập mục là số đo lợi nhuận mà nó có thể mang lại trong kinh doanh, được tính toán dựa trên giá trị khách quan và