

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

TRẦN THỊ THU TRANG

KHAI PHÁ LUẬT KẾT HỢP
TỪ DỮ LIỆU CHUỖI THỜI GIAN

LUẬN VĂN THẠC SĨ: KHOA HỌC MÁY TÍNH

Thái Nguyên - 2012

LỜI CAM ĐOAN

Tôi xin cam đoan luận văn “Khai phá luật kết hợp từ dữ liệu chuỗi thời gian” là công trình nghiên cứu của riêng tôi dưới sự hướng dẫn của PGS.TS. Bùi Thế Hồng. Toàn bộ phần mềm do chính tôi lập trình và kiểm thử. Tôi xin chịu trách nhiệm về lời cam đoan của mình.

Các số liệu và thông tin sử dụng trong luận văn này hoàn toàn là trung thực.

Tác giả

Trần Thị Thu Trang

MỤC LỤC

MỤC LỤC.....	ii
DANH MỤC HÌNH VẼ.....	iv
DANH MỤC CÁC BẢNG.....	v
DANH MỤC CÁC KÝ HIỆU, CÁC TỪ VIẾT TẮT	vi
MỞ ĐẦU	1
CHƯƠNG 1: KHÁI QUÁT VỀ KHAI PHÁ DỮ LIỆU VÀ CHUỖI THỜI GIAN	3
1.1. Khai phá dữ liệu	3
1.1.1. Khai phá dữ liệu là gì?.....	3
1.1.2. Nhiệm vụ của khai phá dữ liệu	3
1.1.3. Triển khai việc khai phá dữ liệu	5
1.1.4. Một số ứng dụng khai phá dữ liệu	6
1.1.5. Quá trình phát hiện tri thức trong cơ sở dữ liệu	7
1.1.6. Các kỹ thuật khai phá dữ liệu	99
1.2. Dữ liệu chuỗi thời gian.....	14
1.2.1. Khái niệm.....	14
1.2.2. Tiền xử lý dữ liệu chuỗi thời gian	17
CHƯƠNG 2: KHAI PHÁ LUẬT KẾT HỢP TỪ DỮ LIỆU CHUỖI THỜI GIAN	20
2.1. Luật kết hợp trong khai phá dữ liệu	20
2.1.1. Khái niệm luật kết hợp	20
2.1.2. Lý thuyết về luật kết hợp	21
2.2. Khai phá luật kết hợp	27

2.2.1. Khai phá luật kết hợp từ cơ sở dữ liệu.....	27
2.2.2. Khai phá luật kết hợp từ dữ liệu chuỗi thời gian.....	28
2.3. Thuật toán khai phá luật kết hợp từ dữ liệu chuỗi thời gian	30
2.3.1. Thuật toán khai phá luật kết hợp từ dữ liệu thường	30
2.3.2. Thuật toán khai phá luật kết hợp từ dữ liệu chuỗi thời gian	40
CHƯƠNG 3: XÂY DỰNG CHƯƠNG TRÌNH THỬ NGHIỆM.....	53
3.1. Phát biểu bài toán	53
3.2. Xây dựng chương trình.....	54
KẾT LUẬN	63
TÀI LIỆU THAM KHẢO.....	64

DANH MỤC HÌNH VẼ

Hình 1.1. Quá trình phát hiện tri thức trong cơ sở dữ liệu.....	8
Hình 1.2. Đồ thị thể hiện thành phần xu hướng dài hạn	15
Hình 1.3. Đồ thị thể hiện thành phần mùa	16
Hình 1.4. Đồ thị thể hiện thành phần chu kỳ	16
Hình 1.5. Trung bình trượt hàm mũ	17
Hình 2.1. Một cây mẫu thường xuyên.....	39
Hình 2.2. FP-Tree và CFP-Tree	42
Hình 2.3: Các khoản mục được ánh xạ.....	44
Hình 2.4: Ví dụ cây CFP-Tree	45
Hình 3.1. Bảng cơ sở dữ liệu.....	55
Hình 3.2. Giao diện chính của chương trình.....	56
Hình 3.3. Thực hiện chọn CSDL	56
Hình 3.4. Thực hiện xóa CSDL	57
Hình 3.5. Tìm tập phổ biến dựa trên thuật toán CFPmine.....	58
Hình 3.6. Thực hiện lệnh Reset.....	59
Hình 3.7. Chọn dữ liệu cho thuật toán tìm luật kết hợp.....	60
Hình 3.8. Thực hiện xóa cơ sở dữ liệu	60
Hình 3.9. Thực hiện luật kết hợp	61
Hình 3.10. Thực hiện lệnh Reset.....	62

DANH MỤC CÁC BẢNG

Bảng 2.1. Ma trận biểu diễn cơ sở dữ liệu	35
Bảng 2.2. Vector biểu diễn nhị phân cho tập 1 thuộc tính	35
Bảng 2.3. Vector biểu diễn nhị phân cho các tập 2 thuộc tính	36
Bảng 2.4. Vector biểu diễn nhị phân cho các tập 3 thuộc tính	36
Bảng 2.5. Vector biểu diễn nhị phân cho các tập 4 thuộc tính	36
Bảng 2.6. Các giao tác cơ sở dữ liệu	38
Bảng 2.7. Khoản mục và số lần xuất hiện trong cơ sở dữ liệu	40

DANH MỤC CÁC KÝ HIỆU, CÁC TỪ VIẾT TẮT

Các từ viết tắt	Nghĩa tiếng anh	Nghĩa tiếng việt
FI	Frequent Itemset	Tập mục thường xuyên
FCI	Frequent Closed Itemset	Tập mục thường xuyên đóng
MFI	Maximally Frequent Itemset	Tập mục thường xuyên lớn nhất
CSDL		Cơ sở dữ liệu
FP-Tree	Frequent Pattern Tree	Cây mẫu thường xuyên
CFP-Tree	Compressed Frequent Pattern Tree	Cây mẫu thường xuyên nén
ITARM	Inter-Transaction Association Rules Mining	Khai phá luật kết hợp liên giao dịch

MỞ ĐẦU

Ngày nay, cuộc cách mạng của kỹ thuật số cho phép số hóa thông tin dễ dàng và chi phí lưu trữ thấp. Với sự phát triển của phần mềm, phần cứng và trang bị nhanh hệ thống máy tính trong kinh doanh. Số lượng dữ liệu khổng lồ được tập trung và lưu trữ trong cơ sở dữ liệu. Dữ liệu sau khi phục vụ cho một mục đích nào đó được lưu lại trong kho dữ liệu và theo ngày tháng khối lượng dữ liệu được lưu trữ ngày càng lớn. Trong khối lượng dữ liệu to lớn này có rất nhiều thông tin có ích mang tính tổng quát, thông tin có tính quy luật vẫn còn đang tiềm ẩn mà chúng ta chưa biết. Từ khối lượng dữ liệu rất lớn cần có những công cụ tự động rút các thông tin và kiến thức có ích. Một hướng tiếp cận có khả năng giúp các công ty khai thác các thông tin có nhiều ý nghĩa từ các tập dữ liệu lớn đó là khai phá dữ liệu.

Với sự bùng nổ và phát triển của công nghệ thông tin đã mang lại nhiều hiệu quả đối với khoa học cũng như các hoạt động thực tế, trong đó khai phá dữ liệu là một trong những lĩnh vực mang lại hiệu quả thiết thực cho con người. Khai phá dữ liệu đã giúp người sử dụng thu được những tri thức hữu ích từ những cơ sở dữ liệu hoặc các kho dữ liệu khổng lồ khác. Luận văn đề cập đến các khái niệm và vấn đề cơ bản trong khai phá luật kết hợp từ dữ liệu chuỗi thời gian được áp dụng trong cơ sở dữ liệu bán hàng.

Luận văn cấu trúc gồm 3 chương:

Chương 1:

Trong chương 1 tìm hiểu khái quát về khai phá dữ liệu và dữ liệu chuỗi thời gian và phương pháp tiền xử lý dữ liệu chuỗi thời gian.

Chương 2:

Trong chương 2 sẽ tìm hiểu phương pháp khai phá dữ liệu từ chuỗi thời gian qua thuật toán ITARM dựa trên cấu trúc cây CFPTree.

Chương 3:

Trong chương 3 tiến hành cài đặt thuật toán ở chương 2 và cài đặt ứng dụng của thuật toán trên cơ sở dữ liệu bán hàng.

Luận văn này được hoàn thành dưới sự hướng dẫn tận tình của PGS.TS **Bùi Thế Hồng**, em xin bày tỏ lòng biết ơn chân thành của mình đối với thầy. Em xin chân thành cảm ơn các thầy, cô giáo Viện Công nghệ thông tin, Trường Đại học Công nghệ thông tin và Truyền thông - Đại học Thái Nguyên đã tham gia giảng dạy, giúp đỡ em trong suốt quá trình học tập nâng cao trình độ kiến thức. Tuy nhiên vì điều kiện thời gian và khả năng có hạn nên luận văn không thể tránh khỏi những thiếu sót. Em kính mong các thầy cô giáo và các bạn đóng góp ý kiến để đề tài được hoàn thiện hơn.

CHƯƠNG 1: KHÁI QUÁT VỀ KHAI PHÁ DỮ LIỆU VÀ CHUỖI THỜI GIAN

1.1. Khai phá dữ liệu

1.1.1. Khai phá dữ liệu là gì?

Khai phá dữ liệu là một khái niệm ra đời vào những năm cuối của thập kỷ 80. Nó bao hàm một loạt các kỹ thuật nhằm phát hiện ra các thông tin có giá trị tiềm ẩn trong các tập dữ liệu lớn (các kho dữ liệu). Về bản chất, khai phá dữ liệu liên quan đến việc phân tích các dữ liệu và sử dụng các kỹ thuật để tìm ra các mẫu hình có tính chính quy trong tập dữ liệu.

Năm 1989, Fayyad, Piatetsky-Shapiro và Smyth đã dùng khái niệm *Phát hiện tri thức trong cơ sở dữ liệu* để chỉ toàn bộ quá trình phát hiện các tri thức có ích từ các tập dữ liệu lớn. Trong đó, *khai phá dữ liệu* là một bước đặc biệt trong toàn bộ quá trình, sử dụng các giải thuật đặc biệt để chiết xuất ra các mẫu (hay các mô hình) từ dữ liệu.

Ở một mức độ trừu tượng nhất định có thể định nghĩa về khai phá dữ liệu: *Khai phá dữ liệu* là một quá trình tìm kiếm, phát hiện các tri thức mới, tiềm ẩn, hữu dụng trong cơ sở dữ liệu lớn.

Khám phá tri thức là mục tiêu chính của khai phá dữ liệu, do vậy hai khái niệm đó được xem như hai lĩnh vực tương đương nhau. Nhưng, nếu phân chia một cách tách bạch thì khai phá dữ liệu là một bước chính trong quá trình khám phá tri thức.

1.1.2. Nhiệm vụ của khai phá dữ liệu

Các bài toán liên quan đến khai phá dữ liệu về bản chất là các bài toán thống kê. Điểm khác biệt giữa các kỹ thuật khai phá dữ liệu và các công cụ phục vụ tính toán thống kê mà chúng ta đã biết là ở khối lượng cần tính toán. Một khi dữ liệu đã trở nên khổng lồ thì những khâu như: thu thập dữ liệu, tiền xử lý và xử lý dữ liệu đều đòi hỏi phải được tự động hóa. Tuy