

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

**NGUYỄN TÚ NAM**

**KHAI PHÁ TẬP MỤC THƯỜNG XUYÊN CÓ  
TRỌNG SỐ TRÊN CƠ SỞ DỮ LIỆU GIAO TÁC**

**LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH**

THÁI NGUYÊN - 2015

## **LỜI CAM ĐOAN**

Tôi xin cam đoan đây là công trình nghiên cứu của tôi dưới sự hướng dẫn của TS Nguyễn Long Giang.

Các số liệu, kết quả nghiên cứu trong luận văn là trung thực và mọi trích dẫn trong báo cáo đều được ghi rõ nguồn gốc. Nếu có sử dụng bất hợp pháp kết quả công trình nghiên cứu của người khác trong báo cáo tôi xin hoàn toàn chịu trách nhiệm.

Tác giả

**Nguyễn Tú Nam**

## LỜI CẢM ƠN

Lời đầu tiên tôi muốn bày tỏ lòng biết ơn sâu sắc và kính trọng của mình tới thầy giáo, TS Nguyễn Long Giang. Trong quá trình tìm hiểu nghiên cứu để hoàn thành luận văn tôi gặp không ít khó khăn, nhưng những lúc như vậy tôi luôn nhận được sự động viên khích lệ của thầy. Thầy đã giúp đỡ tôi rất nhiều trong quá trình nghiên cứu, hướng dẫn tận tình trong cách thức và phương pháp nghiên cứu khoa học cũng như hỗ trợ tôi trong việc tìm tài liệu.

Để có được những kết quả trong luận văn này, tôi xin gửi lời cảm ơn sâu sắc đến Thầy, Cô Trường Đại học Công nghệ thông tin và Truyền thông Thái Nguyên đã tạo điều kiện cho tôi được học hỏi kiến thức thông qua các môn học cũng như hoàn thành khóa học.

Cuối cùng tôi xin bày tỏ lòng cảm ơn chân thành đến gia đình, người thân và bạn bè đồng nghiệp đã khích lệ và động viên tôi hoàn thành luận văn này.!

## MỤC LỤC

|  |      |
|--|------|
| LỜI CAM ĐOAN .....   | i    |
| LỜI CẢM ƠN .....   | iii  |
| MỤC LỤC.....   | iv   |
| Danh mục các ký hiệu, các chữ viết tắt.....  | vi   |
| Danh mục các bảng .....  | vii  |
| Danh mục các hình.....   | viii |
| MỞ ĐẦU.....  | 1    |
| Chương 1: TỔNG QUAN VỀ KHAI PHÁ DỮ LIỆU .....  | 3    |
| 1.1. Các khái niệm cơ bản trong khai phá luật kết hợp .....  | 4    |
| 1.1.1. Cơ sở dữ liệu giao tác.....   | 4    |
| 1.1.2. Tập mục thường xuyên và luật kết hợp.....   | 6    |
| 1.1.3. Bài toán khai phá luật kết hợp.....   | 8    |
| 1.2. Một số thuật toán cơ bản khai phá tập mục thường xuyên.....                                     | 9    |
| 1.2.1. Cách tiếp cận khai phá tập mục thường xuyên .....   | 9    |
| 1.2.2. Thuật toán Apriori .....  | 10   |
| 1.2.3. Thuật toán FP-growth.....   | 15   |
| 1.3. Một số hướng mở rộng bài toán khai phá tập mục thường xuyên.....                                | 24   |
| 1.4. Kết luận chương.....  | 24   |
| Chương 2: KHAI PHÁ TẬP MỤC THƯỜNG XUYÊN CÓ TRỌNG SỐ .....  | 25   |
| 2.1. Thuật toán khai phá tập mục thường xuyên có trọng số MINWAL.....                                | 25   |
| 2.1.1. Các khái niệm cơ bản .....  | 25   |
| 2.1.2. Thuật toán MINWAL khai phá tập mục thường xuyên có trọng số dựa trên thuật toán Apriori ..... | 29   |
| 2.1.3. Ví dụ minh họa thuật toán MINWAL .....  | 32   |

|   |    |
|---|----|
| 2.2. Thuật toán khai phá tập mục thường xuyên có trọng số WFIM..... | 37 |
| 2.2.1. Các khái niệm cơ bản .....                                   | 38 |
| 2.2.2. Thuật toán WFIM dựa trên thuật toán Apriori .....            | 42 |
| 2.2.3. Thuật toán WFIM dựa trên thuật toán FP-Growth.....           | 44 |
| 2.2.4. Ví dụ thuật toán WFIM .....                                  | 46 |
| 2.3. Kết luận chương.....   | 49 |
| Chương 3: ĐÁNH GIÁ CÁC THUẬT TOÁN VÀ ỨNG DỤNG.....                  | 50 |
| 3.1. Đánh giá các giải thuật .....                                  | 50 |
| 3.2. Kiểm tra tập dữ liệu và môi trường thí nghiệm .....            | 51 |
| 3.3. So sánh WFIM với các thuật toán khác .....                     | 52 |
| 3.4. Kiểm tra khả năng phát triển .....                             | 56 |
| 3.5. Ứng dụng chương trình.....                                     | 59 |
| KẾT LUẬN .....  | 63 |
| TÀI LIỆU THAM KHẢO.....   | 65 |
| PHỤ LỤC.....  | 67 |

## Danh mục các ký hiệu, các chữ viết tắt

| <b>Ký hiệu, chữ viết tắt</b> | <b>Diễn giải</b>  |
|------------------------------|---|
| <i>CSDL</i>                  | <i>Cơ sở dữ liệu</i>  |
| <i>TID</i>                   | <i>Transaction Identification</i>                           |
| <i>W</i>                     | <i>Tập các trọng số của các mục</i>                         |
| <i>L</i>                     | <i>Tập tất cả các mục thường xuyên</i>                      |
| $C_k$                        | <i>Tập các k-tập mục ứng viên</i>                           |
| $L_k$                        | <i>Tập các k-tập mục thường xuyên</i>                       |
| $SC(X)$                      | <i>Số đếm hỗ trợ của các tập mục X</i>                      |
| $WFI_k$                      | <i>Tập các k-tập mục thường xuyên có trọng số</i>           |
| $WFI$                        | <i>Tập tất cả các tập mục thường xuyên có trọng số</i>      |
| $MaxW$                       | <i>Trọng số có giá trị lớn nhất trong CSDL giao tác</i>     |
| $MinW$                       | <i>Trọng số có giá trị nhỏ nhất trong tập mục điều kiện</i> |
| $min\_weight$                | <i>Ngưỡng trọng số tối thiểu</i>                            |
| $min\_sup$                   | <i>Ngưỡng hỗ trợ tối thiểu</i>                              |
| $support$                    | <i>Độ hỗ trợ của các tập mục</i>                            |
| $conf$                       | <i>Độ tin cậy</i>   |
| $minconf$                    | <i>Độ tin cậy cực tiểu</i>                                  |
| <i>BFS</i>                   | <i>Breadth First Search</i>                                 |
| <i>DFS</i>                   | <i>Depth First Search</i>                                   |
| <i>WFIM</i>                  | <i>Weighted Frequent Itemset Mining</i>                     |

## Danh mục các bảng

|   |    |
|---|----|
| <b>Bảng 1.1.</b> Biểu diễn ngang của cơ sở dữ liệu giao tác .....                     | 5  |
| <b>Bảng 1.2.</b> Biểu diễn dọc của cơ sở dữ liệu giao tác .....                       | 5  |
| <b>Bảng 1.3.</b> Ma trận giao tác của cơ sở dữ liệu bảng 1.1 .....                    | 6  |
| <b>Bảng 1.4.</b> CSDL giao tác minh họa thực hiện thuật toán Apriori.....             | 13 |
| <b>Bảng 1.5.</b> CSDL giao tác minh họa cho thuật toán FP- growth .....               | 17 |
| <b>Bảng 2.1.</b> CSDL giao tác .....  | 27 |
| <b>Bảng 2.2.</b> Trọng số của các mục.....  | 28 |
| <b>Bảng 2.3.</b> CSDL giao tác D .....  | 32 |
| <b>Bảng 2.4.</b> Trọng số của các mục.....  | 32 |
| <b>Bảng 2.5.</b> CSDL giao tác .....  | 37 |
| <b>Bảng 2.6.</b> Ví dụ các mục với các khoảng trọng số khác nhau .....                | 38 |
| <b>Bảng 2.7.</b> Tập các tập mục thường xuyên với các khoảng trọng số khác nhau ..... | 41 |
| <b>Bảng 2.8.</b> Mục thường xuyên có trọng số (sắp xếp tăng dần theo trọng số).....   | 46 |
| <b>Bảng 3.1.</b> Tổng hợp số liệu thực tế .....                                       | 51 |
| <b>Bảng 3.2.</b> Hiệu năng đối với các ngưỡng trọng số khác nhau .....                | 56 |

## Danh mục các hình

|   |    |
|---|----|
| <b>Hình 1.1.</b> Phân loại các thuật toán khai phá tập mục thường xuyên .....   | 10 |
| <b>Hình 1.2.</b> Cây FP-tree được xây dựng dần khi thêm các giao tác $t_1, t_2, t_3$ . Từ tập dữ liệu ban đầu, ta xây dựng header table của cây FP như sau: ..... | 18 |
| <b>Hình 1.3.</b> Cây FP-tree của CSDL DB trong bảng .....   | 19 |
| <b>Hình 1.4.</b> FP-tree phụ thuộc của $m$ .....  | 22 |
| <b>Hình 1.5.</b> Các FP-tree phụ thuộc của $am, cm$ và $cam$ .....  | 22 |
| <b>Hình 2.1.</b> Cây FP-Tree tổng quát của thuật toán FP-Tree .....   | 47 |
| <b>Hình 2.2.</b> Cây FP-Tree con với tiền tố $\{r\}$ .....  | 48 |
| <b>Hình 3.1.</b> Số lượng tập mục thường xuyên so với FP-Growth (Tập dữ liệu Connect) .....   | 52 |
| <b>Hình 3.2.</b> Thời gian thực hiện so với FP-Growth (Tập dữ liệu Connect) .....   | 53 |
| <b>Hình 3.3.</b> Số lượng tập mục thường xuyên so với các thuật toán khác (Tập dữ liệu Connect) .....   | 53 |
| <b>Hình 3.4.</b> Thời gian thực hiện so với các thuật toán khác (Tập dữ liệu Connect) .....   | 54 |
| <b>Hình 3.5.</b> Thời gian thực hiện so với các thuật toán khác (Tập dữ liệu Mushroom) .....  | 55 |
| <b>Hình 3.6.</b> Thời gian thực hiện so với các thuật toán khác (Tập dữ liệu Mushroom) .....  | 55 |
| <b>Hình 3.7.</b> Khả năng phát triển của WFIM với các ngưỡng hỗ trợ khác nhau (tập dữ liệu T10I4DxK) .....  | 57 |
| <b>Hình 3.8.</b> Khả năng phát triển so với các thuật toán khác (Tập dữ liệu T10I4DxK và ngưỡng hỗ trợ = 0,1%) .....  | 58 |
| <b>Hình 3.9.</b> Khả năng phát triển so với các thuật toán khác (Tập dữ liệu T10I4DxK và ngưỡng hỗ trợ = 0,5%) .....  | 58 |

## MỞ ĐẦU

### Lý do chọn đề tài

Khai phá dữ liệu và khám phá tri thức (Data mining and Knowledge discovery) là một lĩnh vực quan trọng của ngành Công nghệ thông tin. Đây là lĩnh vực đã thu hút đông đảo các nhà khoa học trên thế giới và trong nước tham gia nghiên cứu. Khai phá luật kết hợp (Mining association rules) là bài toán có vai trò quan trọng trong nhiều nhiệm vụ khai phá dữ liệu và có nhiều ứng dụng thực tiễn trong các lĩnh vực khác nhau của đời sống, đặc biệt là trong lĩnh vực kinh doanh.

Khai phá luật kết hợp được giới thiệu bởi Agrawal vào năm 1993 khi phân tích cơ sở dữ liệu bán hàng của siêu thị, phân tích sở thích mua của khách hàng bằng cách tìm ra những mặt hàng khác nhau được khách hàng mua cùng trong một lần mua. Những thông tin như vậy sẽ giúp người quản lý kinh doanh tiếp thị chọn lọc và thu xếp không gian bày hàng hợp lý hơn, giúp cho kinh doanh hiệu quả hơn. Bài toán khai phá luật kết hợp bao gồm hai bài toán con. Bài toán thứ nhất là tìm các *tập mục thường xuyên* (Frequent itemset) thỏa mãn ngưỡng hỗ trợ tối thiểu cho trước, bài toán thứ hai là sinh ra các *luật kết hợp* (Association rule) thỏa mãn ngưỡng tin cậy cho trước từ tập mục thường xuyên tìm được. Mọi khó khăn của bài toán khai phá luật kết hợp tập trung ở bài toán thứ nhất, đó là khai phá tất cả các tập mục thường xuyên thỏa mãn ngưỡng độ hỗ trợ cho trước, và các nghiên cứu về khai phá luật kết hợp tập trung vào bài toán khai phá tập mục thường xuyên.

Xuất phát từ những lợi ích thực tế trên tác giả đã mạnh dạn chọn đề tài “**KHAI PHÁ TẬP MỤC THƯỜNG XUYỀN CÓ TRỌNG SỐ TRÊN CƠ SỞ DỮ LIỆU GIAO TÁC**” làm đề tài nghiên cứu cho luận văn tốt nghiệp của mình.

**Mục tiêu đề tài** tiếp tục nghiên cứu và đề xuất các thuật toán khai phá tập mục thường xuyên có trọng số trong CSDL giao tác

Xây dựng và đề xuất một số giải thuật khai phá tập mục thường xuyên có trọng số.

Lập trình, thử nghiệm các giải thuật khai phá tập mục thường xuyên có trọng số.

**Đối tượng nghiên cứu** các cơ sở dữ liệu giao tác được cập nhật từ kho dữ liệu mẫu UCI

**Phạm vi nghiên cứu** nghiên cứu và thử nghiệm bài toán khai phá tập mục thường xuyên có trọng số trên cơ sở dữ liệu giao tác.

**Phương pháp nghiên cứu** luận văn là nghiên cứu lý thuyết và nghiên cứu thực nghiệm. Về nghiên cứu lý thuyết: các định lý, mệnh đề trong luận văn được chứng minh dựa vào các kiến thức cơ bản và các kết quả nghiên cứu đã công bố. Về nghiên cứu thực nghiệm luận văn thực hiện cài đặt các thuật toán, chạy thử nghiệm thuật toán.

### **Bố cục luận văn**

Luận văn được chia làm 3 chương:

*Chương 1: TỔNG QUAN VỀ KHAI PHÁ DỮ LIỆU*

*Chương 2: KHAI PHÁ TẬP MỤC THƯỜNG XUYỀN CÓ TRỌNG SỐ*

*Chương 3: ĐÁNH GIÁ CÁC THUẬT TOÁN VÀ ỨNG DỤNG*