

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

**BÙI VĂN THẮNG**

**LUẬT KẾT HỢP MỜ VÀ ỨNG DỤNG  
ĐỐI VỚI MỘT SỐ BÀI TOÁN DỰ BÁO**

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

**THÁI NGUYÊN - 2014**

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

**BÙI VĂN THẮNG**

**LUẬT KẾT HỢP MỜ VÀ ỨNG DỤNG  
ĐỐI VỚI MỘT SỐ BÀI TOÁN DỰ BÁO**

**Chuyên ngành: KHOA HỌC MÁY TÍNH**

**Mã số: 60 48 01**

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

**HƯỚNG DẪN KHOA HỌC: TS. VŨ VINH QUANG**

**THÁI NGUYÊN - 2014**

## LỜI CAM ĐOAN

Tôi xin cam đoan đây là công trình nghiên cứu của riêng tôi dưới sự hướng dẫn trực tiếp của **TS. Vũ Vinh Quang**.

Mọi trích dẫn sử dụng trong báo cáo này đều được ghi rõ nguồn tài liệu tham khảo theo đúng qui định.

Mọi sao chép không hợp lệ, vi phạm quy chế đào tạo, hay gian trá, tôi xin chịu hoàn toàn trách nhiệm.

*Thái Nguyên, ngày 27 tháng 8 năm 2014*

**Tác giả**

**Bùi Văn Thắng**

## MỤC LỤC

<b>LỜI CAM ĐOAN .....</b>	<b>i</b>
<b>MỤC LỤC .....</b>	<b>ii</b>
<b>CÁC KÝ HIỆU VÀ CHỮ VIẾT TẮT .....</b>	<b>iii</b>
<b>DANH MỤC BẢNG .....</b>	<b>iv</b>
<b>DANH MỤC HÌNH VẼ.....</b>	<b>v</b>
<b>MỞ ĐẦU.....</b>	<b>1</b>
<b>CHƯƠNG 1. MỘT SỐ KIẾN THỨC CƠ BẢN VỀ KHAI PHÁ DỮ LIỆU.....</b>	<b>3</b>
<b>1.1. Khái niệm cơ bản về khai phá dữ liệu.....</b>	<b>3</b>
1.1.1. Giới thiệu.....	3
1.1.2. Khái niệm khai phá dữ liệu .....	4
<b>1.2. Một số hướng nghiên cứu trong khai phá dữ liệu .....</b>	<b>5</b>
1.2.1. Một số hướng nghiên cứu .....	5
1.2.2. Các dạng dữ liệu có thể khai phá .....	8
<b>1.3. Nhiệm vụ chính của khai phá dữ liệu.....</b>	<b>8</b>
1.3.1. Phân lớp (Classification).....	9
1.3.2. Hồi quy (Regression) .....	9
1.3.3. Khai phá luật kết hợp (Association rule) .....	9
1.3.4. Gom nhóm (Clustering) .....	9
1.3.5. Tổng hợp (Summarization) .....	10
1.3.6. Mô hình ràng buộc (Dependency modeling) .....	10
1.3.7. Dò tìm biến đổi và độ lệch (Change and Deviation Dectection) .....	10
<b>1.4. Bài toán khai phá luật kết hợp.....</b>	<b>10</b>
1.4.1. Bài toán .....	10
1.4.2. Một số thuật toán cơ bản .....	15
<b>1.5. Logic mờ.....</b>	<b>23</b>
1.5.1. Định nghĩa tập mờ.....	23
1.5.2. Độ cao, miền xác định và miền tin cậy của tập mờ .....	25
1.5.3. Các phép toán logic trên tập mờ.....	26
1.5.4. Biến ngôn ngữ và giá trị của nó .....	27
<b>1.6. Kết luận.....</b>	<b>28</b>

<b>CHƯƠNG 2. KHAI PHÁ LUẬT KẾT HỢP MỜ .....</b>	<b>30</b>
<b>2.1. Rời rạc hóa thuộc tính dựa vào tập mờ .....</b>	<b>30</b>
2.1.1. Luật kết hợp với thuộc tính số.....	30
2.1.2. Các phương pháp rời rạc hóa .....	30
<b>2.2. Luật kết hợp mờ .....</b>	<b>33</b>
2.2.1. Rời rạc hóa thuộc tính mờ.....	33
2.2.2. Luật kết hợp mờ .....	35
<b>2.3. Thuật toán khai phá luật kết hợp mờ dựa trên thuật toán Apriori.....</b>	<b>37</b>
<b>2.4. Khai phá luật kết hợp mờ dựa trên thuật toán Fp-Growth.....</b>	<b>40</b>
2.4.1. Thuật toán xây dựng cây CUFP-Tree .....	40
2.4.2. Thuật toán tìm tập phổ biến FP-Growth dựa trên cây CUFP-Tree.....	41
<b>2.5. Ví dụ thử nghiệm.....</b>	<b>42</b>
2.5.1. Xây dựng cây CUFP-Tree.....	42
2.5.2. Thuật toán tìm tập phổ biến .....	45
<b>2.6. Kết luận .....</b>	<b>46</b>
<b>CHƯƠNG 3. ỨNG DỤNG KHAI PHÁ DỮ LIỆU TRONG MÔ HÌNH DỰ</b>	
<b>BÁO .....</b>	<b>48</b>
<b>3.1. Mô hình một số bài toán dự báo .....</b>	<b>48</b>
3.1.1. Giới thiệu.....	48
3.1.2. Một mô hình dự báo là gì? .....	49
3.1.3. Các kỹ thuật mô hình hóa dự báo phổ biến.....	51
<b>3.2. Xây dựng các luật kết hợp mờ trong mô hình dự báo.....</b>	<b>55</b>
<b>3.3. Một số kết quả thực nghiệm.....</b>	<b>55</b>
3.3.1. Môi trường thử nghiệm .....	55
3.3.2. Kết quả thử nghiệm với CSDL gồm 20 giao dịch .....	60
3.3.3. Kết quả thử nghiệm.....	61
<b>PHẦN KẾT LUẬN .....</b>	<b>62</b>
<b>TÀI LIỆU THAM KHẢO.....</b>	<b>63</b>

**CÁC KÝ HIỆU VÀ CHỮ VIẾT TẮT**

CNTT	Công nghệ thông tin
CSDL	Cơ sở dữ liệu
KPDL	Khai phá dữ liệu
KDD	Knowledge Discovery in Databases
ItemSet	Tập mục
Item	Mục

## DANH MỤC BẢNG

Bảng 1.1: Cơ sở dữ liệu giao tác .....	11
Bảng 1.2: Kết quả thuật toán Apriori .....	16
Bảng 1.3: Những biến đổi dữ liệu của FP-Growth.....	19
Bảng 2.1: CSDL thống kê dân số của 10 gia đình [21].....	31
Bảng 2.2: Rời rạc hóa thuộc tính số rời rạc hữu hạn hoặc thuộc tính hạng mục .....	31
Bảng 2.3: Rời rạc hóa thuộc tính số “Tuổi” .....	32
Bảng 2.4: Bảng các ký hiệu sử dụng trong thuật toán khai phá luật kết hợp mờ.....	38
Bảng 2.5: Bảng các ký hiệu sử dụng trong thuật toán.....	40
Bảng 2.6: Cơ sở dữ liệu mờ.....	42
Bảng 2.7: Kết quả sau khi thực hiện Bước 1 .....	42
Bảng 2.8: Header_Table .....	43
Bảng 2.9: CSDL mờ sau khi đã cập nhật .....	43
Bảng 2.10: Tập phổ biến .....	46
Bảng 3.1: Giao tác ví dụ trong CSDL FAM95.....	56
Bảng 3.2: CSDL giao tác Bảng 3.1 sau khi mờ hóa.....	57

## DANH MỤC HÌNH VẼ

Hình 1.1: Quá trình khai phá tri thức trong CSDL.....	3
Hình 1.2: FP-tree của dữ liệu Bảng 1.1 .....	20
Hình 1.3: Thành phần của FP-tree.....	21
Hình 1.4: Hàm thuộc $\mu_{Ax}$ của tập kinh điển A .....	23
Hình 1.5: Hàm thuộc của tập mờ B .....	24
Hình 1.6: Hàm thuộc của tập mờ C .....	24
Hình 1.7: Hàm thuộc $\mu_F(x)$ có mức chuyển đổi tuyến tính .....	25
Hình 1.8: Mô tả giá trị ngôn ngữ bằng tập mờ .....	27
Hình 2.1: Hàm thuộc của các tập mờ “Tuổi_trẻ”, “Tuổi_trung_niên”, và “Tuổi_già” ....	33
Hình 2.2: Kết quả xử lý giao dịch đầu tiên.....	44
Hình 2.3: Kết quả xử lý giao dịch đầu tiên.....	45
Hình 2.4: Cây CUFP-TREE .....	45
Hình 3.1: Hai khách hàng và các đặc tính đầu vào của họ.....	50
Hình 3.2: Dữ liệu khách hàng gồm các đặc tính đầu vào và kết quả đầu ra được cung cấp cho một mô hình dự báo trong quá trình huấn luyện.....	50
Hình 3.3: Khung nhìn hai chiều của một siêu phẳng tối ưu chia tách dữ liệu và các vector hỗ trợ .....	52
Hình 3.4: Khung nhìn hai chiều về kết quả của việc phân cụm một tập dữ liệu đầu vào thành hai cụm: các hình tam giác màu xanh lá cây và các hình vuông màu đỏ.....	53
Hình 3.5: Mạng nơ-ron hướng thuận với tầng đầu vào, tầng ẩn và tầng đầu ra.....	54
Hình 3.6: Giao diện chương trình, 20 giao dịch mờ.....	60
Hình 3.7: Các tập phổ biến tìm được.....	60
Hình 3.8: Luật kết hợp khai phá.....	61
Hình 3.9: Kết quả thử nghiệm với hai thuật toán Apriori mờ và thuật toán CUFP .....	61



## MỞ ĐẦU

### 1. Đặt vấn đề

Khai phá dữ liệu là một lĩnh vực nghiên cứu quan trọng trong lý thuyết về cơ sở dữ liệu, có nhiều ứng dụng trong đời sống xã hội. Mục đích chính là nhằm phát hiện những thông tin mới, các luật mới từ cơ sở dữ liệu đã có hay một cách tổng quát hơn là từ các kho dữ liệu. Rất nhiều lĩnh vực ứng dụng trong thực tiễn đều sử dụng công cụ khai phá dữ liệu và tìm kiếm tri thức. Trong lý thuyết về khai phá dữ liệu, khai phá luật kết hợp đang được quan tâm nghiên cứu nhiều trên thế giới. Một số hướng nghiên cứu đã và đang được các chuyên gia công nghệ thông tin tập chung nghiên cứu là: nghiên cứu thiết kế các hệ mờ cho các ứng dụng cụ thể như hệ trợ giúp quyết định, hệ điều khiển dựa trên hệ tri thức luật, hệ phân loại dựa trên hệ tri thức luật, hệ phân loại dựa trên lập luận dựa trên hệ luật ứng dụng trong các lĩnh vực như: kinh doanh, thị trường chứng khoán và dự đoán thị trường, công nghệ sinh học, giáo dục và đào tạo,...

#### Một số hướng nghiên cứu trong khai phá dữ liệu

- **Luật kết hợp nhị phân:** Đây là hướng nghiên cứu đầu tiên của luật kết hợp. Thuật toán tiêu biểu là Apriori.
- **Luật kết hợp có thuộc tính số và thuộc tính hạng mục:** Nghiên cứu các hệ CSDL có thuộc tính số hoặc thuộc tính hạng mục bằng cách rời rạc hóa dữ liệu cho thuộc tính số để chuyển chúng về thuộc tính nhị phân.
- **Luật kết hợp mờ:** Phương pháp rời rạc hóa dữ liệu có thuộc tính số và thuộc tính hạng mục gặp phải vấn đề “điểm biên gãy”. Để khắc phục điều này, các nhà nghiên cứu đề xuất sử dụng lý thuyết tập mờ và xây dựng các luật kết hợp dạng mờ.
- **Luật kết hợp có trọng số:** Sử dụng phương pháp tính độ hỗ trợ cho các tập mục dựa trên trọng số của các tập mục.

Ngoài ra, còn một số hướng nghiên cứu: khai phá luật kết hợp song song, khai phá luật kết hợp nhiều mức, luật kết hợp tiếp cận theo hướng tập thô,...

Luận văn tập trung nghiên cứu vào khai phá **Luật kết hợp mờ và ứng dụng đối với bài toán dự báo.**

## 2. Hướng nghiên cứu của đề tài

- Nghiên cứu lý thuyết tập mờ.
- Nghiên cứu khai phá dữ liệu và khai phá dữ liệu mờ trên CSDL. Tìm hiểu một số thuật toán trong khai phá dữ liệu: Apriori mờ, thuật toán FP Growth, thuật toán biểu diễn dữ liệu giao dịch mờ dựa trên cây FP-Tree.
- Cài đặt thử nghiệm một số thuật toán khai phá dữ liệu mờ và thử nghiệm trên một số bộ dữ liệu. Đánh giá kết quả sau khi thử nghiệm.

## 3. Đối tượng nghiên cứu

- Nghiên cứu phương pháp luận cho phép phát hiện tri thức dạng luật mờ, như luật kết hợp mờ, luật mờ với thuộc tính có trọng số,... từ các kho dữ liệu.
- Cơ sở lý thuyết của việc nghiên cứu lập luận xấp xỉ dựa trên lý thuyết tập mờ, phương pháp tính toán các thông tin mờ, đánh giá các phương pháp để lấy quyết định.
- Ứng dụng luật kết hợp mờ đối với một số bài toán dự báo.

## 4. Kết quả đạt được

- Tìm hiểu thuật toán nén dữ liệu giao dịch mờ dựa trên cây FP Tree, khai phá tập phổ biến dựa trên cây đã xây dựng. Đây là hướng nghiên cứu mới, giúp làm giảm thời gian khai phá tập phổ biến rất nhiều so với thuật toán Apriori mờ.
- Cài đặt thử nghiệm thuật toán Apriori mờ và thuật toán khai phá luật kết hợp mờ dựa trên thuật toán Fp-Growth.
- Thử nghiệm hai thuật toán trên với một số bộ dữ liệu, so sánh các kết quả đã thu được sau khi thử nghiệm.

## 5. Bố cục của luận văn

Phần mở đầu

Chương 1: Một số kiến thức cơ bản về khai phá dữ liệu

Chương 2: Khai phá luật kết hợp mờ

Chương 3: Ứng dụng khai phá dữ liệu trong mô hình dự báo

Kết luận

Tài liệu tham khảo

Số hóa bởi Trung tâm Học liệu

<http://www.lrc-tnu.edu.vn/>