

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

VŨ VĂN TIỆP

**NGHIÊN CỨU MỘT SỐ THUẬT TOÁN GIA TĂNG CHO VIỆC RÚT
GỌN CÁC THUỘC TÍNH TRONG BẢNG QUYẾT ĐỊNH KHÔNG
ĐẦY ĐỦ**

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

HƯỚNG DẪN: GS.TS VŨ ĐỨC THI

THÁI NGUYÊN 2015

LỜI CẢM ƠN

Em xin chân thành cảm ơn và biết ơn sâu sắc đến GS.TS Vũ Đức Thi, Viện Công nghệ thông tin – Đại học Quốc gia Hà Nội. Người đã tận tình hướng dẫn và giúp đỡ em hoàn thành luận văn này.

Em xin chân thành cảm ơn các Thầy ở Viện Công nghệ thông tin đã dạy bảo, giúp đỡ và truyền đạt kiến thức cho em trong suốt khóa học và quá trình em làm luận văn.

Em xin chân thành cảm ơn các Thầy, các Cô ở trường Đại học Công nghệ thông tin và truyền thông Thái Nguyên đã tận tình dạy bảo, động viên, giúp đỡ và tạo điều kiện cho em trong suốt thời gian học tập và nghiên cứu.

Cuối cùng xin chân thành cảm ơn bạn bè, người thân và gia đình luôn là người đồng hành, động viên, chia sẻ những khó khăn trong suốt thời gian hoàn thành luận văn.

Học viên

Vũ Văn Tiệp

LỜI CAM ĐOAN

Tôi xin cam đoan đề tài "*Nghiên cứu một số thuật toán gia tăng cho việc rút gọn các thuộc tính trong bảng quyết định không đầy đủ*" là công trình nghiên cứu được tôi thực hiện dưới sự hướng dẫn của giáo viên hướng dẫn khoa học.

Một số Định nghĩa, Định lý, Tính chất, Mệnh đề và Thuật toán tôi lấy từ nguồn tài liệu chính xác có trích dẫn tên tài liệu và tên tác giả rõ ràng. Tôi xin chịu trách nhiệm về luận văn của mình.

Học viên

Vũ Văn Tiệp

MỤC LỤC

LỜI CẢM ƠN	i
LỜI CAM ĐOAN	iii
Danh mục các thuật ngữ.....	vi
Danh sách bảng	vii
MỞ ĐẦU	1
Chương 1. TỔNG QUAN.....	4
1.1. Hệ thông tin đầy đủ và mô hình tập thô truyền thống	4
1.1.1. Hệ thông tin đầy đủ.....	4
1.1.2. Bảng quyết định đầy đủ	7
1.1.3. Tập rút gọn và tập lõi.....	7
1.2. Hệ thông tin không đầy đủ và mô hình tập thô dung sai.....	8
1.2.1. Hệ thông tin không đầy đủ.....	9
1.2.2. Bảng quyết định không đầy đủ	10
1.3. Rút gọn thuộc tính trong bảng quyết định không đầy đủ.....	11
1.3.1. Tổng quan về các phương pháp rút gọn thuộc tính	11
1.3.2. Phân nhóm các phương pháp rút gọn thuộc tính	14
1.4. Kết luận chương 1	17
Chương 2. CÁCH TIẾP CẬN GIA TĂNG RÚT GỌN THUỘC TÍNH TRONG BẢNG QUYẾT ĐỊNH KHÔNG ĐẦY ĐỦ KHI BỔ SUNG, LOẠI BỎ TẬP THUỘC TÍNH. 18	
2.1. Rút gọn thuộc tính sử dụng hàm phân biệt mở rộng	18
2.1.1. Ma trận phân biệt và hàm phân biệt mở rộng.....	19
2.1.2. Rút gọn thuộc tính sử dụng hàm phân biệt mở rộng.....	21
2.2. Các thuật toán tiếp cận gia tăng tìm tập rút gọn khi bổ sung, loại bỏ tập thuộc tính	25
2.2.1. Thuật toán tìm tập rút gọn khi bổ sung tập thuộc tính.....	25
2.2.2. Thuật toán tìm tập rút gọn khi loại bỏ tập thuộc tính	29
2.3. Kết luận chương 2.....	Error! Bookmark not defined.
Chương 3. THỬ NGHIỆM VÀ ĐÁNH GIÁ KẾT QUẢ	34
3.1. Bài toán	34
3.2. Phân tích, lựa chọn công cụ	34

3.2.1. Thuật toán tìm tập rút gọn sử dụng hàm phân biệt mở rộng	
3.2.2. Các thuật toán tìm tập rút gọn khi bổ sung và loại bỏ tập thuộc tính	37
3.3. Đánh giá kết quả thử nghiệm	39
3.3.1. Kết quả thử nghiệm thuật toán tìm tập rút gọn sử dụng hàm phân biệt mở rộng. 39	
3.3.2. Kết quả thử nghiệm thuật toán tìm tập rút gọn khi bổ sung tập thuộc tính	41
3.3.3. Kết quả thử nghiệm thuật toán tìm tập rút gọn khi loại bỏ tập thuộc tính	45
KẾT LUẬN	49
Tài liệu tham khảo	50
Phụ lục	52

Danh mục các thuật ngữ

Thuật ngữ tiếng Việt	Thuật ngữ tiếng Anh
<i>Tập thô</i>	<i>Rough Set</i>
<i>Tập thô dung sai</i>	<i>Tolerance Rough Set</i>
<i>Hệ thông tin</i>	<i>Information System</i>
<i>Hệ thông tin đầy đủ</i>	<i>Complete Information System</i>
<i>Hệ thông tin không đầy đủ</i>	<i>Incomplete Information System</i>
<i>Bảng quyết định</i>	<i>Decision Table</i>
<i>Bảng quyết định đầy đủ</i>	<i>Complete Decision Table</i>
<i>Bảng quyết định không đầy đủ</i>	<i>Incomplete Decision Table</i>
<i>Quan hệ không phân biệt được</i>	<i>Indiscernibility Relation</i>
<i>Quan hệ dung sai</i>	<i>Tolerance Relation</i>
<i>Xấp xỉ dưới</i>	<i>Lower Approximation</i>
<i>Xấp xỉ trên</i>	<i>Upper Approximation</i>
<i>Rút gọn thuộc tính</i>	<i>Attribute Reduction</i>
<i>Tập rút gọn</i>	<i>Reduct</i>
<i>Tập lõi</i>	<i>Core</i>

Danh sách bảng

Bảng 1.1. Bảng thông tin về bệnh cúm.....	6
Bảng 1.2. Bảng quyết định không đủ về các xe hơi.....	10
Bảng 1.3. Các phương pháp rút gọn thuộc tính trong công trình [3, 8].....	13
Bảng 2.1. Bảng quyết định không đầy đủ mô tả về các tivi.....	19
Bảng 2.1. Bảng quyết định không đầy đủ mô tả về các tivi (tiếp theo).....	24
Bảng 2.3. Bảng quyết định không đầy đủ về tivi khi bổ sung tập thuộc tính.....	26
Bảng 3.1. Kết quả thực hiện Thuật toán 2.1 và Thuật toán MBAR.....	39
Bảng 3.2. Tập rút gọn của Thuật toán 2.1 và Thuật toán MBAR.....	40
Bảng 3.3. Kết quả thực hiện Thuật toán 2.1 trên bộ số liệu ban đầu.....	42
Bảng 3.4 Kết quả thực hiện Thuật toán 2.1 sau khi lấy ngẫu nhiên 60% số thuộc tính điều kiện.....	42
Bảng 3.5 Kết quả thực hiện Thuật toán 2.2 tìm tập rút gọn khi bổ sung 40% số thuộc tính vào.....	43
Bảng 3.6. Kết quả thực hiện Thuật toán 2.1 trên bộ số liệu ban đầu.....	45
Bảng 3.7 Kết quả thực hiện Thuật toán 2.1 sau khi loại ngẫu nhiên 40% số thuộc tính điều kiện.....	46
Bảng 3.8 Kết quả thực hiện Thuật toán 2.3 tìm tập rút gọn khi loại bỏ 40% số thuộc tính điều kiện.....	47

MỞ ĐẦU

Lý thuyết tập thô - do Zdzislaw Pawlak [10] đề xuất vào những năm đầu thập niên tám mươi của thế kỷ hai mươi - được xem là công cụ hữu hiệu để giải quyết các bài toán phân lớp, phát hiện luật...chứa dữ liệu không đầy đủ, không chắc chắn. Từ khi xuất hiện, lý thuyết tập thô đã được sử dụng hiệu quả trong các bước của quá trình khai phá dữ liệu và khám phá tri thức, bao gồm tiền xử lý số liệu, khai phá dữ liệu và đánh giá kết quả thu được. Rút gọn thuộc tính và trích lọc luật quyết định (luật phân lớp) là hai ứng dụng chính của lý thuyết tập thô trong khai phá dữ liệu. Rút gọn thuộc tính thuộc giai đoạn tiền xử lý dữ liệu còn trích lọc luật thuộc giai đoạn khai phá dữ liệu. Mục tiêu của rút gọn thuộc tính là loại bỏ các thuộc tính dư thừa nhằm tìm tập con nhỏ nhất của tập thuộc tính điều kiện (tập rút gọn) mà bảo toàn thông tin phân lớp của bảng quyết định. Dựa trên tập rút gọn thu được, việc sinh luật và phân lớp đạt hiệu quả cao nhất.

Trong các bài toán thực tế, các bảng quyết định thường thiếu giá trị trên miền giá trị thuộc tính, gọi là các bảng quyết định không đầy đủ. Trên bảng quyết định không đầy đủ, Kryszkiewicz [5] đã mở rộng quan hệ tương đương trong lý thuyết tập thô truyền thống thành quan hệ dung sai và đề xuất mô hình tập thô dung sai nhằm trích lọc luật trực tiếp không qua bước xử lý giá trị thiếu. Dựa trên mô hình tập thô dung sai, một số công trình công bố trong mấy năm gần đây đã đề xuất một số độ đo không chắc chắn nhằm giải quyết bài toán rút gọn thuộc tính và trích lọc luật, đáng chú ý là các công bố được liệt kê trong công trình [8].

Luận văn đặt ra hai mục tiêu chính:

1) Tổng hợp các công bố về các phương pháp rút gọn thuộc tính trong bảng quyết định không đầy đủ theo tiếp cận mô hình tập thô dung sai, trên cơ

sở đó nghiên cứu phương pháp gia tăng rút gọn thuộc tính sử dụng hàm phân biệt mở rộng trong trường hợp bổ sung, loại bỏ tập thuộc tính. Bao gồm:

- Nghiên cứu phương pháp rút gọn thuộc tính trong bảng quyết định không đầy đủ sử dụng hàm phân biệt mở rộng, gồm các bước: xây dựng hàm phân biệt mở rộng; định nghĩa tập rút gọn và độ quan trọng của thuộc tính dựa trên hàm phân biệt mở rộng; xây dựng thuật toán heuristic tìm một tập rút gọn tốt nhất sử dụng hàm phân biệt mở rộng; phân nhóm phương pháp sử dụng hàm phân biệt mở rộng.

- Nghiên cứu hướng tiếp cận gia tăng rút gọn thuộc tính trong bảng quyết định không đầy đủ sử dụng hàm phân biệt mở rộng trong trường hợp bổ sung, loại bỏ tập thuộc tính.

2) Cài đặt thuật toán rút gọn thuộc tính trong bảng quyết định không đầy đủ sử dụng hàm phân biệt mở rộng và các thuật toán gia tăng trong trường hợp bổ sung, loại bỏ tập thuộc tính. Thử nghiệm và đánh giá kết quả trên các bộ số liệu từ kho dữ liệu UCI.

Đối tượng nghiên cứu của luận văn là các *bảng quyết định không đầy đủ* khi bổ sung, loại bỏ tập thuộc tính.

Phạm vi nghiên cứu của luận văn *tập trung* vào bài toán rút gọn thuộc tính ở bước tiền xử lý số liệu trong quá trình khai phá dữ liệu.

Phương pháp nghiên cứu của luận văn là nghiên cứu lý thuyết và nghiên cứu thực nghiệm. Về nghiên cứu lý thuyết: tổng hợp và nắm bắt 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 với các bộ số liệu lấy từ kho dữ liệu UCI [13], so sánh và đánh giá nghiên cứu thực nghiệm với nghiên cứu lý thuyết.

Bố cục của luận văn gồm phần mở đầu và hai chương nội dung, phần kết luận và danh mục các tài liệu tham khảo.

Chương 1 trình bày các khái niệm cơ bản về lý thuyết tập thô của Pawlak [10] và mô hình tập thô mở rộng dựa trên quan hệ dung sai, gọi tắt là mô hình tập thô dung sai [5]. Trình bày tổng quan các kết quả nghiên cứu về các phương pháp rút gọn thuộc tính trong bảng quyết định không đầy đủ theo tiếp cận mô hình tập thô dung sai.

Chương 2 trình bày hai nội dung chính:

- Thứ nhất là phương pháp rút gọn thuộc tính sử dụng hàm phân biệt mở rộng [14], bao gồm: xây dựng hàm phân biệt mở rộng; định nghĩa tập rút gọn và độ quan trọng của thuộc tính dựa trên hàm phân biệt mở rộng; xây dựng thuật toán heuristic tìm một tập rút gọn tốt nhất sử dụng hàm phân biệt mở rộng; phân nhóm phương pháp sử dụng hàm phân biệt mở rộng.

- Thứ hai là xây dựng thuật toán theo hướng tiếp cận gia tăng tìm tập rút gọn của bảng quyết định không đầy đủ sử dụng hàm phân biệt mở rộng trong trường hợp bổ sung, loại bỏ tập thuộc tính [14].

Chương 3 trình bày kết quả thử nghiệm và đánh giá các thuật toán: bao gồm thuật toán tìm tập rút gọn sử dụng hàm phân biệt mở rộng và thuật toán gia tăng tìm tập rút gọn sử dụng hàm phân biệt mở rộng trong trường hợp bổ sung và loại bỏ tập thuộc tính. Thử nghiệm được thực hiện trên các bộ số liệu mẫu từ kho dữ liệu UCI [13].

Cuối cùng, phần kết luận nêu những đóng góp của luận văn và hướng phát triển tiếp theo.