

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

CAO VĂN NGUYỄN

**NGHIÊN CỨU MỘT SỐ KỸ THUẬT KHAI
PHÁ DỮ LIỆU KHÔNG GIAN SỬ DỤNG
CÂY QUYẾT ĐỊNH**

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Thái Nguyên - 2013

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

CAO VĂN NGUYỄN

**NGHIÊN CỨU MỘT SỐ KỸ THUẬT KHAI
PHÁ DỮ LIỆU KHÔNG GIAN SỬ DỤNG
CÂY QUYẾT ĐỊNH**

Chuyên ngành: Khoa học máy tính

Mã số: 60 48 01

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

NGƯỜI HƯỚNG DẪN KHOA HỌC: PGS.TS. ĐẶNG VĂN ĐỨC

Thái Nguyên – 2013

LỜI CẢM ƠN

Tôi xin gửi lời cảm ơn chân thành nhất tới PGS. TS Đặng Văn Đức - người đã hướng dẫn, chỉ bảo tận tình, cung cấp tài liệu và phương pháp luận nghiên cứu khoa học để tôi hoàn thành bản luận văn này.

Tôi xin bày tỏ lòng cảm ơn sâu sắc tới thầy cô, bạn bè cùng khóa, cùng lớp đã giúp đỡ tôi trong suốt những năm học qua.

Xin cảm ơn gia đình, bạn bè, những người luôn khuyến khích, động viên và giúp đỡ tôi trong mọi hoàn cảnh khó khăn.

Tôi xin cảm ơn các thầy cô trong trường Đại học Công nghệ thông tin và Truyền thông, Đại học Thái Nguyên đã hết sức tạo điều kiện cho tôi trong quá trình học và làm luận văn này.

Luận văn được hoàn thành trong thời gian hạn hẹp nên không thể tránh được những thiếu sót. Tôi xin cảm ơn thầy cô, bạn bè, đồng nghiệp đã có những ý kiến đóng góp chân thành cho nội dung của luận văn, để tôi có thể tiếp tục đi sâu tìm hiểu về lĩnh vực này trong tương lai.

Thái Nguyên, 11/2013

Cao Văn Nguyên

caonguyenvp@gmail.com

LỜI CAM ĐOAN

Tôi xin cam đoan 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 nội dung luận văn, những điều đã được trình bày hoặc là của riêng cá nhân tôi hoặc là được tổng hợp từ nhiều nguồn tài liệu. Tất cả các nguồn tài liệu tham khảo được dùng đều có xuất xứ rõ ràng, được trích dẫn hợp pháp.

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

Thái Nguyên, 11/2013

Cao Văn Nguyên

MỤC LỤC

TRANG

Trang phụ bìa	
Lời cảm ơn.....	i
Lời cam đoan	ii
Mục lục	iii
Danh mục các ký hiệu, các chữ viết tắt	iv
Danh mục các bảng	vi
Danh mục các hình (hình vẽ, ảnh chụp, đồ thị...)	vii
MỞ ĐẦU	1
CHƯƠNG I TỔNG QUAN VỀ DỮ LIỆU KHÔNG GIAN VÀ KHAI PHÁ DỮ LIỆU	3
1.1. Tổng quan về dữ liệu không gian địa lý	3
1.1.1. Một số khái niệm	3
1.1.2. Mô hình dữ liệu Vector.....	5
1.1.3. Quan hệ không gian giữa các đối tượng địa lý	8
1.2. Khai phá dữ liệu	8
1.2.1. Định nghĩa khai phá dữ liệu.....	8
1.2.2. Nhiệm vụ chính trong khai phá dữ liệu	9
1.2.3. Các phương pháp khai phá dữ liệu	11
1.3. Cây quyết định.....	13
1.3.1. Khái niệm.....	13
1.3.2. Ưu điểm và nhược điểm của cây quyết định	14
1.3.3. Xây dựng cây quyết định	14
CHƯƠNG 2 KHAI PHÁ DỮ LIỆU KHÔNG GIAN SỬ DỤNG CÂY QUYẾT ĐỊNH	18
2.1. Phân lớp dữ liệu.....	18
2.2. Cây quyết định ứng dụng trong phân lớp dữ liệu.....	20
2.2.1. Thuật toán ID 3	21
2.2.2. Thuật toán C4.5.....	28
2.3. Xây dựng cây quyết định trong khai phá dữ liệu không gian	34
2.3.1. Tư tưởng xây dựng thuật toán.....	34
2.3.2. Thuật toán cây quyết định không gian mở rộng từ ID3.....	36
2.3.3. Ví dụ xây dựng cây quyết định không gian.....	38
2.3.4. Tìm hiểu, đề xuất phân lớp dữ liệu không gian sử dụng cây quyết định.....	43
CHƯƠNG 3 CÀI ĐẶT CHƯƠNG TRÌNH THỬ NGHIỆM.....	55
3.1. Giới thiệu	55
3.2. Lựa chọn công nghệ	55
3.3. Dữ liệu thử nghiệm.....	56

3.4. Thiết kế chương trình	59
3.5. Cài đặt chương trình	60
3.6. Đánh giá kết quả thử nghiệm.....	61
KẾT LUẬN	68
TÀI LIỆU THAM KHẢO	69
PHỤ LỤC	70

DANH MỤC CÁC KÝ HIỆU, CÁC CHỮ VIẾT TẮT

CSDL	Cơ sở dữ liệu
GIS	Geographic information system
Object ID	Identifier of objects
SDT	Spatial Decision Tree
SJI	Spatial Join Index
SJR	Spatial Join Relation
SpatRel	Spatial Relation
SpatMes	Spatial Measure
SQL	Structured Query Language

DANH MỤC CÁC BẢNG

Bảng 1.1: Topology vùng	6
Bảng 1.2: Topology nút	6
Bảng 1.3: Topology cung	6
Bảng 1.4: Dữ liệu tọa độ cung	7
Bảng 1.5: Mô tả dữ liệu đặc trưng cấu trúc Spaghetti.....	7
Bảng 2.1: Dữ liệu thời tiết	24
Bảng 2.2: So sánh Gain của các thuộc tính tại nút gốc	24
Bảng 2.3: So sánh Gain trong nhánh "Quang cảnh" = "Nắng"	25
Bảng 2.4: So sánh Gain trong nhánh "Quang cảnh" = "Mưa"	26
Bảng 2.5: Dữ liệu thời tiết xét thuộc tính độ ẩm dạng số.....	31
Bảng 2.6: Bảng tính Gain	32
Bảng 2.7: Dữ liệu thời tiết xét thuộc tính ngày	33
Bảng 2.8: Bảng quan hệ không gian	40
Bảng 2.9: Bảng quan hệ không gian và độ đo không gian	45
Bảng 2.10: Bảng quan hệ không gian đã lược thuộc tính Object ID.....	48
Bảng 2.11: Bảng quan hệ không gian: khoảng cách đến sông gần nhất	49
Bảng 2.12: Bảng quan hệ không gian rút gọn (đầu vào thuật toán).....	50
Bảng 3.1: Bảng dữ liệu đầu vào	62
Bảng 3.2: Tính Gain cho các thuộc tính dự đoán tại nút gốc	64
Bảng 3.3: Tính Gain các thuộc tính dự đoán nhánh BTScover="Low"	65
Bảng 3.4: Tính Gain các thuộc tính dự đoán nhánh Density="Medium"	66
Bảng 3.5: Tính Gain các thuộc tính dự đoán nhánh Density="Low"	66

DANH MỤC CÁC HÌNH

Hình 1.1: Đối tượng dữ liệu cơ bản điểm, đường vùng	4
Hình 1.2: Biểu diễn đối tượng bằng mô hình dữ liệu Raster.....	4
Hình 1.3: Bản đồ minh họa cấu trúc Topology	5
Hình 1.4: Minh họa dữ liệu Spaghetti	7
Hình 1.5: Các bước của quá trình khai phá dữ liệu	8
Hình 1.6: Cây quyết định	13
Hình 2.1: Phân lớp sử dụng thuộc tính "Quang cảnh"	25
Hình 2.2: Phân nhánh "Quang cảnh" = "Nắng"	25
Hình 2.3: Cây nhánh "Quang cảnh" = "Nắng"	26
Hình 2.4: Cây quyết định tính toán từ thuật toán ID3	26
Hình 2.5: Xác định giá trị phân chia kiểu số	32
Hình 2.6: Chỉ mục kết nối không gian	36
Hình 2.7: Các Layer dự báo cháy rừng	40
Hình 2.8: Layer mục tiêu và Layer phủ bề mặt và mật độ dân số	41
Hình 2.9: Cây quyết định không gian	43
Hình 2.10: Mô tả Object ID các Layer	45
Hình 2.11: Quan hệ không gian giữa Layer mục tiêu và các Layer mô tả	45
Hình 2.12: Thống kê Layer phủ bề mặt theo loại phủ bề mặt	47
Hình 2.13: Thống kê Layer mật độ dân số theo loại mật độ dân số	51
Hình 2.14: Thống kê Layer khoảng cách đến sông gần nhất	52
Hình 2.15: Phân lớp Layer phủ bề mặt theo loại phủ bề mặt	52
Hình 2.16: Nhánh Dryland forest - thống kê Layer mật độ dân số	53
Hình 2.17: Nhánh Dryland forest - thống kê khoảng cách đến sông gần nhất	53
Hình 3.1: Bản đồ các trạm BTS trên địa bàn tỉnh Vĩnh Phúc	56
Hình 3.2: Khoảng cách đến BTS gần nhất	57
Hình 3.3: Bản đồ BTS và các điểm mục tiêu	58
Hình 3.4: Mô tả cấu trúc dữ liệu trạm BTS	59
Hình 3.5: Mô tả cấu trúc dữ liệu bệnh viện, trường học, công sở	59
Hình 3.6: Mô tả cấu trúc dữ liệu vùng dân cư	60
Hình 3.7: Phần mềm ArcMap và ArcCatalog biên tập dữ liệu GIS	60
Hình 3.8: Mô tả kết quả chạy chương trình.....	61
Hình 3.9: File dữ liệu Excel biểu diễn bảng dữ liệu đầu vào	62
Hình 3.10: Kết quả trên phần mềm Weka	63
Hình 3.11: Thống kê Tuple tại các nhánh BTScover	64
Hình 3.12: Thống kê Tuple tại nhánh BTScover="Low" và xét Density	65
Hình 3.13: Biểu diễn kết quả dưới dạng cây quyết định	66

MỞ ĐẦU

1. Đặt vấn đề

Những tiến bộ trong công nghệ CSDL và kỹ thuật thu thập dữ liệu như đọc mã số mã vạch, viễn thám, ghi nhận thông tin từ các vệ tinh,... đã tạo ra một lượng lớn thông tin, dữ liệu. Việc dữ liệu tăng lên nhanh với quy mô lớn đòi hỏi phải được khai phá để trích chọn ra các tri thức hữu ích phục vụ cho công tác chuyên môn. Chính điều này đã dẫn đến sự ra đời của lĩnh vực khai phá dữ liệu hay khai phá tri thức trong các CSDL. Khai phá tri thức trong các CSDL có thể được định nghĩa là khai phá tri thức đáng quan tâm, tiềm ẩn và chưa biết trước trong các CSDL. Khai phá dữ liệu là sự kết hợp của một số lĩnh vực bao gồm học máy, các hệ thống CSDL, thể hiện dữ liệu, thống kê và lý thuyết thông tin.

Đã có nhiều nghiên cứu về khai phá dữ liệu trong các CSDL quan hệ và giao dịch, nhưng đối với các CSDL không gian vấn đề khai phá dữ liệu vẫn còn là những thách thức cần được giải quyết.

Dữ liệu không gian là dữ liệu liên quan đến các đối tượng trong không gian. Một CSDL không gian lưu trữ các đối tượng không gian bao gồm các kiểu dữ liệu không gian và các quan hệ không gian giữa các đối tượng. Dữ liệu không gian mang thông tin hình học và khoảng cách thường được tổ chức theo các cấu trúc chỉ mục không gian và truy cập bằng các phương pháp truy cập không gian. Chính các đặc trưng khác biệt này của các CSDL không gian đã đặt ra nhiều trở ngại nhưng cũng mang đến nhiều cơ hội cho khai phá tri thức từ CSDL không gian. Khai phá dữ liệu không gian hay khai phá tri thức trong CSDL không gian là trích chọn ra các tri thức tiềm ẩn, các quan hệ không gian hay các mẫu chưa rõ lưu trữ trong các CSDL không gian.

Các nghiên cứu trước đây về học máy, các hệ thống CSDL và thống kê đã đặt nền móng cho nghiên cứu khai phá tri thức trong các CSDL. Và những tiến bộ của các CSDL không gian như cấu trúc dữ liệu không gian, lập luận không gian, tính toán hình học,... đã mở đường cho khai phá dữ liệu không gian. Trở ngại lớn nhất trong khai phá dữ liệu không gian là hiệu quả của các thuật toán khai phá dữ liệu không gian do lượng dữ liệu không gian thường có quy mô lớn, các kiểu dữ liệu không gian và các phương pháp truy cập không gian phức tạp.

Các phương pháp khai phá dữ liệu không gian tập trung theo ba hướng chính là khai phá luật kết hợp không gian, phân lớp dữ liệu không gian và phân cụm dữ liệu không gian. Với mong muốn nghiên cứu về phân lớp dữ liệu không gian sử dụng cây quyết định, luận văn đi sâu tìm hiểu một lĩnh vực nhỏ đó là phân lớp dữ liệu không gian sử dụng cây quyết định.