

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

TRẦN HÀ PHƯƠNG

PHÂN CỤM ĐỒ THỊ DỮ LIỆU VÀ ỨNG DỤNG

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

THÁI NGUYÊN - 2016

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

TRẦN HÀ PHƯƠNG

PHÂN CỤM ĐỒ THỊ DỮ LIỆU VÀ ỨNG DỤNG

Chuyên ngành: Khoa học máy tính

Mã số: 60 48 01 01

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Người hướng dẫn khoa học: PGS. TS. ĐOÀN VĂN BAN

THÁI NGUYÊN - 2016

LỜI CAM ĐOAN

Tên tôi là: Trần Hà Phương

Sinh ngày:

Học viên lớp cao học CHK13 - Trường Đại học Công nghệ thông tin và Truyền thông - Đại học Thái Nguyên.

Hiện đang công tác tại:

Xin cam đoan: Đề tài “*Phân cụm đồ thị dữ liệu và ứng dụng*” do Thầy giáo PGS.TS Đoàn Văn Ban hướng dẫn là công trình nghiên cứu của riêng tôi. Tất cả tài liệu tham khảo đều có nguồn gốc, xuất xứ rõ ràng.

Tác giả xin cam đoan tất cả những nội dung trong luận văn đúng như nội dung trong đề cương và yêu cầu của thầy giáo hướng dẫn. Nếu sai tôi hoàn toàn chịu trách nhiệm trước hội đồng khoa học và trước pháp luật.

Thái Nguyên, ngày 14 tháng 4 năm 2016

Tác giả luận văn

Trần Hà Phương

LỜI CẢM ƠN

Sau một thời gian nghiên cứu và làm việc nghiêm túc, được sự động viên, giúp đỡ và hướng dẫn tận tình của Thầy giáo hướng dẫn PGS.TS Đoàn Văn Ban, luận văn với đề tài “*Phân cụm đồ thị dữ liệu và ứng dụng*” đã hoàn thành.

Tôi xin bày tỏ lòng biết ơn sâu sắc đến:

Thầy giáo hướng dẫn **PGS.TS Đoàn Văn Ban** đã tận tình chỉ dẫn, giúp đỡ tôi hoàn thành luận văn này.

Khoa sau Đại học Trường Đại học công nghệ thông tin và truyền thông đã giúp đỡ tôi trong quá trình học tập cũng như thực hiện luận văn.

Tôi xin chân thành cảm ơn bạn bè, đồng nghiệp và gia đình đã động viên, khích lệ, tạo điều kiện giúp đỡ tôi trong suốt quá trình học tập, thực hiện và hoàn thành luận văn này.

Thái Nguyên, ngày 16 tháng 6 năm 2016

Tác giả luận văn

Trần Hà Phương

MỤC LỤC

LỜI CAM ĐOAN	i
LỜI CẢM ƠN	ii
MỤC LỤC	iii
DANH MỤC CÁC TỪ VIẾT TẮT	v
DANH MỤC BẢNG	vi
DANH MỤC CÁC HÌNH ẢNH	vii
MỞ ĐẦU	1
1. Tính khoa học và cấp thiết của đề tài	1
2. Mục tiêu, đối tượng và phạm vi nghiên cứu của đề tài	2
3. Phương pháp luận nghiên cứu	2
4. Nội dung và bố cục của luận văn	2
CHƯƠNG 1 TỔNG QUAN VỀ PHÂN CỤM DỮ LIỆU	4
1.1 Khái niệm, mục tiêu và các bước cơ bản của phân cụm dữ liệu	4
1.1.1 Phân cụm dữ liệu là gì?	4
1.1.2 Các mục tiêu của phân cụm dữ liệu	5
1.1.3 Các bước cơ bản để phân cụm	7
1.2 Một số khái niệm cần thiết khi tiếp cận phân cụm dữ liệu	8
1.2.1 Phân loại các kiểu dữ liệu	8
1.2.2 Độ đo tương tự và phi tương tự	9
1.3 Những kỹ thuật tiếp cận trong phân cụm dữ liệu	11
1.3.1 Phương pháp phân cụm phân hoạch	12
1.3.2 Phương pháp phân cụm phân cấp	12
1.3.3 Phương pháp phân cụm dựa trên mật độ	13
1.3.4 Phương pháp phân cụm dựa trên lưới	14
1.3.5 Phương pháp phân cụm dựa trên mô hình	15
1.3.6 Phương pháp phân cụm dữ liệu có liên kết	15
1.4 Các ứng dụng của phân cụm dữ liệu	16
1.5 Các yêu cầu và những vấn đề còn tồn tại trong phân cụm dữ liệu	18
1.5.1 Các yêu cầu của phân cụm dữ liệu	18
1.5.2 Những vấn đề còn tồn tại trong phân cụm dữ liệu	19
1.6 Tổng kết chương	20

CHƯƠNG 2 THUẬT TOÁN PHÂN CỤM ĐỒ THỊ DỮ LIỆU	22
2.1 Tổng quan về lý thuyết đồ thị	22
2.1.1 Giới thiệu chung	22
2.1.2 Biểu diễn đồ thị trên máy tính	23
2.2 Mô hình đồ thị dữ liệu.....	27
2.3 Độ đo trong phân cụm đồ thị dữ liệu	28
2.3.1 Độ đo cho phân cụm dữ liệu nói chung.....	28
2.3.2 Độ đo cho phân cụm đồ thị.....	30
2.4 Một số thuật toán phân cụm dữ liệu dựa trên đồ thị	31
2.4.1 Thuật toán CHAMELEON.....	31
2.4.2 Thuật toán phân cụm quang phổ.....	33
2.4.3 Thuật toán phân cụm phân cấp	35
2.5 Kết luận chương	46
CHƯƠNG 3. ỨNG DỤNG THUẬT TOÁN ĐỒ THỊ QUANG PHỔ TRONG VIỆC PHÂN LOẠI KẾT QUẢ HỌC TẬP CỦA HỌC SINH	47
3.1 Đặt vấn đề.....	47
3.2 Xây dựng chương trình ứng dụng	49
3.2.1 Mục đích chương trình	49
3.2.2 Cơ sở dữ liệu.....	49
3.2.3 Các bước thực hiện	49
3.2.4 Môi trường cài đặt	50
3.2.5 Cài đặt.....	50
3.3 Các chức năng chính của chương trình	51
3.3.1 Chương trình chính.....	51
3.3.2 Biểu diễn dữ liệu theo đồ thị.....	52
3.3.3 Phân cụm dữ liệu đồ thị quang phổ	52
3.4 Đánh giá hiệu quả của thuật toán phân cụm dữ liệu đồ thị quang phổ	54
3.5 Kết luận chương	58
KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN	59
TÀI LIỆU THAM KHẢO	61

DANH MỤC CÁC TỪ VIẾT TẮT

Từ hoặc cụm từ	Từ tiếng Anh	Từ tiếng Việt
BDCM	Binding data Clustering Methods	Phương pháp phân cụm dữ liệu có liên kết
CA	Continuous Attribute	Thuộc tính liên tục
CSDL		Cơ sở dữ liệu
DA	Discrete Attribute	Thuộc tính rời rạc
DBM	Density-Based Methods	Phương pháp dựa trên mật độ
GBM	Grid-Based Methods	Phương pháp dựa trên lưới
HM	Hierarchical Methods	Phương pháp phân cấp
MBCM	Model-Based Clustering Methods	Phương pháp dựa trên mô hình phân cụm
MC	Markov Clustering	Phân cụm theo mô hình Markov
MST	Minimum Spanning Tree	Cây khung nhỏ nhất
PM	Partitioning Methods	Phương pháp phân hoạch
RWA	Random Walk Algorithm	Thuật toán bước đi ngẫu nhiên
SC	Star Clustering	Phân cụm hình sao
SCA	Spectral Clustering Algorithm	Thuật toán phân cụm quang phổ
SOM	Self-Organizing Map	Mạng tự tổ chức

DANH MỤC BẢNG

Bảng 3.1. Các module chính của chương trình.....	51
--	----

DANH MỤC CÁC HÌNH ẢNH

Hình 1.1. Ví dụ về phân cụm dữ liệu [7]	5
Hình 1.2. Ví dụ phân cụm các đối tượng dựa trên khoảng cách [7]	5
Hình 1.3. Ví dụ phân cụm các đối tượng dựa trên kích cỡ [7]	6
Hình 1.4. Các bước trong quá trình phân cụm	8
Hình 1.5. Các chiến lược phân cụm phân cấp [11]	13
Hình 1.6. Cấu trúc phân cụm dữ liệu dựa trên lưới	14
Hình 2.1. Ví dụ về mô hình đồ thị.....	22
Hình 2.2. Phân loại đồ thị.....	23
Hình 2.3. Ma trận kề vô hướng (trên) và có hướng (dưới)	25
Hình 2.4. Ma trận trọng số vô hướng (trên) và có hướng (dưới)	26
Hình 2.5. Ma trận liên thuộc vô hướng (trên) và có hướng (dưới)	27
Hình 2.6. Minh họa thuật toán CHAMELEON	32
Hình 2.7. Nguyên lý chung của AntTree	36
Hình 2.8. Kiến trúc khác nhau giữa SOM và SOMTree	41
Hình 2.9. Phân việc từ cây $tree_{c_{old}}$ cho treec	44
Hình 2.10. Tách subtrees khỏi cây $tree_{c_{old}}$ và đưa vào list	45
Hình 2.11. Tái liên kết subtrees vào treec	45
Hình 3.1. Màn hình chính của chương trình	51
Hình 3.2. Biểu diễn dữ liệu theo đồ thị	52
Hình 3.3. Phân cụm dữ liệu đồ thị quang phổ với dữ liệu vào là dữ liệu kiểm tra ...	53
Hình 3.4. Phân cụm dữ liệu đồ thị quang phổ với dữ liệu vào là điểm học sinh	54
Hình 3.5. Kết quả phân cụm dữ liệu dạng ba cụm Gaussian với 1000 mẫu dữ liệu.	55
Hình 3.6. Kết quả phân cụm dữ liệu dạng ba cụm Gaussian với độ lớn lần lượt là 100, 1000, 3000 mẫu dữ liệu.....	55
Hình 3.7. Kết quả phân cụm dữ liệu dạng hai nửa vàng trắng với kích thước dữ liệu là ba cụm Gaussian với độ lớn lần lượt là 7500 mẫu dữ liệu.....	56
Hình 3.8. Kết quả phân cụm dữ liệu dạng hai nửa vàng trắng với hai thuật toán K mean (trái) và đồ thị quang phổ (phải).....	56
Hình 3.9. Kết quả phân cụm dữ liệu điểm học sinh với số cụm khác nhau.....	57

MỞ ĐẦU

1. Tính khoa học và cấp thiết của đề tài

Phân cụm là một trong những vấn đề cơ bản phổ biến trong các lĩnh vực nhận dạng mẫu, học máy và khai thác dữ liệu. Hiện tại, trên thực tế có rất nhiều thuật toán phân cụm được công bố. Tuy nhiên, do không tồn tại một thuật toán phân cụm duy nhất cho tất cả các loại bộ dữ liệu, những thuật toán phân cụm mới vẫn liên tục được đề xuất. Kết quả là, người dùng phải chọn thuật toán thích hợp nhất từ nhiều ứng viên để đạt được kết quả chính xác.

Trong thực tế, việc lựa chọn thuật toán phân cụm dữ liệu phù hợp là rất khó khăn do người sử dụng thường không có một kiến thức tiên nghiệm về sự đa dạng và phức tạp của dữ liệu. Để phần nào giảm bớt nhược điểm trên, các thuật toán phân cụm dựa trên đồ thị được đề xuất do ưu điểm ở khả năng xử lý các bộ dữ liệu đa dạng và có cấu trúc. Bản chất của các thuật toán này là biểu diễn dữ liệu dựa trên đồ thị và phân cụm các thành phần theo các thuật toán thiết kế riêng.

Đồ thị là những cấu trúc toán học được sử dụng để đại diện cho mối quan hệ giữa cặp đối tượng từ một tập hợp xác định. Đồ thị chứa đỉnh (đại diện cho các đối tượng) và các cạnh nối các đỉnh (đại diện cho mối quan hệ giữa các đối tượng cặp). Đây là phương pháp cấu trúc dữ liệu quan trọng được sử dụng trong rất nhiều lĩnh vực như khai thác dữ liệu, xử lý ngôn ngữ tự nhiên, tìm kiếm thông tin và khai thác thông tin. Trong phân cụm, sự tương đồng giữa các đối tượng được phân cụm có thể được diễn tả như một đồ thị có trọng số. Trong đó, các đối tượng là các đỉnh và sự tương đồng là trọng số của các cạnh. Bài toán phân cụm sẽ được đơn giản hóa về bài toán phân cụm đồ thị mà nhiệm vụ chính là tách các đồ thị phụ dày đặc và kết nối thưa thớt khỏi nhau dựa trên khái niệm của mật độ nội cụm so với khoảng cách liên cụm.

Với những lý do trên, tác giả đã chọn đề tài ***“Phân cụm đồ thị dữ liệu và ứng dụng”*** làm đề tài nghiên cứu luận văn tốt nghiệp thạc sĩ chuyên ngành Khoa học máy tính.