

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

**TRỊNH ANH TUẤN**

**PHÂN CỤM DỮ LIỆU VÀ ỨNG DỤNG PHÂN CỤM  
SINH HỌC TRƯỜNG TRUNG HỌC CƠ SỞ CHU VĂN  
AN**

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

**Người hướng dẫn: TS.Nguyễn Long Giang**

**THÁI NGUYÊN - 2016**

## **LỜI CAM ĐOAN**

Tác giả Trịnh Anh Tuấn xin cam kết rằng nội dung của Luận văn này chưa được nộp cho bất kỳ một chương trình cấp bằng cao học nào cũng như bất kỳ một chương trình đào tạo cấp bằng nào khác.

Ngoài ra, tác giả cũng xin cam kết Luận văn thạc sĩ này là nỗ lực riêng của cá nhân tác giả. Các kết quả, phân tích, kết luận trong Luận văn thạc sĩ này (ngoài các phần được trích dẫn) đều là kết quả làm việc của cá nhân tác giả.

*Thái Nguyên, tháng 6 năm 2016*

**Tác Giả**

**Trịnh Anh Tuấn**

## LỜI CẢM ƠN

Để hoàn thành được luận văn này, trước hết tôi xin gửi lời cảm ơn sâu sắc nhất tới TS. Nguyễn Long Giang, Viện Công nghệ thông tin - Viện Hàn lâm Khoa học và Công nghệ Việt Nam đã tận tình hướng dẫn, chỉ bảo, định hướng, đóng góp những ý kiến quý báu trong suốt quá trình thực hiện luận văn.

Tôi xin chân thành cảm ơn các thầy, cô giáo trong Bộ môn Khoa học máy tính, Khoa Công nghệ thông tin, Phòng Đào tạo Sau đại học - Nghiên cứu Khoa học, Trường Đại học Công nghệ thông tin và truyền thông Thái Nguyên đã tạo mọi điều kiện tốt nhất để tôi hoàn thành khóa học.

Xin cảm ơn đồng nghiệp tại trường Trung học cơ sở Chu Văn An thành phố Thái Nguyên đã trợ giúp rất nhiều trong thời gian qua.

Trong quá trình thực hiện Luận văn, mặc dù đã cố gắng hết mình, song chắc chắn luận văn của em vẫn còn nhiều thiếu sót. Em rất mong nhận được sự chỉ bảo vào đóng góp tận tình của các thầy cô để luận văn của em được hoàn thiện hơn.

*Thái Nguyên, tháng 6 năm 2016*

**Tác Giả**

**Trịnh Anh Tuấn**

## MỤC LỤC

<b>LỜI CAM ĐOAN</b> .....	i
<b>LỜI CẢM ƠN</b> .....	iii
<b>MỤC LỤC</b> .....	iv
<b>DANH MỤC CÁC BẢNG</b> .....	vii
<b>DANH MỤC CÁC HÌNH</b> .....	viii
<b>MỞ ĐẦU</b> .....	1
1.1. Sự cần thiết lựa chọn đề tài .....	1
1.2. Mục tiêu đề tài .....	2
1.3. Đối tượng và phạm vi nghiên cứu .....	2
1.4. Phương pháp nghiên cứu.....	3
1.5. Cấu trúc của luận văn .....	3
<b>Chương 1. TỔNG QUAN</b> .....	4
1.1. Quá trình khám phá tri thức .....	4
1.2. Khai phá dữ liệu .....	5
1.2.1. <i>Khái niệm khai phá dữ liệu</i> .....	5
1.2.2. <i>Các kỹ thuật khai phá dữ liệu</i> .....	6
1.3. Phân cụm dữ liệu.....	8
1.3.1. <i>Khái niệm về phân cụm dữ liệu</i> .....	8
1.3.2. <i>Một số vấn đề trong phân cụm dữ liệu</i> .....	9
1.3.3. <i>Mục tiêu của phân cụm dữ liệu</i> .....	10
1.3.4. <i>Các bước cơ bản trong phân cụm dữ liệu</i> .....	10
1.3.5. <i>Yêu cầu của phân cụm dữ liệu</i> .....	11
1.3.6. <i>Ứng dụng của phân cụm dữ liệu</i> .....	12
1.4. Kết luận chương .....	13
<b>Chương 2. CÁC PHƯƠNG PHÁP PHÂN CỤM DỮ LIỆU</b> .....	14
2.1. Kiểu dữ liệu .....	14
2.1.1. <i>Phân loại kiểu dữ liệu dựa trên kích thước miền</i> .....	14
2.1.2. <i>Phân loại kiểu dữ liệu dựa trên hệ đo</i> .....	14
2.2. Phép đo độ tương tự và phép đo khoảng cách .....	16

2.2.1.	<i>Khái niệm tương tự và không tương tự</i> .....	16
2.2.2.	<i>Phép đo khoảng cách</i> .....	17
2.3.	<i>Phương pháp phân cụm phân hoạch</i> .....	18
2.3.1.	<i>Giới thiệu phương pháp</i> .....	18
2.3.2.	<i>Thuật toán K-MEANS</i> .....	19
2.3.3.	<i>Thuật toán PAM</i> .....	21
2.4.	<i>Phương pháp phân cụm phân cấp</i> .....	24
2.4.1.	<i>Giới thiệu phương pháp</i> .....	24
2.4.2.	<i>Thuật toán HERACHICAL</i> .....	25
2.4.3.	<i>Thuật toán BIRCH</i> .....	28
2.5.	<i>Phương pháp phân dựa trên mật độ</i> .....	31
2.5.1.	<i>Giới thiệu phương pháp</i> .....	31
2.5.2.	<i>Thuật toán DBSCAN</i> .....	32
2.6.	<i>Phương pháp phân cụm dựa trên lưới</i> .....	36
2.6.1.	<i>Giới thiệu phương pháp</i> .....	36
2.6.2.	<i>Thuật toán STING</i> .....	37
2.7.	<i>Kết luận chương</i> .....	40
<b>Chương 3. PHÂN CỤM KẾT QUẢ HỌC TẬP TẠI TRƯỜNG TRUNG HỌC CƠ SỞ CHU VĂN AN</b> .....		41
3.1.	<i>Bài toán phân cụm kết quả học tập của học sinh tại trường trung học cơ sở Chu Văn An</i> .....	41
3.1.1.	<i>Giới thiệu trường Trung học cơ sở Chu Văn An</i> .....	41
3.1.2.	<i>Bảng dữ liệu kết quả học tập của học sinh</i> .....	42
3.1.3.	<i>Bài toán phân cụm kết quả học tập của học sinh</i> .....	43
3.2.	<i>Lựa chọn phương pháp, công cụ</i> .....	44
3.2.1.	<i>Lựa chọn ngôn ngữ R thực hiện phân cụm</i> .....	44
3.2.2.	<i>Các bước thực hiện phân cụm bằng ngôn ngữ R</i> .....	46
3.3.	<i>Kết quả phân cụm bằng thuật toán K-means</i> .....	48
3.3.1.	<i>Phân cụm học sinh dựa trên kết quả học tập</i> .....	48
3.3.2.	<i>Phân cụm học sinh dựa trên điểm trung bình các môn</i> .....	52
3.3.3.	<i>Phân cụm dựa trên điểm trung bình môn toán và môn văn</i> .....	53

3.4. Kết luận chương .....	54
KẾT LUẬN.....	55
TÀI LIỆU THAM KHẢO .....	57

**DANH MỤC CÁC BẢNG**

<b>Bảng 3.1.</b> Bảng dữ liệu kết quả học tập của học sinh.....	43
<b>Bảng 3.2.</b> Phân cụm theo kết quả học tập .....	48
<b>Bảng 3.3.</b> Thống kê phân cụm theo địa bàn hành chính .....	49
<b>Bảng 3.4.</b> Thống kê phân cụm theo hoàn cảnh gia đình .....	50
<b>Bảng 3.5.</b> Thống kê phân cụm theo dân tộc .....	51
<b>Bảng 3.6.</b> Thống kê phân cụm theo giới tính .....	52
<b>Bảng 3.7.</b> Phân cụm theo điểm trung bình môn toán .....	52
<b>Bảng 3.8.</b> Phân cụm theo điểm trung bình môn văn .....	53

## DANH MỤC CÁC HÌNH

Hình 1.1. Quá trình khám phá tri thức .....	4
Hình 1.2. Quy trình phân cụm .....	8
Hình 2.1. Khởi tạo các đối tượng medoid.....	22
Hình 2.2. Cây CF được dùng trong thuật toán BIRCH .....	29
Hình 2.3. Ý tưởng của thuật toán phân cụm phân cấp.....	31
Hình 2.4. Lân cận với ngưỡng $\epsilon$ của điểm $p$ .....	32
Hình 2.5. Mật độ liên lạc .....	33
Hình 2.6. Mật độ liên thông.....	34
Hình 2.7. Các mức ô lưới khác nhau trong quá trình truy vấn .....	38
Hình 3. 1. Website của trường Trung học cơ sở Chu Văn An.....	41
Hình 3.2. Cơ cấu tổ chức của trường Trung học cơ sở Chu Văn An.....	42
Hình 3.3. Thống kê số học sinh theo điểm toán .....	47
Hình 3.4. Kết quả phân cụm .....	48



## MỞ ĐẦU

### 1.1. Sự cần thiết lựa chọn đề tài

Sự phát triển của nhanh chóng các ứng dụng công nghệ thông tin và Internet vào nhiều lĩnh vực đời sống xã hội, quản lý kinh tế, khoa học kỹ thuật... trong mấy năm gần đây đã tạo ra nhiều cơ sở dữ liệu khổng lồ. Để khai thác hiệu quả nguồn thông tin từ các cơ sở dữ liệu khổng lồ đó nhằm mục đích dự báo, hỗ trợ ra quyết định, bên cạnh các phương pháp khai thác thông tin truyền thống, các nhà nghiên cứu đã phát triển các phương pháp, kỹ thuật và phần mềm mới hỗ trợ tiến trình khám phá, phân tích, tổng hợp thông tin, lĩnh vực này được gọi là khai phá dữ liệu và khám phá tri thức (Data mining and Knowledge discovery)

Khai phá dữ liệu và khám phá tri thức là một lĩnh vực quan trọng của ngành Công nghệ thông tin với mục tiêu là tìm kiếm các tri thức có ích, cần thiết, tiềm ẩn và chưa được biết trước trong cơ sở dữ liệu lớn. Đây là lĩnh vực đã và đang thu hút đông đảo các nhà khoa học trên thế giới và trong nước tham gia nghiên cứu. Khai phá dữ liệu có thể xem là nhiệm vụ quan trọng trong quá trình khám phá tri thức từ cơ sở dữ liệu, bao gồm ba bước chính: thu thập và tiền xử lý dữ liệu; lựa chọn các thuật toán khai phá dữ liệu; đánh giá kết quả và biểu diễn tri thức. Các bài toán quan trọng trong khai phá dữ liệu bao gồm: phân lớp (classification); hồi quy (regression); phân cụm (clustering); khai phá luật kết hợp (rule association)... Các kỹ thuật, công cụ sử dụng trong khai phá dữ liệu bao gồm: cây quyết định; mạng nơron nhân tạo; thuật toán di truyền; các kỹ thuật phân lớp, phân cụm; các phương pháp trong thống kê như phân tích tương quan, phân tích phương sai, hồi quy đơn biến, đa biến... Khai phá dữ liệu có nhiều ứng dụng trong các lĩnh vực khác nhau của đời sống như: phân tích, dự báo trong kinh tế, tài chính; chuẩn đoán bệnh trong y tế; tin sinh học; hỗ trợ quá trình sản xuất, kinh doanh...

Phân cụm (clustering) là bài toán có vai trò quan trọng trong khai phá dữ liệu và có nhiều ứng dụng trong thực tiễn. Mục tiêu của phương pháp phân cụm dữ liệu là quá trình nhóm các đối tượng tương tự nhau trong cơ sở dữ liệu vào các cụm sao cho các đối tượng trong cùng một cụm là tương đồng, còn các đối tượng thuộc các

cụm khác nhau sẽ không tương đồng. Điểm mạnh của phân cụm dữ liệu là đưa ra được những cấu trúc có ích hoặc những cụm các đối tượng tìm thấy trực tiếp từ dữ liệu mà không cần bất kì một tri thức cơ sở nào. Giống như cách tiếp cận học máy, phân cụm dữ liệu được hiểu như là phương pháp học không có thầy (unsupervised learning). Không giống như phân lớp dữ liệu, phân cụm dữ liệu không đòi hỏi phải định nghĩa trước các mẫu dữ liệu huấn luyện. Vì thế, có thể coi phân cụm dữ liệu là một cách học bằng quan sát (learning by observation), trong khi phân lớp dữ liệu là học bằng ví dụ (learning by example). Trong phương pháp này sẽ không thể biết kết quả các cụm thu được sẽ như thế nào khi bắt đầu quá trình. Vì vậy, cần có một chuyên gia để đánh giá các cụm thu được. Phân cụm dữ liệu được sử dụng nhiều trong các ứng dụng về phân đoạn thị trường, phân đoạn khách hàng, nhận dạng mẫu, phân loại trang Web, phân loại, đánh giá học sinh, sinh viên trong các trường học... Ngoài ra, phân cụm dữ liệu còn có thể được sử dụng như một bước tiền xử lí cho các thuật toán khai phá dữ liệu khác.

Là một người công tác trong lĩnh vực giáo dục phổ thông, với mong muốn áp dụng các kiến thức đã học về các phương pháp phân cụm vào bài toán thực tiễn là phân cụm học sinh của trường Trung học cơ sở Chu Văn An, thành phố Thái Nguyên dựa vào kết quả học tập, tác giả luận văn chọn đề tài: “*Phân cụm dữ liệu và ứng dụng phân cụm học sinh trường Trung học cơ sở Chu Văn An*”.

## **1.2. Mục tiêu đề tài**

Nắm bắt được một cách tổng thể các phương pháp phân cụm trong khai phá dữ liệu. Trên cơ sở đó, áp dụng các kỹ thuật phân cụm vào giải quyết bài toán thực tiễn tại địa phương nơi tác giả làm việc là phân cụm kết quả học tập của học sinh trường Trung học cơ sở Chu Văn An, thành phố Thái Nguyên.

## **1.3. Đối tượng và phạm vi nghiên cứu**

Đối tượng nghiên cứu của đề tài là các phương pháp phân cụm dữ liệu trong khai phá dữ liệu và cơ sở dữ liệu về kết quả học tập của học sinh trường Trung học cơ sở Chu Văn An, thành phố Thái Nguyên.