

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

---

**TẠ DUY KHÁNH**

**PHƯƠNG PHÁP LAN TRUYỀN ĐỘ TƯƠNG TỰ TRONG  
PHÂN CỤM DỮ LIỆU VÀ ỨNG DỤNG**

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

**Thái Nguyên - 2015**

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

---

**TẠ DUY KHÁNH**

**PHƯƠNG PHÁP LAN TRUYỀN ĐỘ TƯƠNG TỰ TRONG  
PHÂN CỤM 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 NGUYỄN ĐÌNH HÓA**

**Thái Nguyên - 2015**

## LỜI CẢM ƠN

Đầu tiên, tôi xin gửi lời cảm ơn sâu sắc đến thầy giáo, cán bộ hướng dẫn khoa học PGS.TS Nguyễn Đình Hóa đã tận tình hướng dẫn tôi từ những buổi đầu tiên khi tiếp cận với đề tài khoa học. Trong quá trình làm luận văn, tôi cũng nhận được sự giúp đỡ rất nhiệt tình từ nhóm nghiên cứu của TS. Lê Hoàng Sơn tại Trung tâm tính toán hiệu năng cao, trường ĐH KHTN và đề tài NCKH cấp ĐHQG, mã số GG.14.60.

Tôi xin bày tỏ lòng biết ơn đến các thầy cô giáo ở trường Đại học Công nghệ thông tin và Truyền thông – Đại học Thái Nguyên, các cán bộ Trung tâm Đông Đô - Hà Nội, đã tận tình giảng dạy và tạo mọi điều kiện cho tôi học tập, nghiên cứu và hoàn thành luận văn này.

Tôi xin chân thành cảm ơn các bạn học viên lớp CK12H, CK13H – Khoa học máy tính đã giúp đỡ, tạo điều kiện cho tôi trong suốt quá trình học tập và thực hiện luận văn.

Cuối cùng, tôi xin gửi lời cảm ơn sâu sắc nhất đến gia đình, đồng nghiệp và bạn bè tôi, những người đã động viên, tạo mọi điều kiện cho tôi lao động và học tập trong suốt thời gian qua.

Tôi xin cam đoan luận văn là công trình nghiên cứu của riêng cá nhân tôi, không sao chép của ai. Luận văn là do tôi tự nghiên cứu, đọc, dịch tài liệu, tổng hợp và thực hiện. Nội dung lý thuyết trong luận văn có sử dụng một số tài liệu tham khảo như đã trình bày trong phần tài liệu tham khảo. Chương trình phần mềm và những kết quả trong luận văn là trung thực và chưa được công bố trong bất kỳ một hệ thống nào khác.

*Một lần nữa, xin chân thành cảm ơn!*

## MỤC LỤC

LỜI CẢM ƠN.....	i
MỤC LỤC.....	ii
DANH MỤC CÁC TỪ VIẾT TẮT.....	iv
DANH MỤC CÁC HÌNH VẼ.....	v
MỞ ĐẦU .....	1
<b>CHƯƠNG 1: HỆ THỐNG THÔNG TIN ĐỊA LÝ VÀ PHÂN CỤM DỮ LIỆU ĐỊA LÝ .....</b>	<b>6</b>
<b>1.1 Tổng quan về hệ thống thông tin địa lý.....</b>	<b>6</b>
1.1.1 Lịch sử ra đời.....	6
1.1.2 Định nghĩa.....	7
1.1.3 Các thành phần của hệ thống thông tin địa lý .....	8
1.1.4 Dữ liệu trong hệ thống thông tin địa lý .....	10
<b>1.2 Phân cụm dữ liệu địa lý.....</b>	<b>11</b>
1.2.1 Phân cụm dữ liệu .....	11
1.2.2 Một số kỹ thuật phân cụm dữ liệu .....	12
1.2.2.1 Thuật toán phân cụm theo cây phân cấp.....	13
1.2.2.2 Thuật toán phân cụm phân hoạch : Phân cụm k-means .....	14
1.2.2.3 Phân cụm mờ .....	16
<b>1.3 Dữ liệu địa lý và vấn đề phân cụm đối tượng địa lý.....</b>	<b>17</b>
1.3.1 Cấu trúc dữ liệu trong GIS.....	18
1.3.1.1 Hai mô hình dữ liệu không gian .....	18
1.3.1.2 Dữ liệu thuộc tính .....	19
1.3.2 Các vấn đề trong phân cụm dữ liệu địa lý .....	20
<b>CHƯƠNG 2: PHÂN CỤM BẰNG THUẬT TOÁN LAN TRUYỀN ĐỘ TƯƠNG TỰ .....</b>	<b>22</b>
<b>2.1 Các khái niệm cơ sở.....</b>	<b>22</b>
2.1.1 Ý tưởng thuật toán .....	22
<b>2.1.2 Các công thức chính trong thuật toán AP .....</b>	<b>24</b>
2.1.3 Thuật toán AP nguyên thủy .....	25
<b>2.2 Thuật toán lan truyền AP tự thích nghi.....</b>	<b>27</b>
2.2.1 Phương pháp thích ứng giảm dần .....	28
2.2.2 Kỹ thuật thích nghi p-scanning.....	30

<b>2.3 Thuật toán lan truyền AP với tập dữ liệu hỗn hợp kiểu số và kiểu phân loại .....</b>	<b>31</b>
2.3.1 Khoảng cách và ý nghĩa.....	32
2.3.2 Phương pháp .....	32
2.3.3 Cải thiện độ đo tương tự .....	34
2.3.4 Thích nghi thuật toán lan truyền.....	36
<b>CHƯƠNG 3: XÂY DỰNG ỨNG DỤNG PHÂN CỤM DỮ LIỆU ĐỊA LÝ</b>	<b>39</b>
<b>3.1 Bài toán thực tế và cách tiếp cận phân cụm dữ liệu.....</b>	<b>39</b>
3.1.1 Bài toán khai thác các dữ liệu quan trắc khí tượng .....	39
3.1.2 Lựa chọn giải pháp kỹ thuật công nghệ.....	40
<b>3.2 Các phần mềm GIS .....</b>	<b>40</b>
<b>3.3 Tìm hiểu về phần mềm mã nguồn mở MapWindow .....</b>	<b>42</b>
<b>3.4 Thiết kế một plug-in trên phần mềm mã nguồn mở Mapwindow ...</b>	<b>44</b>
3.4.1 Thêm một plug-ins từ Visual Studio vào MapWindow .....	44
3.4.2 Xây dựng ứng dụng với Active X map control trong Visual Studio..	45
<b>Kết quả chạy thử nghiệm.....</b>	<b>53</b>
<b>KẾT LUẬN .....</b>	<b>56</b>
1. Một số kết quả đạt được của luận văn .....	56
2. Những hạn chế và hướng phát triển.....	56
<b>TÀI LIỆU THAM KHẢO .....</b>	<b>57</b>

## DANH MỤC CÁC TỪ VIẾT TẮT

STT	Từ viết tắt	Từ tiếng Anh	Ý nghĩa
1	GIS	Geographical Information System	Hệ thống thông tin địa lý
2	AP	affinity propagation	Thuật toán lan truyền độ tương tự
3	CSDL	Database	Cơ sở dữ liệu
4	SIL	Silhouette	Công thức Silhouette
5	DEM	Digital Elevation Model	Mô hình kỹ thuật số độ cao
6	DTM	Digital Terrain Model	Mô hình kỹ thuật số các địa hình
7	TIN	Triangulated Irregular Network	Lưới tam giác không đều
8	SQL	Structured Query Language	Ngôn ngữ truy vấn có cấu trúc

## DANH MỤC CÁC HÌNH VẼ

Hình 1.1: Các thành phần của hệ thống thông tin địa lý Gis

Hình 1.2: Thuật toán phân cụm K-means

Hình 1.3: Cấu trúc vector và raster

Hình 2.1: Đồ thị Affinity Propagation (AP)

Hình 2.2: Minh họa hiệu năng của ba kỹ thuật rời rạc hóa khác nhau

Hình 3.1: Phần mềm mã nguồn mở Mapwindow

Hình 3.2: Kiểm tra plug-ins vừa add trong MapWindow

Hình 3.3: Kéo thả Map Control vào form

Hình 3.4: Kéo Legend vào form

Hình 3.5: Kéo thêm DataGridView vào form Table

Hình 3.6: Giao diện plugin APCluster

## MỞ ĐẦU

### 1. Đặt vấn đề

Nguồn dữ liệu dồi dào cung cấp nhiều thông tin, từ đó nhân loại đúc rút thành tri thức trong quá trình phát triển xã hội loài người. Với sự phát triển của công nghệ điện toán và hệ thống lưu trữ dữ liệu thì khối lượng tài nguyên số ngày càng trở nên phong phú và đồ sộ. Trong xã hội hiện đại, thông tin đóng một vai trò then chốt. Nhu cầu xử lý dữ liệu, trích rút thông tin, kịp thời khai thác chúng để mang lại những hiệu quả thiết thực cho công tác quản lý, hoạt động sản xuất kinh doanh,... ngày càng trở nên cấp thiết.

Khai phá dữ liệu nói chung để trích rút thông tin và phân cụm dữ liệu nói riêng là một trong những trọng tâm nghiên cứu của khoa học máy tính. Phân cụm dữ liệu là một trong những biện pháp để tìm kiếm tri thức, khi ta chưa biết nhiều thông tin về miền ứng dụng. Phân cụm được coi như một công cụ độc lập để xem xét phân bố dữ liệu, là bước tiền xử lý cho các bước sau. Phân cụm dữ liệu hiện có nhiều ứng dụng trong hầu hết các lĩnh vực hoạt động kinh tế, xã hội. Có nhiều phương pháp và thuật toán phân cụm dữ liệu khác nhau, tùy theo cách tiếp cận bài toán dưới góc độ nào. Một phương pháp mới được đề xuất tương đối gần đây là *Phương pháp lan truyền độ tương tự*.

Thuật toán lan truyền độ tương tự (*Affinity Propagation* - AP) là thuật toán phân cụm dữ liệu được đưa ra bởi Frey & Dueck vào năm 2007 dựa trên ý tưởng thuật toán lan truyền độ tin cậy trong suy diễn trên mạng xác suất Bayes, dựa trên cơ sở toán học của lý thuyết xác suất. Thuật toán lan truyền làm việc dựa trên sự tương đồng (*affinity* nghĩa là sự giống nhau, sự tương thích, sự hấp dẫn) giữa các cặp điểm dữ liệu và đồng thời xem xét tất cả các điểm dữ liệu như các tâm cụm tiềm năng, theo thuật ngữ ở đây là tất cả các điểm dữ liệu đều là hình mẫu (*exemplar*) tiềm năng, và trao đổi các thông điệp giá trị thực cho đến khi có được tập hình mẫu tốt (phân cụm tương ứng).



Thuật toán phân cụm AP có một số ưu điểm: cho kết quả phân cụm tốt, đặc biệt là trong trường hợp có số lượng lớn các cụm, phát hiện cụm có hình dáng bất kỳ, không yêu cầu phải xác định trước số cụm. Nó cũng cho phép dễ dàng thực hiện phân cụm thỏa mãn một số điều kiện xác định trước nào đó, tức là phân cụm bán giám sát. Đặc tính này thích hợp cho phân cụm dữ liệu trong GIS vì những ràng buộc điều kiện địa hình tự nhiên hoặc quản lý hành chính cần tính đến trong các bài toán thực tế.

Luận văn chọn đề tài “*Phương pháp lan truyền độ tương tự trong phân cụm dữ liệu và ứng dụng*” là hướng nghiên cứu chính, với mục tiêu khám phá những điểm mạnh, điểm yếu của phương pháp này, hiểu biết sâu thêm về một cách tiếp cận phân cụm, đồng thời nâng cao kỹ năng thực hành triển khai ứng dụng.

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

Đối tượng nghiên cứu là các phương pháp phân cụm dữ liệu, tập trung vào thuật toán lan truyền độ tương tự.

Phạm vi nghiên cứu là các điểm mạnh, điểm yếu và tiềm năng ứng dụng phương pháp lan truyền độ tương tự trong phân cụm dữ liệu địa lý.

## **3. Hướng nghiên cứu của đề tài**

Luận văn dự kiến hướng nghiên cứu là:

Nghiên cứu lý thuyết: tìm hiểu sâu hơn về thuật toán lan truyền độ tương tự, trên cơ sở nắm vững bản chất của phương pháp lan truyền độ tin cậy trong suy diễn trên mạng xác suất Bayes. dựa trên cơ sở toán học của lý thuyết xác suất.

Nghiên cứu ứng dụng: Cài đặt thử nghiệm thuật toán với dữ liệu mô phỏng để đánh giá, phân tích đánh giá kết quả; thử với dữ liệu thực tế.

## **4. Những nội dung nghiên cứu chính**

Nội dung nghiên cứu của luận văn bao gồm:

Tìm hiểu tổng quan về phân cụm dữ liệu; các điểm đặc thù của bài toán phân cụm dữ liệu địa lý; Một số đặc điểm của thuật toán lan truyền độ tương tự, trên cơ sở lý thuyết toán học hoặc phân tích thực nghiệm.

Về thực hành: Cài đặt thử nghiệm thuật toán với dữ liệu mô phỏng để đánh giá, phân tích đánh giá kết quả; thử với dữ liệu thực tế.

Làm quen với hệ thống thông tin địa lý nguồn mở; cơ sở dữ liệu địa lý; cách viết plugin tích hợp phép phân tích dữ liệu địa lý.

## **5. Phương pháp nghiên cứu**

Phương pháp nghiên cứu lý thuyết: tổng quan, phân tích các kết quả nghiên cứu đã có, nhận biết các ưu nhược điểm, lựa chọn cách tiếp cận phù hợp nhất để giải quyết bài toán ứng dụng.

Nghiên cứu thực nghiệm qua phân tích kết quả thử nghiệm với dữ liệu mô phỏng; dữ liệu thực tế; so sánh đánh giá và kết luận.

## **6. Ý nghĩa khoa học của đề tài**

Đề tài nghiên cứu có ý nghĩa khoa học, góp phần làm hiểu biết sâu sắc hơn phương pháp lan truyền độ tương tự trong phân cụm dữ liệu.

Ứng dụng thực tế: phân cụm dữ liệu môi trường, không khí, thời tiết... nhận được từ các trạm quan trắc khí tượng để xác định những tiểu vùng môi trường khí tượng trong một địa phương, khu vực.

## **7. Bố cục của luận văn**

Luận văn bao gồm 3 chương cùng với phần Mở đầu, phần Kết luận, phần Mục lục, phần Tài liệu tham khảo.

Chương 1: Tổng quan về phân cụm dữ liệu GIS và phân cụm dữ liệu địa lý: Trình bày các khái niệm cơ bản, các cách tiếp cận, phương pháp, thuật toán;