

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

PHẠM THỊ LIÊN

**BÀI TOÁN ĐỒ THỊ CON ĐẰNG CẤU TRONG
KHAI PHÁ DỮ LIỆU ĐỒ THỊ VÀ ỨNG DỤNG
PHÁT HIỆN ĐỒ THỊ CON PHỔ BIẾN**

LUẬN VĂN THẠC SĨ: KHOA HỌC MÁY TÍNH

Thái Nguyên - 2020

LỜI CAM ĐOAN

Tôi xin cam đoan luận văn này do tự bản thân tôi tìm hiểu, nghiên cứu dưới sự hướng dẫn của PGS. TS. Đoàn Văn Ban. Các chương trình thực nghiệm do chính bản thân tôi lập trình, các kết quả là hoàn toàn trung thực. Các tài liệu tham khảo được trích dẫn và chú thích đầy đủ.

TÁC GIẢ LUẬN VĂN

Phạm Thị Liên

LỜI CẢM ƠN

Tôi xin bày tỏ lời cảm ơn chân thành tới tập thể các thầy cô giáo Viện công nghệ thông tin – Viện Hàn lâm Khoa học và Công nghệ Việt Nam, các thầy cô giáo Trường Đại học Công nghệ thông tin và truyền thông - Đại học Thái Nguyên đã dạy dỗ chúng tôi trong suốt quá trình học tập chương trình cao học tại trường.

Đặc biệt tôi xin bày tỏ lòng biết ơn sâu sắc tới thầy giáo PGS.TS. Đoàn Văn Ban đã quan tâm, định hướng và đưa ra những góp ý, chỉnh sửa quý báu cho tôi trong quá trình làm luận văn tốt nghiệp. Cũng như các bạn bè, đồng nghiệp, gia đình và người thân đã quan tâm, giúp đỡ và chia sẻ với tôi trong suốt quá trình làm luận văn tốt nghiệp.

Dù đã có nhiều cố gắng nhưng chắc chắn sẽ không tránh khỏi những thiếu sót vì vậy rất mong nhận được sự đóng góp ý kiến của các thầy, cô và các bạn để luận văn này được hoàn thiện hơn.

Tôi xin chân thành cảm ơn!

Thái Nguyên, tháng 08 năm 2020

Phạm Thị Liên

MỤC LỤC

	<i>Trang</i>
MỞ ĐẦU.....	1
CHƯƠNG 1: KHAI PHÁ ĐỒ THỊ.....	3
1.1. Cấu trúc đồ thị.....	3
1.2. Các dạng biểu diễn cấu trúc dữ liệu đồ thị	6
1.2.1. Danh sách liên thuộc	6
1.2.2. Danh sách liên kề	7
1.2.3. Ma trận liên thuộc	8
1.2.4. Ma trận liên kề	9
1.2.5. Dạng chính tắc của đồ thị.....	10
1.3. Các kỹ thuật khai phá đồ thị	14
1.3.1. Phát hiện cấu trúc cộng đồng mạng xã hội	15
1.3.2. Khai phá đồ thị con thường xuyên đóng.....	19
1.4. Tổng kết chương 1	20
CHƯƠNG 2: BÀI TOÁN ĐỒ THỊ ĐẲNG CẤU VÀ KHAI PHÁ ĐỒ THỊ CON PHỔ BIẾN	21
2.1. Bài toán đồ thị đẳng cấu.....	21
2.2. Thuật toán kiểm tra đồ thị đẳng cấu	24
2.2.1. Thuật toán Dijkstra tìm đường đi ngắn nhất.....	24
2.2.2. Thuật toán tính khoảng cách $d(u, v)$ trong các đồ thị phụ thêm và đồ thị kết đôi	24
2.2.3. Thuật toán xác ma trận dấu và dạng chính tắc của nó.....	25
2.2.4. Thuật toán sắp xếp các đỉnh của hai đồ thị để kiểm tra tính đẳng cấu của chúng dựa vào dạng chính tắc	26
2.2.5. Một số tính chất của đồ thị đẳng cấu	26
2.3. Bài toán đẳng cấu đồ thị con SGI	30
2.3.1. Một số khái niệm cơ sở và ký hiệu	31
2.3.2. Cây quyết định của đồ thị	32
2.3.3. Thuật toán xây dựng cây quyết định.....	36
2.4. Khai phá đồ thị con phổ biến	40
2.4.1. Cây các đồ thị con dạng chính tắc	40
2.4.2. Phép kết nối N-Join hai đồ thị	41

2.4.3. Phép N-Extension	43
2.5. Thuật toán FFSM cho khai phá đồ thị con phổ biến trong CSDL đồ thị	44
2.6. Kết luận chương 2	47
CHƯƠNG 3 THỬ NGHIỆM VÀ ĐÁNH GIÁ	48
3.1. Dữ liệu và môi trường thử nghiệm	48
3.1.1. Bộ dữ liệu thử nghiệm	48
3.1.2. Môi trường thử nghiệm.....	49
3.2. Cài đặt và thử nghiệm thuật toán tìm kiếm tra đồ thị đẳng cấu.....	50
3.2.1. Mô tả yêu cầu bài toán kiểm tra đồ thị đẳng cấu	50
3.2.2. Kết quả thử nghiệm.....	50
3.3. Thử nghiệm thuật toán FFSM cho khai phá đồ thị con phổ biến	53
3.3.1. Mô tả yêu cầu bài toán khai phá đồ thị con phổ biến	53
3.3.2. Phân tích đánh giá kết quả	53
3.4. Kết luận chương 3	55
KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN	56
TÀI LIỆU THAM KHẢO	57

THUẬT NGỮ VÀ TỪ VIẾT TẮT

Thuật ngữ	Diễn giải
CAM	Canonical Adjacency Matrix
Community Social Structure	Cấu trúc cộng đồng mạng xã hội
CSDL	Cơ sở dữ liệu
FFSM	FastFrequent Subgraph Mining
GI	Graph isomorphism - Đồng cấu đồ thị
Graph isomorphism problem	Bài toán đồ thị đồng cấu
SFI	Viện Santa Fe
SGI	SubGraph Isomorphism – Đồng cấu đồ thị con

DANH MỤC CÁC HÌNH VẼ

Hình 1.1 Các đồ thị vô hướng và có hướng.....	4
Hình 1.2. Hợp của hai đồ thị.....	5
Hình 1.3. G_1' là phân bù của G_1	5
Hình 1.4. Danh sách liên thuộc của đồ thị.....	7
Hình 1.5. Đồ thị vô hướng.....	7
Hình 1.6. Biểu diễn danh sách liên kề của đồ thị hình 1.5.....	8
Hình 1.7. Ma trận liên thuộc của đồ thị vô hướng.....	8
Hình 1.8. Ma trận liên kề của đồ thị trong ví dụ 1.6.....	9
Hình 1.9. Các đồ thị con cực đại.....	11
Hình 1.10. Mạng lưới cộng tác của các nhà khoa học làm việc tại SFI [1].....	17
Hình 2.1. Các phần tử hàng - cột của ma trận liên kề.....	33
Hình 2.2. Cây quyết định để phân lớp các ma trận liên kề của G	34
Hình 2.3. Cấu trúc từ điển và các chỉ mục cho cây quyết định.....	35
Hình 2.4. Cây quyết định của G và G'	35
Hình 2.5. Cây quyết định compact để phân lớp các ma trận liên kề $\{A, B, \dots, F\}$ và $\{A', B', \dots, F'\}$ của đồ thị G và G'	37
Hình 2.6. Hai phép toán N-Join và N-Extension.....	43
Hình 3.1. Mô tả cấu trúc bộ dữ liệu thử nghiệm với kho gồm 2 đồ thị.....	48
Hình 3.2. Bộ dữ liệu đơn giản thử nghiệm thuật toán kiểm tra 2 đồ thị đẳng cấu.....	51
Hình 3.3. Kết quả đồ thị đẳng cấu với bộ dữ liệu đơn giản.....	51
Hình 3.4. Kết quả đồ thị đẳng cấu với đồ thị target có kích thước lớn.....	52
Hình 3.5. Ví dụ về không tồn tại đẳng cấu đồ thị con.....	52
Hình 3.6. Thời gian chạy thuật toán FFSM trên bộ dữ liệu 10000 đồ thị.....	53
Hình 3.7. Số đồ thị con tìm được ứng với các độ hỗ trợ tối thiểu khác nhau trên bộ dữ liệu 10.000 đồ thị.....	54

DANH MỤC CÁC BẢNG

Bảng 3.1. Thông tin phân cứng được sử dụng thử nghiệm	49
Bảng 3.2. Tổng hợp kết quả chạy thuật toán FFSM trên 2 bộ dữ liệu	54

MỞ ĐẦU

1. Lý do chọn đề tài

Khai phá dữ liệu là lĩnh vực đang được nhiều người tập trung nghiên cứu và phát triển nhiều ứng dụng phổ biến. Trong đó, bài toán khai phá đồ thị con thường xuyên đã và đang thu hút được nhiều sự quan tâm nghiên cứu bởi phạm vi ứng dụng quan trọng trong nhiều lĩnh vực khác nhau như trong tin-sinh (bioinformatics), tin-hóa (cheminformatics) khai thác đăng nhập trang web, lập chỉ mục video, hay lập chỉ mục cơ sở dữ liệu hiệu quả, ... Vấn đề khai phá mẫu thường xuyên là từ một tập dữ liệu các đối tượng với một ngưỡng độ hỗ trợ tối thiểu minsup. Dữ liệu có thể rất đa dạng từ dữ liệu nhị phân, dữ liệu số nguyên, số thực hoặc các dữ liệu có cấu trúc phức tạp hơn như cấu trúc cây, đồ thị, ... Cho đến nay vẫn chưa có lời giải hiệu quả cho bài toán này do độ phức tạp của bài toán là rất lớn khi đồ thị có số đỉnh lớn và mật độ các cạnh dày. Tuy nhiên, sự phức tạp của những vấn đề này sẽ giảm khi cơ sở dữ liệu (CSDL) đồ thị có thêm thông tin về các đỉnh và các cạnh đã được gán nhãn. Có thể sử dụng các nhãn để hạn chế các đỉnh có thể tạo thành các cặp trong quá trình kiểm tra sự đẳng cấu của đồ thị con.

Và đó chính là lý do em lựa chọn đề tài “*Bài toán đồ thị con đẳng cấu trong khai phá dữ liệu đồ thị và ứng dụng phát hiện đồ thị con phổ biến*” để nghiên cứu làm luận văn thạc sĩ của mình.

2. Mục đích nghiên cứu

Nghiên cứu thuật toán đẳng cấu đồ thị và thuật toán phát hiện đồ thị con phổ biến trong CSDL đồ thị.

3. Đối tượng nghiên cứu

- Bài toán đồ thị con đẳng cấu
- Khai phá dữ liệu đồ thị
- Bài toán khai phá đồ thị con phổ biến trong CSDL đồ thị
- Thuật toán FFSM
- Thuật toán SGI Decision Tree.

4. Phạm vi nghiên cứu

+ Lý thuyết:

- Tìm hiểu lý thuyết đồ thị và các phương pháp khai phá dữ liệu đồ thị.
- Nghiên cứu các thuật toán tìm đồ thị con đẳng cấu và các thuật toán khai phá đồ thị con phổ biến.

+ Thực nghiệm:

- Xây dựng chương trình thực nghiệm thuật toán.
- Áp dụng chương trình trên một số bộ cơ sở dữ liệu trong khai phá dữ liệu đồ thị.

5. Ý nghĩa khoa học

Đây là một hướng nghiên cứu lý thuyết, xây dựng các thuật toán tìm kiếm đồ thị con đẳng cấu và đồ thị con phổ biến và kiểm thử trên các bộ dữ liệu về đồ thị.

6. Phương pháp nghiên cứu

- *Phương pháp nghiên cứu lý thuyết:* Suu tập các tài liệu có liên quan đến đề tài, nghiên cứu để hiểu sâu các nội dung vấn đề cần nghiên cứu.

- *Phương pháp nghiên cứu thực nghiệm:* Cài đặt chương trình thử nghiệm phương pháp tìm kiếm đồ thị con đẳng cấu và thuật toán FFSSM để liệt kê các đồ thị con phổ biến trên các bộ dữ liệu và đánh giá hiệu suất của phương pháp.

- *Phương pháp trao đổi khoa học:* Trao đổi nội dung nghiên cứu với giáo viên hướng dẫn, với các đồng nghiệp để đề xuất và giải quyết các nội dung luận văn đề ra.

7. Cấu trúc đề tài

Luận văn được chia thành các phần chính như sau:

Chương 1: Khai phá dữ liệu đồ thị

Chương 2: Đồ thị đẳng cấu và đồ thị con phổ biến

Chương 3: Thử nghiệm thuật toán tìm đồ thị con phổ biến

Kết luận và hướng phát triển.