

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

LÊ THỊ MAI HƯƠNG

NGHIÊN CỨU MỘT SỐ THUẬT TOÁN
PHÂN CỤM DỮ LIỆU NỬA GIÁM SÁT VÀ ỨNG DỤNG
PHÂN ĐOẠN ẢNH X-QUANG

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

THÁI NGUYÊN - 2017

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

LÊ THỊ MAI HƯƠNG

**NGHIÊN CỨU MỘT SỐ THUẬT TOÁN
PHÂN CỤM DỮ LIỆU NỬA GIÁM SÁT VÀ ỨNG DỤNG
PHÂN ĐOẠN ẢNH X-QUANG**

Chuyên ngành: Khoa học máy tính

Mã số: 60 48 01

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Giáo viên hướng dẫn: TS.Nguyễn Đình Dũng

THÁI NGUYÊN - 2017

LỜI CAM ĐOAN

Tôi xin cam đoan luận văn này do chính tôi thực hiện, dưới sự hướng dẫn khoa học của TS. Nguyễn Đình Dũng, các kết quả lý thuyết được trình bày trong luận văn là sự tổng hợp từ các kết quả đã được công bố và có trích dẫn đầy đủ, kết quả của chương trình thực nghiệm trong luận văn này được tác giả thực hiện là hoàn toàn trung thực, nếu sai tôi hoàn toàn chịu trách nhiệm.

Thái Nguyên, tháng 6 năm 2016

Học viên

Lê Thị Mai Hương

LỜI CẢM ƠN

Luận văn này được hoàn thành tại Trường Đại học Công nghệ Thông tin và Truyền thông dưới sự hướng dẫn của TS. Nguyễn Đình Dũng. Tác giả xin bày tỏ lòng biết ơn tới các thầy cô giáo thuộc Trường Đại học Công nghệ Thông tin và Truyền thông, các thầy cô giáo thuộc Viện Công nghệ Thông tin – Viện Hàn lâm Khoa học và Công nghệ Việt Nam đã tạo điều kiện, giúp đỡ tác giả trong quá trình học tập và làm luận văn tại Trường, đặc biệt tác giả xin bày tỏ lòng biết ơn tới TS. Nguyễn Đình Dũng đã tận tình hướng dẫn và cung cấp nhiều tài liệu cần thiết để tác giả có thể hoàn thành luận văn đúng thời hạn.

Xin chân thành cảm ơn anh chị em học viên cao học và bạn bè đồng nghiệp đã trao đổi, khích lệ tác giả trong quá trình học tập và làm luận văn tại Trường Đại học Công nghệ Thông tin và Truyền thông – Đại học Thái Nguyên.

Cuối cùng tác giả xin gửi lời cảm ơn đến gia đình, những người đã luôn bên cạnh, động viên và khuyến khích tôi trong quá trình thực hiện đề tài.

Thái Nguyên, ngày tháng năm 2017

Học viên

Lê Thị Mai Hương

MỤC LỤC

LỜI CẢM ƠN	i
LỜI CAM ĐOAN.....	i
DANH MỤC TỪ VIẾT TẮT.....	v
DANH MỤC HÌNH VẼ.....	vi
LỜI MỞ ĐẦU	1
CHƯƠNG 1: TỔNG QUAN VỀ PHÂN CỤM DỮ LIỆU VÀ BÀI TOÁN PHÂN ĐOẠN ẢNH X-QUANG NHA KHOA	3
1.1. Khai phá dữ liệu.....	3
1.1.1. Khái niệm khai phá dữ liệu	3
1.1.2. Quá trình khai phá tri thức trong cơ sở dữ liệu	3
1.1.3. Các kỹ thuật tiếp cận trong khai phá dữ liệu:.....	5
1.2. Phân cụm dữ liệu	6
1.2.1. Khái niệm phân cụm dữ liệu	6
1.2.2. Các bước cơ bản để phân cụm dữ liệu	6
1.2.3. Các kiểu dữ liệu và độ đo tương tự, độ đo phi tương tự.....	7
1.2.3.1. Phân loại kiểu dữ liệu dựa trên kích thước miền	7
1.2.3.2. Phân loại kiểu dữ liệu dựa trên hệ đo	7
1.2.3.3. Khái niệm và phép đo độ tương tự	9
1.2.4. Các yêu cầu đối với kỹ thuật phân cụm dữ liệu	12
1.2.5. Ứng dụng của phân cụm dữ liệu	14
1.3. Cấu trúc giải phẫu răng.....	15
1.3.1. Cấu trúc giải phẫu răng	15
1.3.2. Phân loại ảnh X - quang nha khoa	17
1.4. Bài toán phân đoạn ảnh X - quang nha khoa.....	19
1.4.1. Phân đoạn ảnh	19
1.4.2. Phân loại các phương pháp phân đoạn ảnh	20
1.4.3. Phân đoạn ảnh X – quang nha khoa.....	21
KẾT LUẬN CHƯƠNG 1.....	23
CHƯƠNG 2: MỘT SỐ THUẬT TOÁN PHÂN CỤM NỬA GIÁM SÁT ..	24
2.1. Phân cụm mờ	24
2.1.1. Các khái niệm cơ bản về tập mờ	24

2.1.2. Thuật toán phân cụm mờ FCM (Fuzzy C-Means)	28
2.2. Thuật toán phân cụm nửa giám sát mờ bằng phương pháp học tích cực	31
2.3. Thuật toán phân cụm nửa giám sát mờ chuẩn (SSSFC).....	33
2.4. Thuật toán phân cụm nửa giám sát mờ theo quy tắc entropy (eSFCM)	35
2.5. Thuật toán nửa giám sát mờ lai ghép	36
2.5.1. Lược đồ tổng quan lai ghép.....	36
2.5.2. Thuật toán tách ngưỡng Otsu	38
2.5.3. Thuật toán phân cụm nửa giám sát mờ lai ghép	40
KẾT LUẬN CHƯƠNG 2	41
CHƯƠNG 3: XÂY DỰNG ỨNG DỤNG PHÂN ĐOẠN ẢNH X – QUANG	
NHA KHOA	42
3.1. Đặc tả yêu cầu.....	42
3.1.1. Yêu cầu thực tế.....	42
3.1.2. Mục đích của ứng dụng.....	43
3.2. Đặc tả dữ liệu.....	43
3.3. Các bước phân đoạn ảnh.....	44
3.4. Thiết kế hệ thống	45
3.4.1. Chức năng phân đoạn ảnh X – quang nha khoa.....	45
3.4.2. Chức năng xem chi tiết kết quả.....	46
3.4.3. Chức năng đánh giá chất lượng phân đoạn	47
3.5. Minh họa các chức năng của ứng dụng	48
3.5.1. Giao diện chính của ứng dụng.....	48
3.5.2. Chọn ảnh cần phân đoạn	49
3.5.3. Phân đoạn ảnh bằng thuật toán FCM.....	49
3.5.4. Phân đoạn ảnh bằng thuật toán nửa giám sát mờ.....	50
3.5.5. Chọn độ đo đánh giá kết quả phân cụm	50
3.6. Đánh giá kết quả phân đoạn	51
KẾT LUẬN CHƯƠNG 3	52
KẾT LUẬN	53
TÀI LIỆU THAM KHẢO	54
PHỤ LỤC	57
CODE MATLAB CỦA ỨNG DỤNG PHÂN ĐOẠN ẢNH BẰNG THUẬT	
TOÁN BÁN GIÁM SÁT MỜ LAI GHÉP	57

DANH MỤC TỪ VIẾT TẮT

Từ viết tắt	Từ viết đầy đủ
DB	Davies-Bouldin
eSFCM	Semi-supervised Entropy regularized Fuzzy Clustering
FCM	Fuzzy C-Mean
PBM	Pakhira, Bandyopadhyay and Maulik
SSFCM	Semi-Supervised Fuzzy C-Mean
SSSFC	Semi-Supervised Standard Fuzzy Clustering
SSWC	Simplified Silhouette Width Criterion
CSDL	Cơ sở dữ liệu
PCDL	Phân cụm dữ liệu

DANH MỤC HÌNH VẼ

Hình 1.1. Quá trình khám phá tri thức trong CSDL	4
Hình 1.2. Cơ quan răng (răng và nha chu).....	15
Hình 1.3. Một số loại ảnh X-Quang nha khoa	19
Hình 1.4. Những khó khăn trong việc phân đoạn ảnh nha khoa.....	22
Hình 2.1. Hàm thuộc tuyến tính.....	25
Hình 2.2. Hàm thuộc dạng sin.....	25
Hình 2.3. Hàm thuộc Gauss	26
Hình 2.4. Bao trong của tập mờ	26
Hình 2.5. Phép hợp tập mờ dạng 1	27
Hình 2.6. Phép giao tập mờ dạng 1	28
Hình 2.7. Phần bù của tập mờ trung bình	28
Hình 2.8. Lược đồ tổng quan của thuật toán lai ghép.....	37
Hình 3.1: Ảnh dữ liệu đầu vào của ứng dụng	44
Hình 3.2: Biểu đồ usecase mô tả chức năng của ứng dụng	45
Hình 3.3: Biểu đồ trình tự chức năng phân đoạn ảnh	46
Hình 3.4: Biểu đồ trình tự chức năng xem kết quả	47
Hình 3.5: Biểu đồ trình tự chức năng đánh giá kết quả	48
Hình 3.6: Giao diện chính của phần mềm.....	48
Hình 3.7: Chọn ảnh cần phân đoạn	49
Hình 3.8. Kết quả phân đoạn bằng FCM	49
Hình 3.9. Kết quả phân đoạn bằng SSSFC	50
Hình 3.10. Đánh giá kết quả phân đoạn	50

LỜI MỞ ĐẦU

Khai phá dữ liệu là một tập hợp các kỹ thuật được sử dụng để tự động khai thác và tìm ra các mối quan hệ lẫn nhau của dữ liệu trong một tập hợp dữ liệu khổng lồ và phức tạp đồng thời cũng tìm ra các mẫu tiềm ẩn trong dữ liệu đó. Hiện nay việc khai phá dữ liệu được nghiên cứu theo các hướng mô tả khái niệm, luật kết hợp, phân lớp và dự đoán, phân cụm (xem [1], [2], [7]) và có nhiều ứng dụng trong thực tế, trong đó phân đoạn ảnh X-Quang trong lĩnh vực y tế là một ứng dụng điển hình [13]. Ngày nay, việc xử lý các hình ảnh y tế có vai trò quan trọng trong việc tự động hóa phân tích, hỗ trợ chẩn đoán và điều trị các bệnh khác nhau. Trong đó, quá trình phân đoạn thường được yêu cầu như là giai đoạn sơ bộ. Tuy nhiên các phân vùng trong hình ảnh y tế rất phức tạp nên việc phân đoạn chính xác là rất quan trọng.

Trong các phương pháp phân đoạn ảnh hiện có, phân cụm là một phương pháp được sử dụng rộng rãi bởi tính đơn giản và hiệu quả mà nó mang lại (xem [8]-[12]). Phân cụm dữ liệu là lĩnh vực học máy không giám sát, nó có chức năng tổ chức một tập đối tượng dữ liệu thành các cụm sao cho những đối tượng trong cùng một cụm thì tương tự như nhau còn các đối tượng ở các cụm khác nhau thì kém tương tự nhau hơn. Nhược điểm chung của thuật toán phân cụm là chất lượng phân cụm phụ thuộc nhiều vào các tham số và thông tin khởi tạo. Để giảm thiểu các hạn chế này, gần đây đã có nhiều tác giả (xem [8]-[12]) giải quyết theo cách tiếp cận nửa giám sát, trong đó việc phân cụm được thực hiện dựa vào các thông tin hỗ trợ đóng vai trò điều khiển quá trình phân cụm, nhờ đó mà chất lượng phân cụm được nâng lên đáng kể.

Mục tiêu của luận văn là nghiên cứu, tìm hiểu một số thuật toán phân cụm nửa giám sát và xây dựng được một ứng dụng thử nghiệm cho thuật toán phân đoạn ảnh X-quang hỗ trợ chuẩn đoán bệnh trong lĩnh vực nha khoa. Các kết quả đạt được trong luận văn này là kết quả trong quá trình học tập và nghiên cứu của tác giả tại Trường Đại học Công nghệ Thông tin và Truyền thông. Ngoài phần

mở đầu, kết luận và tài liệu tham khảo, nội dung luận văn được trình bày thành ba chương:

Chương 1 trình bày các khái niệm cơ bản về phân cụm dữ liệu và bài toán phân đoạn ảnh X-quang nha khoa.

Chương 2, tác giả tìm hiểu một số thuật toán phân cụm dữ liệu trong đó tập trung nghiên cứu thuật toán phân cụm dữ liệu nửa giám sát.

Chương 3 là kết quả thực nghiệm cho thuật toán phân cụm nửa giám sát đối với bài toán phân đoạn ảnh X-quang nha khoa.