

**ĐẠI HỌC QUỐC GIA HÀ NỘI
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ**

PHẠM TUẤN DŨNG

**NGHIÊN CỨU PHƯƠNG PHÁP PHÂN LOẠI VÀ
XÂY DỰNG CƠ SỞ DỮ LIỆU LỚP PHỦ ĐỒ THỊ
TẠI VIỆT NAM SỬ DỤNG DỮ LIỆU ĐA NGUỒN**

Chuyên ngành: Hệ thống thông tin

Mã số: 9480101.01

**TÓM TẮT LUẬN ÁN TIẾN SĨ
CÔNG NGHỆ THÔNG TIN**

Hà Nội – 2021

Công trình được hoàn thành tại: Trường Đại học Công nghệ, Đại học Quốc gia Hà Nội

Tập thể cán bộ hướng dẫn:

Hướng dẫn chính: PGS.TS. Doãn Minh Chung

Cơ quan công tác: Viện Công nghệ Vũ trụ, Viện Hàn lâm KH&CN VN

Hướng dẫn phụ: TS Bùi Quang Hưng

Cơ quan công tác: Trường Đại học Công nghệ, ĐHQGHN

Phản biện:

.....

Phản biện:

.....

Phản biện:

.....

L luận án sẽ được bảo vệ trước Hội đồng cấp Đại học Quốc gia
chấm luận án tiến sĩ họp tại

vào hồi giờ ngày tháng năm

Có thể tìm hiểu luận án tại:

- Thư viện Quốc gia Việt Nam
- Trung tâm Thông tin - Thư viện, Đại học Quốc gia Hà Nội

1. Lý do chọn đề tài

Trên thực tế, nghiên cứu các phương pháp phân loại lớp phủ đô thị trên phạm vi toàn cầu là một công việc tương đối khó khăn bởi quá trình thu thập, xử lý dữ liệu gặp nhiều thách thức. Khi sử dụng các bộ dữ liệu lớp phủ toàn cầu cho các nghiên cứu về khu vực, nếu không có các dữ liệu mặt đất tin cậy, thì độ chính xác của các bản đồ lớp phủ giảm xuống rõ rệt. Các vấn đề đặt ra khi xây dựng bản đồ lớp phủ đô thị cho Việt Nam dựa trên các bộ dữ liệu lớp phủ toàn cầu đó là: thiếu các đặc trưng cho khu vực nghiên cứu; sự suy giảm độ chính xác do dữ liệu đầu vào có độ phân giải không cao; các thách thức khi xây dựng các ứng dụng đáp ứng được nhu cầu của các nhà nghiên cứu về độ chính xác, linh hoạt, cập nhật, chia sẻ dữ liệu nhanh chóng.

Bên cạnh đó, các nghiên cứu về phương pháp phân loại lớp phủ đô thị từ dữ liệu viễn thám gặp phải một số thách thức về độ phân giải không gian của dữ liệu, sự thay đổi về bề mặt và sóng bức xạ theo các mùa trong năm, các vấn đề nảy sinh khi xử lý dữ liệu đa nguồn, cụ thể:

Thứ nhất, các dữ liệu viễn thám được sử dụng trong các nghiên cứu về đô thị thường có độ phân giải cao hoặc trung bình thu nhận từ các cảm biến đa phổ trên các vệ tinh viễn thám. Đối với các dữ liệu có độ phân giải không gian trung bình, mỗi điểm ảnh tương ứng trên mặt đất có thể chứa nhiều loại lớp phủ khác nhau. Điều này khiến cho các dữ liệu được thu nhận bởi cảm biến không đồng nhất, từ đó dẫn đến sự suy giảm độ chính xác của các phương pháp phân loại lớp phủ mặt đất nói chung và lớp phủ đô thị nói riêng. Đối với các dữ liệu có độ phân giải cao, đòi hỏi phải thu thập, lưu trữ, xử lý lượng dữ liệu rất lớn, do vậy không phù hợp với các bài toán phân loại lớp phủ trên

phạm vi rộng. Bên cạnh đó, đây đều là những vệ tinh thương mại, vì vậy chi phí dành cho việc mua các dữ liệu rất cao, không phù hợp với các nghiên cứu khoa học. Ngoài ra, đối với bài toán phân loại lớp phủ đô thị, các dữ liệu quang phổ có độ phân giải cao thường bị ảnh hưởng bởi hiệu ứng đổ bóng từ các tòa nhà cao tầng, ảnh hưởng đến độ chính xác của các dữ liệu khu vực xung quanh và kết quả đầu ra của phương pháp phân loại.

Thứ hai, các lớp phủ mặt đất có sự thay đổi theo từng thời điểm trong năm do ảnh hưởng bởi mặt trời, nhiệt độ, độ ẩm,... Ví dụ: thực vật, bề mặt nước bị ảnh hưởng bởi mùa mưa, mùa khô; lớp phủ thực vật thay đổi mạnh theo mùa và theo chu trình phát triển, thu hoạch. Đối với khu vực nhiệt đới gió mùa như Việt Nam, thời tiết chia thành bốn mùa rõ rệt, các đối tượng trên mặt đất cũng có sự phản xạ ánh nắng mặt trời khác nhau trong từng mùa, dẫn đến tín hiệu thu được trên cảm biến vệ tinh của cùng một đối tượng cũng khác nhau tùy theo từng thời điểm.

Thứ ba, việc kết hợp nhiều nguồn dữ liệu khác nhau nhằm mục đích nâng cao độ chính xác của kết quả phân loại, bằng cách tận dụng những ưu điểm của từng loại dữ liệu, cũng nảy sinh những thách thức khi xử lý các dữ liệu trong bài toán phân loại lớp phủ đô thị. Do dữ liệu được thu thập từ nhiều nguồn sẽ có sự khác nhau về kiểu dữ liệu, độ phân giải, thời điểm thu thập,... đòi hỏi phải có quá trình tiền xử lý trước khi dùng làm dữ liệu đầu vào của các phương pháp phân loại. Quá trình tiền xử lý phải sử dụng các phương pháp tái lấy mẫu phù hợp với từng loại dữ liệu khác nhau, tuy nhiên các phương pháp này

cũng ảnh hưởng tới chất lượng dữ liệu và độ chính xác của các phương pháp phân loại lớp phủ đô thị.

Ngoài ra, quá trình đô thị hoá nhanh chóng cũng dẫn tới những tác động tới cảnh quan thiên nhiên, khí hậu, môi trường,... trong đó có những tác động theo chiều hướng tiêu cực, đặc biệt là môi trường. Trong các hậu quả không mong muốn đó, ô nhiễm không khí là một trong những vấn đề được quan tâm hàng đầu bởi nó ảnh hưởng đến nhiều mặt của đời sống kinh tế - xã hội của con người. Sự phát triển bùng nổ của các đô thị trên thế giới trong vài thập niên gần đây đã đặt ra những thách thức cho các nhà hoạch định chính sách phát triển đô thị và các nhà nghiên cứu về lớp phủ đô thị. Những ảnh hưởng tiêu cực của quá trình đô thị hoá đến môi trường như ô nhiễm không khí, nguồn nước, biến đổi khí hậu,... đã được quan tâm trong nhiều nghiên cứu của các nhà khoa học trên thế giới. Sự liên hệ giữa quá trình đô thị hoá và các chỉ số môi trường có thể được phân tích thông qua việc phân loại lớp phủ đô thị và tính toán sự mở rộng đô thị dựa trên các cơ sở dữ liệu viễn thám.

Chính vì các lý do trên, nghiên cứu sinh đã lựa chọn đề tài “Nghiên cứu phương pháp phân loại và xây dựng cơ sở dữ liệu lớp phủ đô thị tại Việt Nam sử dụng dữ liệu đa nguồn” làm đề tài nghiên cứu trong luận án của mình.

2. Mục tiêu nghiên cứu của luận án

- Nghiên cứu cơ sở khoa học của phương pháp phân loại lớp phủ đô thị Việt Nam sử dụng dữ liệu đa nguồn. Luận án tập trung phân tích các phương pháp phân loại lớp phủ mặt đất, lớp phủ đô thị trên phạm vi toàn cầu và khu vực; các phương pháp tái lấy mẫu dữ liệu viễn thám.

- Nghiên cứu và cải tiến phương pháp phân loại lớp phủ đô thị của GLCNMO cho khu vực Việt Nam trên cơ sở lựa chọn dữ liệu và tính toán các ngưỡng phù hợp.

- Nghiên cứu các phương pháp tái lấy mẫu đối với dữ liệu viễn thám đa nguồn trong bài toán phân loại lớp phủ đô thị tại Việt Nam.

- Xây dựng cơ sở dữ liệu lớp phủ đô thị tại Việt Nam ứng dụng trong đánh giá sự ảnh hưởng của quá trình phát triển đô thị tới vấn đề ô nhiễm không khí tại Việt Nam.

3. Phạm vi nghiên cứu của luận án

Luận án tập trung nghiên cứu các vấn đề liên quan đến cải tiến phương pháp phân loại lớp phủ đô thị của GLCNMO cho khu vực Việt Nam, so sánh các phương pháp tái lấy mẫu trên dữ liệu viễn thám trong bài toán phân lớp đô thị tại Việt Nam; xây dựng cơ sở dữ liệu lớp phủ đô thị tại Việt Nam ứng dụng trong đánh giá sự ảnh hưởng của quá trình phát triển đô thị tới vấn đề ô nhiễm không khí tại Việt Nam.

4. Đóng góp của luận án

- Cải tiến phương pháp phân loại lớp phủ toàn cầu của GLCNMO nhằm xây dựng bản đồ lớp phủ đô thị cho Việt Nam. Các nội dung cải tiến đó là: (i) đề xuất một phương pháp lấy mẫu ngẫu nhiên dựa trên việc tổng hợp các bộ dữ liệu lớp phủ toàn cầu, kết hợp với các ảnh có độ phân giải cao như Google Earth và Landsat ETM+ và công cụ trong ArcGIS và Python, (ii) đề xuất phương pháp tính ngưỡng dựa trên histogram của tập mẫu. Kết quả đánh giá cho thấy việc hiệu chỉnh dữ liệu đầu vào, lựa chọn giá trị ngưỡng phù hợp với các thông số thu thập tại Việt Nam giúp nâng cao độ chính xác của dữ liệu lớp phủ đô thị tại Việt Nam.

- Đánh giá sự ảnh hưởng của quá trình tái lấy mẫu tới chất lượng của ảnh viễn thám và tác động của quá trình này đến độ chính xác của phương pháp phân loại lớp phủ mặt đất tại Việt Nam.

- Xây dựng cơ sở dữ liệu lớp phủ đô thị tại Việt Nam ứng dụng trong đánh giá sự ảnh hưởng của quá trình phát triển đô thị tới vấn đề ô nhiễm không khí tại Việt Nam. Luận án xây dựng và quản lý cơ sở dữ liệu lớp phủ đô thị Việt Nam trên nền nền tảng xử lý, quản lý và phân tích dữ liệu không gian lớn SEAP. Nghiên cứu dựa trên dữ liệu viễn thám và dữ liệu thu thập được tại các trạm mặt đất để xây dựng bản đồ đô thị và bản đồ PM2.5 tại Việt Nam trong các năm 2004, 2008, 2012 và 2015. Trên cơ sở đó, luận án đã đạt được các kết quả sau: (i) Đánh giá sự mở rộng đô thị của Việt Nam từ năm 2004 đến năm 2015; (ii) Phân tích tình trạng ô nhiễm không khí của Việt Nam từ năm 2004 đến năm 2015; (iii) Tìm hiểu tác động của việc mở rộng đô thị đối với ô nhiễm không khí ở Việt Nam từ năm 2004 đến năm 2015.

CHƯƠNG 1: TỔNG QUAN VỀ LỚP PHỦ MẶT ĐẤT, LỚP PHỦ ĐÔ THỊ VÀ CÁC PHƯƠNG PHÁP PHÂN LOẠI LỚP PHỦ MẶT ĐẤT SỬ DỤNG DỮ LIỆU VIỄN THÁM

1.1. Tổng quan về lớp phủ mặt đất, lớp phủ đô thị

1.1.1. Nghiên cứu về lớp phủ mặt đất trên thế giới

1.1.2. Các cơ sở dữ liệu lớp phủ mặt đất toàn cầu

1.1.3. Nghiên cứu về lớp phủ đô thị trên thế giới

1.1.4. Các nghiên cứu về lớp phủ mặt đất và lớp phủ đô thị tại

Việt Nam

1.2. Phân loại lớp phủ mặt đất sử dụng dữ liệu viễn thám

1.2.1. Dữ liệu viễn thám sử dụng trong nghiên cứu về phân loại lớp phủ mặt đất

1.2.2. Quy trình xử lý dữ liệu viễn thám trong nghiên cứu về phân loại lớp phủ mặt đất

1.2.2.1 Tiền xử lý dữ liệu viễn thám trong nghiên cứu về phân loại lớp phủ mặt đất

Có hai quá trình tái lấy mẫu ảnh viễn thám phổ biến đó là tổng hợp giá trị (biến đổi ảnh có độ phân giải cao thành ảnh có độ phân giải thấp) và nội suy giá trị (biến đổi ảnh có độ phân giải thấp thành ảnh có độ phân giải cao).

a, Các phương pháp tổng hợp giá trị ảnh viễn thám

- * Phương pháp dựa trên luật đa số
- * Phương pháp lấy giá trị ngẫu nhiên
- * Phương pháp lấy giá trị điểm trung tâm
- * Phương pháp lấy giá trị trung bình
- * Phương pháp lấy giá trị cực đại hoặc cực tiểu
- * Phương pháp lấy giá trị trung bình dựa trên trọng số

b, Các phương pháp nội suy giá trị ảnh viễn thám

- * Nội suy láng giềng gần nhất
- * Nội suy song tuyến tính
- * Nội suy xoắn bậc ba

Các chỉ số đánh giá phương pháp tái lấy mẫu ảnh viễn thám

a, Sai số bình phương trung bình (MSE)

b, Tỷ số tín hiệu cực đại/nhiều (PSNR)

c, Chỉ số so sánh sự tương đồng cấu trúc (SSIM)

1.2.2.2 Các phương pháp phân loại lớp phủ mặt đất sử dụng dữ liệu viễn thám

Các phương pháp phân loại lớp phủ mặt đất

Các chỉ số đánh giá phương pháp phân loại lớp phủ mặt đất

CHƯƠNG 2: NGHIÊN CỨU PHƯƠNG PHÁP PHÂN LOẠI LỚP PHỦ ĐÔ THỊ TẠI VIỆT NAM

2.1. Đặt vấn đề

2.2. Phương pháp phân loại lớp phủ toàn cầu GLCNMO

Bộ dữ liệu lớp phủ toàn cầu (Global Land Cover by National Mapping Organizations - GLCNMO) được Trung tâm Viễn thám môi trường (Trung tâm CEReS), Đại học Chiba – Nhật Bản phát triển từ năm 2003 trong khuôn khổ dự án Xây dựng bản đồ toàn cầu (Global Mapping Project - GMP) do Nhật Bản đề xuất tại Hội nghị về Môi trường và Phát triển của Liên hiệp quốc diễn ra tại Rio de Janeiro năm 1992. Hệ thống CEReS Gaia được Trung tâm CEReS phát triển từ năm 2003-2013 với sự tài trợ của JSPS (Japan Society for the Promotion of Science). Chức năng chính của hệ thống này là tích hợp, quản lý, chia sẻ dữ liệu không gian địa lý toàn cầu và khu vực. Năm 2003, dự án cơ sở dữ liệu lớp phủ mặt đất toàn cầu GLCNMO được xây dựng dựa trên dữ liệu MODIS 500m, cung cấp dữ liệu lớp phủ toàn cầu với độ chính xác cao với sự cộng tác của 40 quốc gia trên thế giới trong việc cung cấp dữ liệu địa phương và kiểm chứng phương pháp phân loại lớp phủ tại quốc gia của mình. Phiên bản 2 năm 2008 có thêm 14 quốc gia tham gia. Phiên bản 3 công bố trong năm 2017 với các dữ liệu

được thu thập từ nhiều nguồn khác nhau trong đó các dữ liệu MODIS được thu thập trong năm 2013.

2.2.1. Các nguồn dữ liệu được sử dụng trong phương pháp.

2.2.1.1. Dữ liệu mật độ dân số toàn cầu năm 2008

2.2.1.2. Dữ liệu ánh sáng ban đêm toàn cầu DMSP-OLS

2.2.1.3. Dữ liệu bề mặt không thấm nước toàn cầu EstISA 2010

2.2.1.4. Dữ liệu thu nhập bình quân đầu người của các quốc gia năm 2008

2.2.1.5. Dữ liệu MODIS-NDVI năm 2008

2.2.2. Phương pháp phân loại lớp phủ đô thị của GLCNMO

2.2.2.1. Quy trình xử lý dữ liệu

Quá trình tạo bản đồ lớp phủ đô thị toàn cầu của bộ dữ liệu GLCNMO gồm có 5 bước cơ bản:

- Bước 1: Các dữ liệu đầu vào như bản đồ phân bố dân cư LandScan, bản đồ ánh sáng ban đêm DMSP-OLS và bản đồ bề mặt không thấm nước EstISA có độ phân giải không gian 1km được biến đổi bằng các phương pháp tái lấy mẫu cho kết quả là các bản đồ có độ phân giải 500m.

- Bước 2: Dữ liệu thu nhập bình quân đầu người của các quốc gia trên thế giới năm 2008 được sử dụng để chia các nước vào bốn nhóm dựa theo mức độ phát triển kinh tế.

- Bước 3: Từ dữ liệu NDVI cao nhất nhận được từ quá trình xử lý ảnh MODIS, các khu vực chứa nhiều thực vật (như các công viên lớn trong lòng thành phố, các khu vực sân golf) được loại bỏ khỏi bản đồ đô thị. Chỉ số thực vật cao nhất được tính toán bằng cách so sánh chỉ

số NDVI của 23 ảnh MODIS tổ hợp 16 ngày khoảng thời gian từ 01/01/2008 đến 02/01/2009.

- Bước 4: Với từng khu vực (Đại lục Á-Âu, Châu Phi, Bắc Mỹ, Nam Mỹ và Châu Đại Dương) dựa trên các nhóm thu nhập, các ngưỡng thông số về mật độ dân số, ánh sáng ban đêm, mật độ bề mặt không thấm nước và chỉ số thực vật được tính toán dựa trên các ảnh vệ tinh có độ phân giải cao Landsat ETM+ và Google Earth.

- Bước 5: Các bản đồ ánh sáng ban đêm và mật độ bề mặt không thấm nước được sử dụng để loại bỏ các khu vực ngoại ô, nông thôn ra khỏi bản đồ đô thị theo nguyên tắc: khu vực ngoại ô, nông thôn thông thường có tỷ lệ ánh sáng ban đêm và bề mặt không thấm nước thấp hơn khu vực đô thị.

2.2.2.2. Đánh giá kết quả.

Bộ dữ liệu bản đồ đô thị GLCNMO có độ chính xác tương đối tốt trên phạm vi toàn cầu. Tuy nhiên tại các khu vực phát triển như Châu Âu, một vài khu vực đô thị với nhiều cây xanh, hoặc nằm sát công viên bị loại bỏ khỏi bản đồ đô thị. Ngược lại ở khu vực đang phát triển như Châu Á hay Châu Phi, một vài thành phố nhỏ cũng không được thể hiện.

2.3. Cải tiến phương pháp GLCNMO để phát hiện sự mở rộng đô thị tại Việt Nam

2.3.1. Thu thập dữ liệu đầu vào cho phương pháp phân loại lớp phủ đô thị của GLCNMO đối với Việt Nam.

2.3.2. Trích xuất bản đồ lớp phủ mặt đất của Việt Nam trên nền GLCNMO.

2.3.3. Phát triển phương pháp phân loại lớp phủ đô thị cho Việt Nam trên cơ sở kế thừa và cải tiến thuật toán của GLCNMO

Tập mẫu được lấy dựa trên phương pháp lấy mẫu ngẫu nhiên theo lớp (stratified random sampling) đối với các điểm ảnh không phải đô thị và lấy mẫu có hệ thống (systematic sampling) đối với các điểm ảnh thuộc lớp đô thị [189] với cùng độ phân giải 500m. Để tính toán được các ngưỡng phù hợp, các đa giác mẫu chứa các vùng đô thị được lựa chọn trên toàn bộ lãnh thổ Việt Nam, với 100 đa giác được lấy mẫu. Các điểm ảnh thuộc lớp đô thị được lấy mẫu nằm trong các đa giác đã được lựa chọn. Các điểm ảnh thuộc các lớp khác được lấy ngẫu nhiên trên toàn bộ lãnh thổ Việt Nam bằng các công cụ trong ArcGIS và Python, để đảm bảo việc lấy mẫu là chính xác và không phụ thuộc vào đối tượng lấy mẫu. Số lượng các điểm ảnh đối với từng lớp (ngoại trừ lớp đô thị) được tính dựa trên tỷ lệ các lớp trên bản đồ của GLCNMO.

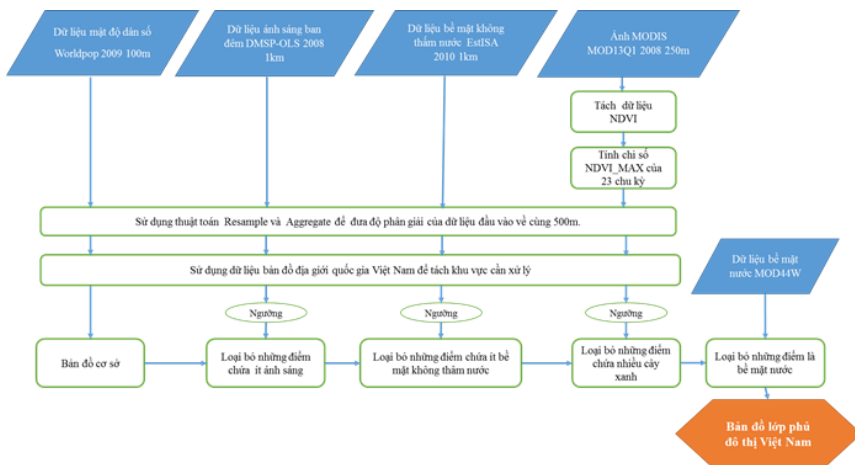
Tổng cộng có 620 điểm ảnh thuộc lớp đô thị và 1046 điểm ảnh thuộc các lớp khác được chọn để tính giá trị theo nguyên tắc: các giá trị thuộc lớp đô thị được ưu tiên cao nhất. Sau đó, các điểm ảnh này được chuyển đổi thành các shapefile nhằm mục đích so sánh với các ảnh có độ phân giải cao hơn như Google Earth và Landsat ETM+ để loại bỏ các điểm không phù hợp, kết quả là có 618 điểm đô thị và 1039 điểm thuộc các lớp khác đảm bảo yêu cầu. Các điểm này được chia thành hai tập: tập học (training set) gồm 425 điểm đô thị và 839 điểm thuộc các lớp khác, tập kiểm tra (testing set) chứa 193 điểm đô thị và 200 điểm thuộc các lớp khác.

Từ tập học, các ngưỡng phù hợp với từng dữ liệu đầu vào được tính toán dựa trên nguyên tắc: ngưỡng tốt nhất là ngưỡng có thể phân chia

nhiều nhất các điểm đô thị và các điểm thuộc lớp khác thành 2 phần tách biệt. Đầu tiên, tính toán biểu đồ tần suất (histogram) của các dữ liệu EstISA, DMSP-OLS và MOD13Q1 NDVI dựa trên tập học. Tiếp theo, các ngưỡng thích hợp của từng dữ liệu được tính toán dựa theo hàm sau:

```
thresholding(urban_histogram, non_urban_histogram, total_non_urban_points):
```

- 1: for i in range(data_size_value):
- 2: sum_urban = sum_urban + urban_histogram[i]
- 3: sum_non_urban = sum_non_urban + non_urban_histogram[i]
- 4: oa = sum_urban + (total_non_urban_points - sum_non_urban)
- 5: if oa > training_accuracy:
- 6: training_accuracy = oa
- 7: threshold = i
- 8: return threshold, training_accuracy



Phương pháp phân loại bao gồm 2 bước:

- Bước tiền xử lý dữ liệu: Các bản đồ được biến đổi về cùng độ phân giải không gian 500m và tách vùng phân tích bằng cách sử dụng bản đồ ranh giới của Việt Nam.

- Bước xử lý dữ liệu: Các bản đồ dữ liệu đầu vào được xử lý qua từng bước để tách được bản đồ lớp phủ đô thị

2.3.4. Đánh giá độ chính xác của phương pháp cải tiến

	Đối với phương pháp GLCNMO v2	Đối với phương pháp được đề xuất
Độ bao phủ	85.71%	89.29%
Độ chính xác	57%	70%
Chỉ số F1	68.47%	78.48%

CHƯƠNG 3: NGHIÊN CỨU CÁC PHƯƠNG PHÁP TÁI LẤY MẪU ĐỐI VỚI DỮ LIỆU VIỄN THÁM ĐA NGUỒN TRONG BÀI TOÁN PHÂN LOẠI LỚP PHỦ ĐÔ THỊ TẠI VIỆT NAM

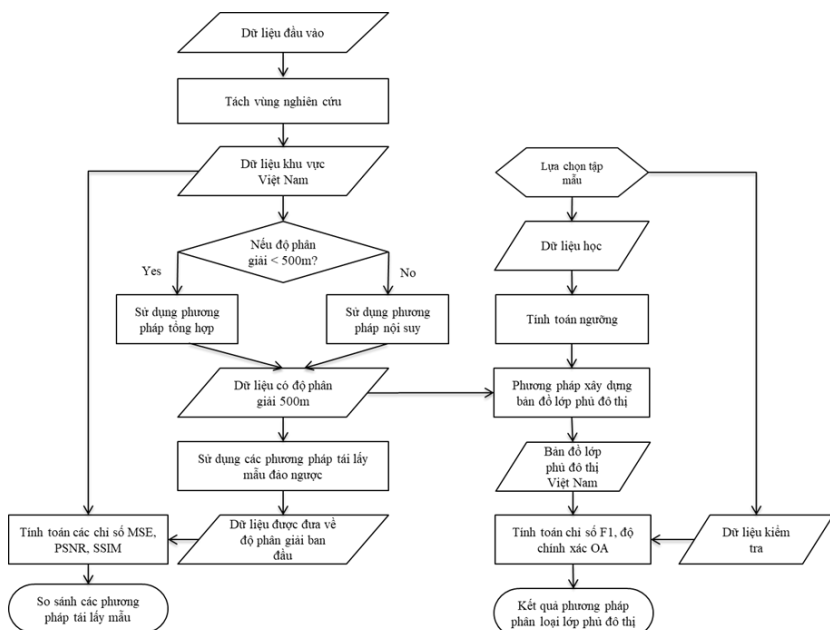
3.1. Đặt vấn đề

3.2. So sánh các phương pháp tái lấy mẫu trên dữ liệu viễn thám đa nguồn trong bài toán phân lớp đô thị tại Việt Nam

3.2.1. Dữ liệu dùng trong nghiên cứu

3.2.2. Quy trình xử lý dữ liệu

Bài toán được chia thành hai phần: Phần thứ nhất so sánh các phương pháp tái lấy mẫu, phần thứ hai đánh giá sự ảnh hưởng của các phương pháp tái lấy mẫu đến dữ liệu viễn thám dùng trong phân loại lớp phủ đô thị tại Việt Nam.



3.2.3. Đánh giá các phương pháp tái lấy mẫu

3.2.3.1. Các bước tái lấy mẫu

Các dữ liệu đầu vào được tái lấy mẫu sử dụng các phương pháp phù hợp với từng loại dữ liệu. Dữ liệu DMSP-OLS và EstISA được nội suy bằng các phương pháp nội suy láng giềng gần nhất, nội suy song tuyến tính và nội suy xoắn bậc ba. Dữ liệu MOD13Q1 NDVI được tổng hợp dựa vào các phương pháp lấy giá trị trung bình, giá trị năm giữa, giá trị cực tiểu và giá trị cực đại.

Do các dữ liệu Worldpop và MOD44W có các đặc tính riêng, vì vậy dữ liệu Worldpop được tổng hợp bằng phương pháp cộng tổng giá trị các điểm thành phần, dữ liệu MOD44W được tổng hợp bằng phương pháp dựa trên luật đa số.

Sau bước tái lấy mẫu này các dữ liệu được đưa về cùng độ phân giải 500m và sử dụng làm dữ liệu đầu vào cho bài toán phân loại lớp phủ đô thị.

Quá trình tái lấy mẫu được chia làm hai pha: trong pha đầu tiên các dữ liệu được tái lấy mẫu về cùng độ phân giải 500m, trong pha thứ hai các dữ liệu không có dữ liệu kiểm chứng có cùng độ phân giải 500m sẽ được tái lấy mẫu một lần nữa với các phương pháp ngược của pha thứ nhất để đưa về độ phân giải gốc và so sánh với dữ liệu gốc.

3.2.3.2. So sánh các phương pháp tái lấy mẫu

Do các dữ liệu Worldpop và MOD44W chỉ sử dụng một phương pháp tái lấy mẫu, vì vậy việc so sánh là không cần thiết đối với các phương pháp này.

Các dữ liệu EstISA, DMSP-OLS, và MOD13Q1 NDVI sau khi được tái lấy mẫu qua hai pha sẽ được so sánh với dữ liệu gốc bằng cách sử dụng các chỉ số đánh giá MSE, PSNR, và SSIM. Chỉ số MSE càng thấp thì càng tốt, các chỉ số PSNR và SSIM càng cao thì càng tốt.

Đối với dữ liệu EstISA, phương pháp nội suy láng giềng gần nhất kết hợp với các phương pháp tổng hợp khác cho kết quả tốt nhất với chỉ số MSE nhỏ nhất, các chỉ số PSNR và SSIM cao nhất, lần lượt là 0, $+\infty$, và 1. Trong khi đó phương pháp nội suy song tuyến tính kết hợp với các phương pháp tổng hợp khác cho kết quả kém nhất với chỉ số MSE cao nhất, các chỉ số PSNR và SSIM thấp nhất, lần lượt là 0.0026, 25.7724, và 0.9779.

Đối với dữ liệu DMSP-OLS, phương pháp nội suy láng giềng gần nhất kết hợp với các phương pháp tổng hợp khác cho kết quả tốt nhất với chỉ số MSE nhỏ nhất, các chỉ số PSNR và SSIM cao nhất, lần

lượt là $0, +\infty$, và 1 . Trong khi đó phương pháp nội suy xoắn bậc ba kết hợp với phương pháp tổng hợp giá trị cực tiểu cho kết quả kém nhất với chỉ số MSE cao nhất, các chỉ số PSNR và SSIM thấp nhất, lần lượt là 0.0112 , 19.5249 , và 0.9455 .

Đối với dữ liệu MOD13Q1, phương pháp tổng hợp giá trị trung bình kết hợp với phương pháp nội suy xoắn bậc ba cho kết quả tốt nhất với chỉ số MSE nhỏ nhất, các chỉ số PSNR và SSIM cao nhất, lần lượt là 0.0008 , 37.0509 , và 0.98 . Trong khi đó phương pháp tổng hợp giá trị cực đại kết hợp với phương pháp nội suy láng giềng gần nhất cho kết quả kém nhất với chỉ số MSE cao nhất, các chỉ số PSNR và SSIM thấp nhất, lần lượt là 0.0011 , 35.6711 , và 0.9715 .

3.2.4. Đánh giá sự ảnh hưởng của các phương pháp tái lấy mẫu trên dữ liệu viễn thám trong bài toán phân lớp đô thị tại Việt Nam

3.2.4.1. Tính toán ngưỡng cho các dữ liệu đầu vào

3.2.4.2. Phương pháp phân loại lớp phủ đô thị

Từ các dữ liệu đầu vào và các ngưỡng được tính toán phù hợp dựa trên tập học, bản đồ lớp phủ đô thị được tính toán bằng phương pháp phân loại lớp phủ đô thị GLCNMO v2.

3.2.4.3. Phương pháp đánh giá độ chính xác

3.2.4.4. Ảnh hưởng của các phương pháp tái lấy mẫu đối với việc phân loại lớp phủ đô thị tại Việt Nam

Do các dữ liệu được sử dụng làm dữ liệu đầu vào có độ phân giải tương đối thấp, bên cạnh đó số lượng của tập kiểm tra tương đối nhỏ, vì vậy độ chính xác của nhiều bản đồ đầu ra cho kết quả như nhau. Kết quả tốt nhất của tổ hợp các phương pháp tái lấy mẫu có chỉ số F1 là 0.9842 đối với sáu tổ hợp khác nhau. Ví dụ, tổ hợp các phương pháp

tái lấy mẫu bao gồm: tính tổng giá trị đối với dữ liệu Worldpop, nội suy láng giềng gần nhất đối với dữ liệu DMSP-OLS, nội suy song tuyến tính đối với dữ liệu EstISA, tổng hợp dựa trên tính toán giá trị trung bình đối với dữ liệu MOD13Q1 và phương pháp dựa trên luật đa số đối với dữ liệu MOD44W. Kết quả cũng chỉ ra rằng độ chính xác của bản đồ đầu ra phụ thuộc chủ yếu vào phương pháp tổng hợp dựa trên tính toán giá trị trung bình đối với dữ liệu MOD13Q1.

CHƯƠNG 4: XÂY DỰNG CƠ SỞ DỮ LIỆU LỚP PHỦ ĐÔ THỊ ỨNG DỤNG TRONG NGHIÊN CỨU ẢNH HƯỞNG CỦA ĐÔ THỊ HOÁ TỚI Ô NHIỄM KHÔNG KHÍ TẠI VIỆT NAM

4.1. Đặt vấn đề

4.2. Xây dựng cơ sở dữ liệu lớp phủ đô thị ứng dụng trong nghiên cứu ảnh hưởng của đô thị hoá tới ô nhiễm không khí tại Việt Nam

4.2.1 Thiết kế cơ sở dữ liệu

Để xử lý các dữ liệu viễn thám trong bài toán phân loại lớp phủ đô thị và tính toán chỉ số ô nhiễm không khí tại Việt Nam, hệ thống phân tích và xử lý dữ liệu phải đạt được các yêu cầu sau: (i) Có khả năng tính toán lớn: xử lý được các dữ liệu đa nguồn với các thuật toán tính toán phức tạp; (ii) Thông lượng truy xuất dữ liệu cao: tránh hiện tượng nút cổ chai trong quá trình xử lý; (iii) Có thể sử dụng nhiều công nghệ lập trình: phù hợp với nhiều thuật toán khác nhau; (iv) Có khả năng thực hiện đa nhiệm: cho phép nhiều thuật toán hoặc nhiều tiến trình của một thuật toán cùng xử lý; (v) Có khả năng xử lý độc lập: các công

việc xử lý cần cần được thiết kế như các tiến trình độc lập, tránh hiện tượng hệ thống bị treo.

Để đáp ứng tối đa các yêu cầu trên, kiến trúc của hệ thống được thiết kế phân tán với các mục tiêu tối ưu hóa sau: (i) Chia sẻ tài nguyên: với các yêu cầu lưu trữ và tính toán lớn, việc sử dụng một phần cứng duy nhất là không khả thi bởi hiệu năng sẽ giảm khi yêu cầu tăng cao và không tối ưu chi phí; (ii) Có khả năng mở rộng: khả năng mở rộng cao hơn nhiều lần so với sự hữu hạn của hệ thống tập trung; (iii) Tính thành phần: hệ thống phân tán có thể sử dụng các thành phần khác nhau và dễ dàng thay thế; (iv) Độ tin cậy: hệ thống phân tán cung cấp khả năng chịu lỗi, thậm chí là khả năng tự phục hồi.

Chính vì vậy, hệ thống cơ sở dữ liệu được xây dựng bao gồm hai loại cơ sở dữ liệu được sử dụng kết hợp: (i) Cơ sở dữ liệu quan hệ: các dữ liệu có cấu trúc như cây thư mục, các thông tin hỗ trợ quản lý thư mục và tệp được lưu trữ trên PostgreSQL; (ii) Cơ sở dữ liệu phi quan hệ: các dữ liệu dạng tài liệu, dữ liệu và siêu dữ liệu viễn thám được lưu trữ trên MongoDB (hệ quản trị dữ liệu kiểu NoSQL) để tăng hiệu năng của hệ thống và tương thích với công nghệ sử dụng trên máy chủ dịch vụ. Ngoài ra còn có hệ thống quản lý tệp phân tán Hadoop.

Cơ sở dữ liệu phi quan hệ bao gồm hai nhóm là:

Dữ liệu ảnh vệ tinh: được thiết kế để lưu trữ các loại dữ liệu từ ảnh vệ tinh, bao gồm các trường dữ liệu là id, sensor, level, satellite, level0, level1, level2 (lưu trữ các loại ảnh tương ứng với các cấp độ), đây là các dữ liệu đầu vào để thực hiện các thuật toán tính toán hay phục vụ cho thao tác tìm kiếm ảnh vệ tinh của người dùng.

Các dữ liệu khác: được thiết kế để lưu trữ thông tin về các loại ảnh vệ tinh, đường dẫn của các tệp ảnh và các thông tin khác của ảnh như tên vệ tinh, kiểu cảm biến, thời gian thu thập ảnh, siêu dữ liệu của ảnh.

4.2.2. Thiết kế hệ thống chức năng

Để quản lý dữ liệu và thực hiện các chức năng của hệ thống, nền tảng SEAP (Nền tảng phân tích và khám phá dữ liệu không gian lớn) được sử dụng để phân tích và lưu trữ dữ liệu

4.3. Nghiên cứu sự ảnh hưởng của quá trình phát triển đô thị đến vấn đề ô nhiễm không khí tại Việt Nam

4.3.1. Các nghiên cứu về ảnh hưởng của quá trình phát triển đô thị tới vấn đề ô nhiễm không khí tại Việt Nam

4.3.2. Dữ liệu đầu vào của nghiên cứu

4.3.3. Phương pháp nghiên cứu sự ảnh hưởng của quá trình phát triển đô thị đến môi trường tại Việt Nam

a, Phương pháp phân loại lớp phủ đô thị

b, Phương pháp tính toán nồng độ PM2.5

Trong nghiên cứu, mô hình PM2.5 tại Việt Nam được tính toán bằng phương pháp hồi quy.

Tính toán đồng nhất AOD và nhiệt độ. Mục tiêu của phần này là hợp nhất nhiều sản phẩm AOD và Nhiệt độ từ các ảnh vệ tinh để có được một bộ dữ liệu AOD nhất quán với chất lượng và độ phủ dữ liệu cao. Phương pháp Terra Regression được sử dụng để tích hợp dữ liệu aerosol (sử dụng dữ liệu Terra AOD). Phương pháp này dựa trên chất lượng dữ liệu, trong đó dữ liệu vệ tinh có chất lượng cao nhất được chọn trong tương quan với dữ liệu AERONET được sử dụng như biến phản hồi trong mô hình hồi quy tuyến tính.

Dựa trên kết quả xác thực, dữ liệu MODIS Terra AOD có chất lượng tốt nhất sẽ được chọn dưới dạng biến phản hồi. Dữ liệu AOD từ các vệ tinh MODIS Aqua và VIIRS NPP được tính toán dựa trên mô hình hồi quy theo phương trình sau:

$$AOD_{MODIS\ Terra} = a * AOD_{MODIS\ Aqua} + b \quad (4.1)$$

$$AOD_{MODIS\ Terra} = c * AOD_{NPP\ VIIRS} + d \quad (4.2)$$

$AOD_{MODIS\ Terra}$, $AOD_{MODIS\ Aqua}$ và $AOD_{NPP\ VIIRS}$ lần lượt là các giá trị AOD của các cảm biến trên các vệ tinh viễn thám MODIS Terra, MODIS Aqua và VIIRS NPP, (a,b) và (c,d) là (độ dốc và điểm chặn) của mô hình hồi quy.

Sau khi đã xây dựng được mô hình hồi quy, dữ liệu hồi quy MODIS Aqua và VIIRS NPP được tính toán như sau :

$$AOD_{MODIS\ Aqua\ regress} = a * AOD_{MODIS\ Aqua} + b \quad (4.3)$$

$$AOD_{MODIS\ VIIRS\ regress} = c * AOD_{NPP\ VIIRS} + d \quad (4.4)$$

Cuối cùng, dữ liệu hồi quy MODIS Terra, MODIS Aqua và VIIRS NPP được kết hợp bằng cách sử dụng phương pháp Maximum Likelihood.

Tương tự như dữ liệu sol khí, dữ liệu nhiệt độ từ ảnh vệ tinh cũng được kết hợp bằng phương pháp NHC MF Regression (sử dụng dữ liệu Nhiệt độ từ Trung tâm Dự báo Khí tượng Thủy văn Quốc gia Hoa Kỳ làm biến mục tiêu). Dữ liệu nhiệt độ từ vệ tinh MODIS Terra và Aqua được tính toán bằng mô hình hồi quy dưới dạng phương trình sau:

$$Temp_{NHC MF} = a * Temp_{Aqua} + b \quad (4.5)$$

$$Temp_{NHC MF} = c * Temp_{Terra} + d \quad (4.6)$$

c, Phương pháp tính toán nồng độ PM

Dữ liệu AOD và nhiệt độ từ vệ tinh đã được kiểm tra với dữ liệu AOD mặt đất từ các trạm AERONET và nhiệt độ từ các trạm CEM và các tính toán cho thấy mối tương quan cao giữa các dữ liệu AOD. Nghiên cứu trước đây cho thấy AOD và PM có mối tương quan cao với nhiệt độ. Từ đó phát triển mô hình giúp dự đoán giá trị nồng độ PM từ dữ liệu AOD và nhiệt độ đã biết. Mô hình cho thấy mối quan hệ của PM và các tham số khác thông qua một hàm f với sai số chấp nhận được ε .

$$PM=f(AOD, Temp)+\varepsilon \quad (4.7)$$

4.3.4. Đánh giá sự ảnh hưởng của quá trình phát triển đô thị đến môi trường tại Việt Nam

a, Đánh giá mở rộng của lớp phủ đô thị

Dựa trên phương pháp phân loại lớp phủ đô thị, thu được các bản đồ lớp phủ đô thị Việt Nam năm 2004, 2008, 2012 và 2015. Từ những bản đồ này, có thể thấy rằng tốc độ mở rộng các đô thị ở Việt Nam trong giai đoạn 2004-2015 là rất nhanh chóng. Khu vực đô thị tăng từ 4623 km² vào năm 2004 đến 5094 km² vào năm 2015 tức là tăng khoảng 471 km² trong giai đoạn này. Vùng lõi đô thị Việt Nam tập trung vào hai thành phố lớn nhất: thủ đô Hà Nội và thành phố Hồ Chí Minh.

b, Đánh giá sự ô nhiễm không khí

Giá trị PM_{2.5} cao chứa những thông tin về sự thay đổi quy mô khu vực, điều này đặc biệt quan trọng đối với các khu vực đô thị đông dân. Nồng độ PM_{2.5} có sự khác biệt rõ ràng giữa các khu vực thành thị và ngoài đô thị với giá trị trung bình cao của PM_{2.5} tập trung ở một số khu vực đô thị lớn như Hà Nội và thành phố Hồ Chí Minh. Khu vực

miền Bắc Việt Nam có giá trị PM2.5 cao hơn các vùng khác do ảnh hưởng của gió mùa, nguyên nhân là do gió mùa bắc/đông bắc (gió mùa đông) đã mang không khí ô nhiễm từ lục địa Trung Quốc vào lãnh thổ Việt Nam.

c, Mối tương quan giữa mở rộng đô thị và ô nhiễm không khí

Quá trình đô thị hóa nhanh chóng cùng với sự gia tăng dân số và kinh tế ở Việt Nam đã và đang ảnh hưởng đến môi trường sống, đặc biệt là ô nhiễm không khí.

Đánh giá tác động của việc mở rộng đô thị đến chất lượng không khí được thực hiện bằng cách phân tích nồng độ PM2.5 vệ tinh riêng biệt cho các khu vực thành thị và nông thôn ở Việt Nam. Có sự khác biệt rõ ràng giữa giá trị nồng độ PM2.5 trung bình của khu vực thành thị và khu vực ngoài thành thị, giá trị PM2.5 ở khu vực thành thị nhìn chung cao hơn khu vực nông thôn trong nhiều năm. Năm 2004, PM2.5 của vệ tinh trung bình ở khu vực thành thị là 25,142, trong khi mức trung bình ở khu vực nông thôn chỉ là 21,056. Trong năm 2015, giá trị PM2.5 lần lượt là 22,945 và 20,054 đối với khu vực thành thị và ngoài thành thị.

Tại hai thành phố lớn nhất của Việt Nam là Hà Nội và Thành phố Hồ Chí Minh, giá trị nồng độ PM2.5 thường cao hơn giá trị trung bình của cả nước. Ví dụ, trong năm 2015, giá trị PM2.5 của các khu vực thành thị ở Hà Nội và thành phố Hồ Chí Minh là 25,782 và 25,462, so với 22,954 của Việt Nam.

Tuy nhiên, ở các vùng nhỏ hơn như Hà Nội hoặc Thành phố Hồ Chí Minh, sự khác biệt của giá trị PM2.5 ở khu vực thành thị và ngoài đô thị là rất nhỏ. Vì vậy, có thể nói, đối với dữ liệu có độ phân giải

thấp như dữ liệu MODIS, khu vực nghiên cứu phải đủ lớn để có thể đánh giá sự khác biệt về các chỉ số ô nhiễm không khí.

KẾT LUẬN VÀ KIẾN NGHỊ VỀ NHỮNG NGHIÊN CỨU TIẾP THEO

Các dữ liệu viễn thám ngày càng đóng vai trò quan trọng trong những nghiên cứu về phương pháp phân loại lớp phủ mặt đất và lớp phủ đô thị. Trong nội dung luận án, tác giả đã nghiên cứu những nội dung lý thuyết liên quan đến dữ liệu viễn thám, các phương pháp xử lý, phân tích dữ liệu viễn thám. Dữ liệu viễn thám được biến đổi từ tín hiệu thu nhận được bởi các cảm biến đặt trên vật mang, qua quá trình tiền xử lý sẽ trở thành dữ liệu đầu vào của các bài toán khác nhau trong đó có phương pháp phân loại lớp phủ đô thị. Trong quá trình tiền xử lý dữ liệu, nhiều phương pháp tái lấy mẫu khác nhau được sử dụng như các phương pháp tổng hợp hoặc nội suy ảnh viễn thám do sự khác biệt về kiểu dữ liệu, độ phân giải không gian của từng loại dữ liệu. Kết quả thu được của các phương pháp phân loại lớp phủ đô thị là các bản đồ đô thị khu vực Việt Nam. Trên cơ sở các bản đồ nền này, luận án đã đánh giá những ảnh hưởng của quá trình phát triển đô thị tới vấn đề ô nhiễm không khí tại Việt Nam.

Những đóng góp khoa học chính của luận án bao gồm:

- Cải tiến phương pháp phân loại lớp phủ toàn cầu của GLCNMO nhằm xây dựng bản đồ lớp phủ đô thị cho Việt Nam. Các nội dung cải tiến đó là: (i) đề xuất một phương pháp lấy mẫu ngẫu nhiên dựa trên việc tổng hợp các bộ dữ liệu lớp phủ toàn cầu, kết hợp với các ảnh có độ phân giải cao như Google Earth và Landsat ETM+ và công cụ trong

ArcGIS và Python, (ii) đề xuất phương pháp tính ngưỡng dựa trên histogram của tập mẫu. Kết quả đánh giá cho thấy việc hiệu chỉnh dữ liệu đầu vào, lựa chọn giá trị ngưỡng phù hợp với các thông số thu thập tại Việt Nam giúp nâng cao độ chính xác của dữ liệu lớp phủ đô thị tại Việt Nam.

- Đánh giá sự ảnh hưởng của quá trình tái lấy mẫu tới chất lượng của ảnh viễn thám và tác động của quá trình này đến độ chính xác của phương pháp phân loại lớp phủ mặt đất tại Việt Nam.

- Đánh giá sự ảnh hưởng của quá trình phát triển đô thị tới vấn đề ô nhiễm không khí tại Việt Nam. Nghiên cứu dựa trên dữ liệu viễn thám và dữ liệu thu thập được tại các trạm mặt đất để xây dựng bản đồ đô thị và bản đồ PM_{2.5} tại Việt Nam trong các năm 2004, 2008, 2012 và 2015. Trên cơ sở đó, luận án đã đạt được các kết quả sau: (i) ước tính sự mở rộng đô thị của Việt Nam từ năm 2004 đến năm 2015; (ii) Phân tích tình trạng ô nhiễm không khí của Việt Nam từ năm 2004 đến năm 2015; (iii) Tìm hiểu tác động của việc mở rộng đô thị đối với ô nhiễm không khí ở Việt Nam từ năm 2004 đến năm 2015.

Tuy nhiên, các nghiên cứu trong luận án vẫn còn một vài hạn chế đó là:

- Do các dữ liệu đầu vào có độ phân giải tương đối thấp, vì vậy bản đồ lớp phủ đô thị thu được vẫn thiếu các thông tin chi tiết và chứa các điểm khó có thể phân biệt rạch ròi là đô thị hay không phải đô thị. Bên cạnh đó việc xây dựng tập học và tập kiểm tra chưa đủ lớn cũng ảnh hưởng đến việc phân biệt các điểm của kết quả đầu ra. Chính các yếu tố này đã đem đến những ảnh hưởng không mong muốn đến độ chính xác của dữ liệu đầu ra.

- Kết quả đầu ra của phương pháp phân loại lớp phủ đô thị là những bản đồ nền các khu vực đô thị của Việt Nam. Các bản đồ này có thể được ứng dụng trong nhiều lĩnh vực khác nhau như quản lý sử dụng đất, quy hoạch đô thị, thúc đẩy phát triển kinh tế - xã hội,... Tuy nhiên trong nội dung luận án mới chỉ đề cập tới khía cạnh nghiên cứu ảnh hưởng của quá trình đô thị hoá tới vấn đề ô nhiễm không khí tại Việt Nam.

Từ những hạn chế trên, một số hướng nghiên cứu tiếp theo của nghiên cứu sinh đó là:

- Tích hợp các nguồn dữ liệu viễn thám có độ phân giải cao như ảnh Landsat 8, Sentinel 2, ảnh radar, ảnh Lidar, không ảnh, ... để có thể thu được các bản đồ lớp phủ đô thị có độ phân giải cao hơn, chi tiết hơn.

- Nghiên cứu, ứng dụng các bản đồ lớp phủ đô thị trong các lĩnh vực liên ngành khác nhau.

DANH MỤC CÔNG TRÌNH KHOA HỌC CỦA TÁC GIẢ LIÊN QUAN ĐẾN LUẬN ÁN

1. Dung, P. T., Chuc, M. D., Thanh, N. T. N., Hung, B. Q., & Chung, D. M. (2016). Optimizing GLCNMO version 2 method to detect Vietnam's urban expansion. 2016 Eighth International Conference on Knowledge and Systems Engineering (KSE), 309–314. <https://doi.org/10.1109/KSE.2016.7758072>. (Scopus index).

2. Dung, P. T., Chuc, M. D., Thanh, N. T. N., Hung, B. Q., & Chung, D. M. (2019). Comparison of Resampling Methods on Different Remote Sensing Images for Vietnam's Urban Classification. Research and Development on Information and Communication Technology. <https://doi.org/10.32913/rd-ict.vol2.no15.663>.

3. Pham, T. D., Pham, V. H., Luu, Q. T., Ngo, X. T., Nguyen, T. N. T., & Bui, Q. H. (2019). Analyzing the impacts of urban expansion on air pollution in Vietnam using the SEAP platform. IOP Conference Series: Earth and Environmental Science. <https://doi.org/10.1088/1755-1315/266/1/012008>. (Scopus index).

4. Bui, Q. H., Luu, Q. T., Ha, D. Van, Pham, T. D., Praseuth, S., & Laffly, D. (2020). Spatial Data Infrastructure. In TORUS 2 – Toward an Open Resource Using Services (Vol. 7, pp. 247–261). (Book chapter) <https://doi.org/10.1002/9781119720553.ch9>. (Wiley publishing).

Danh mục này gồm 04 công trình./.

