

# TRÍCH XUẤT THÔNG TIN NGUỒN NƯỚC BỀ MẶT TRONG ĐÔ THỊ TẠI THÀNH PHỐ HỒ CHÍ MINH TỪ ẢNH SENTINEL-2 SỬ DỤNG THUẬT TOÁN PHÂN CỤM K-MEANS

PHAN THỊ ANH THU<sup>1</sup>, NGUYỄN TUẤN KIỆT<sup>2</sup>, TRẦN THÀNH LONG<sup>3</sup>

<sup>1</sup>Khoa Kỹ thuật Xây dựng, Trường Đại học Bách khoa TP.HCM

<sup>2</sup>Đại học Quốc gia Thành phố Hồ Chí Minh

## Tóm tắt:

Sự hiện diện của các nguồn nước bề mặt như sông, hồ và các nguồn nước khác có vai trò quan trọng trong việc điều hòa nhiệt độ và duy trì cân bằng sinh thái trong hệ thống đô thị. Việc quan trắc và đánh giá chính xác sự phân bố của các nguồn nước bề mặt trong đô thị, đặc biệt với độ chính xác cao, trở thành một yếu tố quan trọng để quản lý hiệu quả môi trường đô thị. Trong bối cảnh biến đổi khí hậu đang gia tăng, đô thị đang phải đối mặt với nhiều thách thức, bao gồm hiện tượng đảo nhiệt đô thị. Sự hiện diện của các nguồn nước bề mặt có thể giúp giảm thiểu hiện tượng này bằng cách tạo ra các vùng mát mẻ và cân bằng nhiệt độ. Điều này không chỉ mang lại lợi ích cho môi trường mà còn góp phần vào sự bền vững và an toàn của đô thị trong tương lai. Nghiên cứu này trình bày một thuật toán tự động trích xuất thông tin các nguồn nước bề mặt đô thị. Phương pháp đề xuất đã được kiểm tra trên dữ liệu Sentinel-2 với độ phân giải không gian 20 mét và có thể áp dụng trên phạm vi rộng lớn. Phương pháp sử dụng thuật toán phân cụm K-means để tiến hành phân loại tự động các ảnh được tính theo chỉ số NDWI, MNDWI và tỷ số tích hợp. Nghiên cứu đã được thực hiện tại khu vực Thành phố Hồ Chí Minh. Số lượng lớp tối ưu cho K-means trong theo phương pháp đề xuất là 6. Kết quả đánh giá độ chính xác cho thấy rằng phương pháp đề xuất phù hợp để trích lọc thông tin các nguồn nước bề mặt nhanh chóng và chính xác.

Từ khóa: Sentinel -2, K-means, nước mặt đô thị.

## 1. Giới thiệu chung

Theo ước tính 70.9% diện tích bề mặt Trái đất được bao phủ bởi nước. Nước là một tài nguyên tự nhiên quý giá được tồn tại dưới nhiều hình thức và phân bố ở nhiều khu vực khác nhau trên thế giới. Nước đa phần được lưu trữ tại các đại dương. Nước tại lục địa chủ yếu tập trung vào nước mặt, nước ngầm và nước đóng băng tại các khu vực lạnh giá như vùng cực hay đỉnh núi cao. Sự khai thác quá

mức nguồn nước cùng với tác động do thay đổi trong sử dụng đất và biến đổi khí hậu ảnh hưởng tiêu cực đến chu trình nước. Nước đóng vai trò quan trọng đối với hệ sinh thái đô thị và khí hậu đô thị đặc biệt trong việc giảm thiểu tác động của đảo nhiệt đô thị [1-2]. Do đó thông tin về phân bố của các dòng sông và hồ theo không gian và thời gian là yếu tố quan trọng để hiểu về chu trình nước và là cơ sở để tìm các giải pháp khắc phục những vấn đề như

thoát nước tại khu vực đô thị. Ước tính đáng tin cậy về nước mặt rất quan trọng đối với các lĩnh vực khoa học khác nhau.

Hiện nay, truy cập và sử dụng dữ liệu vệ tinh đã trở nên đơn giản hơn. Dữ liệu này cho phép chúng ta trích xuất và giám sát tài nguyên nước mặt để đáp ứng nhu cầu của con người và cung cấp thông tin cho các cơ quan quản lý để thúc đẩy sự bền vững. Thay đổi về diện tích của các nguồn nước có thể được phát hiện bằng cách so sánh hình ảnh vệ tinh từ các giai đoạn thời gian khác nhau hoặc sử dụng các thuật toán phân loại đa dạng [3]. Bằng cách phân tích các đặc trưng phản xạ từ ảnh vệ tinh đa phổ, chúng ta có thể cải thiện độ chính xác của việc tự động trích xuất thông tin về nước mặt đô thị. Để đáp ứng nhu cầu ngày càng tăng về phân tích thông tin về bề mặt nước, đã có nhiều phương pháp được đề xuất dựa trên thuật toán phân loại, bao gồm cả phân loại giám định và phi giám định trực tiếp từ ảnh tổ hợp màu, sử dụng các chỉ số tính từ kênh phổ, hoặc dựa vào các phương pháp ngưỡng để trích xuất nguồn nước từ hình ảnh vệ tinh. Đối với việc tự động hóa trích xuất thông tin về bề mặt từ ảnh vệ tinh, nhiều tác giả đã sử dụng các phương pháp học máy khác nhau, như random forest [4], support vector machines [5]. Ngoài ra, các phương pháp phân loại không giám định không đòi hỏi mẫu đào tạo và phù hợp hơn để phát triển các thuật toán tự động. Phương pháp K-means là một trong những phương pháp phân loại phi giám định phổ biến được sử dụng rộng rãi trong nhiều mục đích nghiên cứu, đặc biệt trong lĩnh vực tự động hóa trích xuất thông tin địa lý [6].

Nhiều nhà nghiên cứu đã đề xuất các chỉ số phổ khác nhau, chẳng hạn như chỉ số khác biệt chuẩn hóa nước [7], chỉ số khác biệt nước được sửa đổi (MNDWI) [8], và chỉ số khai

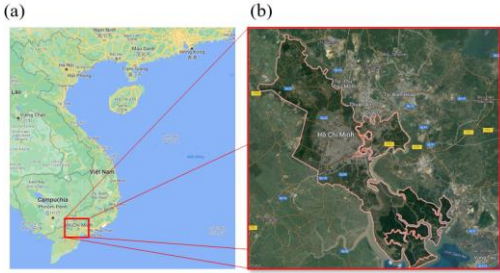
thác nước tự động (AWEI) [9] để xác định vùng nước mặt. Mặc dù nhiều tác giả đã nghiên cứu, tuy nhiên các thuật toán được đề xuất gần như luôn dựa trên các phương pháp phân loại giám định yêu cầu dữ liệu thực tế từ mặt đất, kiến thức và công sức chuyên gia. Điều này gây cản trở rất lớn trong việc phát triển các thuật toán hoàn toàn tự động.

Mục tiêu chính của nghiên cứu này là kiểm tra việc sử dụng một thuật toán phi giám định cho việc tự động trích thông tin bề mặt từ ảnh vệ tinh Sentinel-2. Nội dung cụ thể của phương pháp đề xuất bao gồm kiểm tra việc sử dụng các chỉ số nước để làm nổi bật vùng nước mặt đô thị. Từ đó, sử dụng thuật toán phân loại phi giám định K-means để tự động phân nhóm. Cuối cùng là đánh giá kết quả tách bề mặt nước dựa trên phương pháp phân loại K-means từ giá trị các chỉ số nước khác nhau.

## **2. Khu vực và dữ liệu nghiên cứu**

Xem xét tầm quan trọng của các nguồn nước bề mặt đô thị đối với chất lượng cuộc sống trong các thành phố lớn, và để đánh giá tính hiệu quả của phương pháp đề xuất, Thành phố Hồ Chí Minh được lựa chọn làm khu vực nghiên cứu (Hình 1). Thành phố Hồ Chí Minh là một trung tâm quan trọng về kinh tế, giải trí, văn hóa và giáo dục tại Việt Nam. Thành phố này thuộc loại đô thị đặc biệt của Việt Nam và nằm ở vùng chuyển tiếp giữa Đông Nam Bộ và Tây Nam Bộ. Với 16 quận, 1 thành phố và 5 huyện, diện tích của thành phố là 2.095 km<sup>2</sup>. Mật độ dân số trung bình ở đây là 4.375 người/km<sup>2</sup>. Do đó mật độ xây dựng ở khu vực này rất cao. Ngoài ra, dân số đông còn gây ra áp lực rất lớn đối với việc sử dụng tài nguyên nước tại khu vực này. Trong nghiên cứu này, bảy hình ảnh Sentinel-2 với độ che phủ mây ít hơn 15% được sử dụng để bao quát hết tòa bộ

phạm vi của thành phố (Bảng 1). Các hình ảnh Sentinel-2 đã được tải xuống thông qua Cổng thông tin Copernicus Open Access Hub (<https://scihub.copernicus.eu/dhus/>). Tất cả các hình ảnh Sentinel-2 được thu thập bằng vệ tinh Sentinel-2 và tải xuống dưới dạng sản phẩm cấp độ L2A cung cấp giá trị phản xạ đầy của bầu khí quyển tương ứng từng kênh phổ (Bảng 2).



Hình 1: Khu vực TP. HCM. (a) Vị trí địa lý của TP. HCM và (b) hình ảnh vệ tinh tại khu vực TP. HCM

Bảng 1: Danh sách ảnh Sentinel 2

STT	Ngày thu thập	Độ phủ mây
1	17/01/2023	9.25
2	09/02/2023	14.02
3	16/02/2023	7.63
4	08/03/2023	0.13
5	08/03/2023	0.23
6	04/02/2023	12.4
7	24/02/2023	0.02

Bảng 2: Đặc trưng của các kênh phổ

Kênh phổ	Bước sóng trung tâm (nm)		
	Sentinel-2A	Sentinel-2B	Độ phân giải (m)
Band 1	442.7	442.2	60
Band 2	492.4	492.1	10
Band 3	559.8	559.0	10
Band 4	664.6	664.9	10
Band 5	704.1	703.8	20
Band 6	740.5	739.1	20
Band 7	782.8	779.7	20
Band 8	832.8	832.9	10
Band 8A	864.7	864.0	20
Band 9	945.1	943.2	60
Band 10	1373.5	1376.9	60
Band 11	1613.7	1610.4	20
Band 12	2202.4	2185.7	20

### 3. Nguồn dữ liệu thực nghiệm

#### 3.1. Chỉ số nhận diện bề mặt nước

Để tăng cường sự thể hiện của một đối tượng so với các đối tượng khác trên ảnh việc sử dụng các chỉ số được tính từ giá trị phản xạ phổ là điều cần cần thiết. Căn cứ vào đặc trưng phản xạ phổ của đối tượng nước bề mặt, McFeeters (1996) đề xuất chỉ số chênh lệch chuẩn hóa của nước (NDWI), và Xu (2006) đề xuất sử dụng chỉ số chênh lệch chuẩn hóa nước được sửa đổi (MNDWI), tương ứng. Các chỉ số này dựa trên đặc trưng phản xạ phổ của kênh xanh lá và đặc trưng hấp thụ của kênh hồng ngoại sóng ngắn đối với đối tượng mặt nước. Trong nghiên cứu này, cả hai chỉ số đều được sử dụng để phục vụ công tác xác định vùng có nước phủ mặt. Bên cạnh sử dụng chỉ số nước khác biệt bình quân (NDWI) và chỉ số khác biệt nước bình quân hiệu chỉnh (MNDWI), việc lập tỷ số tích hợp giữa kênh hồng ngoại gần, hồng ngoại sóng ngắn và kênh xanh lá cũng được sử dụng để tăng cường sự hiển thị đường mép nước của khu vực nghiên cứu.

$$NDWI = \frac{Green - NIR}{Green + NIR} \quad (1)$$

$$MNDWI = \frac{Green - SWIR1}{Green + SWIR2} \quad (2)$$

$$Tỷ số = \left( \frac{NIR * SWIR1}{Green^2} \right) \quad (3)$$

#### 3.2. Ghép ảnh

Do khu vực thành phố Hồ Chí Minh được thể hiện trong phạm vi của nhiều tấm ảnh Sentinel-2. Do đó, sau khi tải ảnh về cần tiến hành ghép ảnh. Quá trình ghép ảnh được thực hiện tự động do các ảnh đã được nắn chỉnh hình học ở hệ tọa độ WGS 84, múi chiếu 48. Sau khi tiến hành ghép ảnh, khu vực Thành phố Hồ Chí Minh được cắt ảnh tự động dựa

trên dữ liệu ranh giới hành chính ở định dạng shape file (Hình 1b).

### **3.3. Giải thuật phân cụm K-means**

Mục tiêu chính của nghiên cứu này là phát triển thuật toán phi giám định và hoàn toàn tự động để trích thông tin bề mặt nước đô thị từ dữ liệu Sentinel-2 sử dụng giải thuật K-means. Thuật toán K-means là một phương pháp phân cụm dữ liệu, trong đó mỗi cụm dữ liệu được đặc trưng bởi một tâm. Tâm đại diện cho một cụm và có giá trị bằng trung bình của tất cả các quan sát nằm trong cụm đó. Thuật toán này sử dụng khoảng cách từ mỗi quan sát tới các tâm để xác định nhãn cho chúng, thuộc về cụm nào gần nhất. Ban đầu thuật toán sẽ khởi tạo ngẫu nhiên một số lượng xác định trước tâm cụm. Sau đó tiến hành xác định nhãn cho từng điểm dữ liệu và tiếp tục cập nhật lại tâm cụm. Thuật toán sẽ dừng cho tới khi toàn bộ các điểm dữ liệu được phân về đúng cụm hoặc số lượt cập nhật tâm chạm ngưỡng. Dữ liệu đầu vào cho thuật toán lần lượt là giá trị NDWI, MNDWI và tỷ số tích hợp.

Số lượng nhóm là thông tin rất quan trọng đối với phương pháp phân loại K-means. Vì vậy, cần tiến hành xác định số lượng nhóm tối ưu cho giải thuật này. Phương pháp Elbow được sử dụng để lựa chọn số lượng cụm phù hợp dựa trên biểu đồ hàm biến dạng. Để thực hiện điều này, chúng ta huấn luyện thuật toán K-means với một loạt các giá trị khác nhau cho số lượng cụm và tính toán hàm biến dạng cho từng trường hợp. Hàm biến dạng thể hiện mức độ biến đổi của dữ liệu khi chia thành các cụm khác nhau. Trên biểu đồ hàm biến dạng, chúng ta tìm điểm gọi là khuỷu tay (elbow point), đây là điểm mà sau đó tốc độ giảm của hàm biến dạng trở nên không đáng kể. Điều này ngụ ý rằng số lượng cụm tại điểm này có thể coi là lựa chọn tốt nhất cho việc phân chia

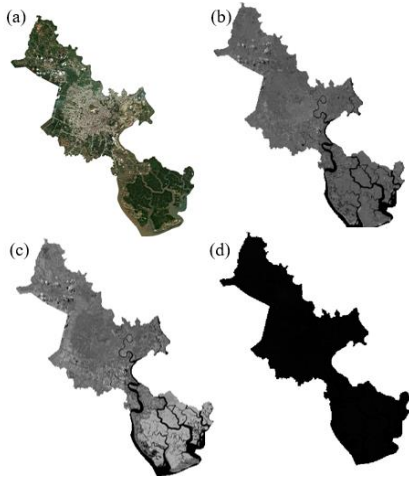
dữ liệu. Căn cứ vào vị trí điểm elbow, nhóm tác giả tiến hành lựa chọn số lượng cụm được phân cho phù hợp với dữ liệu.

### **3.4. Đánh giá kết quả**

Độ chính xác của kết quả tách vùng nước mặt được đánh giá dựa trên sự so sánh kết quả sử dụng ba loại chỉ số tách mặt nước khác nhau. Tất cả các pixel trong khu vực nghiên cứu được sử dụng để đánh giá độ chính xác. Dựa trên đồ thị elbow cho từng loại dữ liệu số lượng nhóm tối ưu được xác định cho khu vực Thành phố Hồ Chí Minh. Phân lớp nước được sử dụng trong đánh giá độ chính xác đã được trích xuất từ từng loại ảnh chỉ số khác nhau bằng giải thuật K-means. Cuối cùng, các ảnh được chồng lên nhau để tiến hành đánh giá kết quả. Thông tin lớp nước mặt cũng được chồng lớp trực tiếp lên ảnh tổ màu màu tự nhiên của ảnh vệ tinh để tiến hành đánh giá trực quan.

## **4. Kết quả**

Chỉ số tính từ giá trị phản xạ của các kênh phổ có thể cung cấp thông tin đặc trưng cho một số đối tượng nhất định. Tuy nhiên trong một số trường hợp các đối tượng khác cũng có giá trị tương đồng. Chỉ số NDWI, N\MNDWI và tỷ số tích hợp dựa trên đặc trưng phản xạ của nước và các đối tượng không phải là nước trên từng kênh ảnh. Về cơ bản, việc tính giá trị các chỉ số này sẽ cung cấp kết quả là một tờ ảnh cấp độ xám với các giá trị của mỗi pixel là giá trị đặc trưng tổng hợp của toàn bộ đối tượng trong phạm vi của pixel đó (Hình 2). Đối với khu vực của pixel tổng hợp chứa cả nước và các đối tượng khác thì giá trị được tính có thể không thể hiện rõ nét sự khác biệt của các đối tượng. Quá trình phân nhóm theo K-means cân nhắc đến cả yếu tố giá trị lẫn mối tương quan về vị trí của các pixel. Do đó, phương pháp này được kỳ vọng mang lại độ chính xác cao.

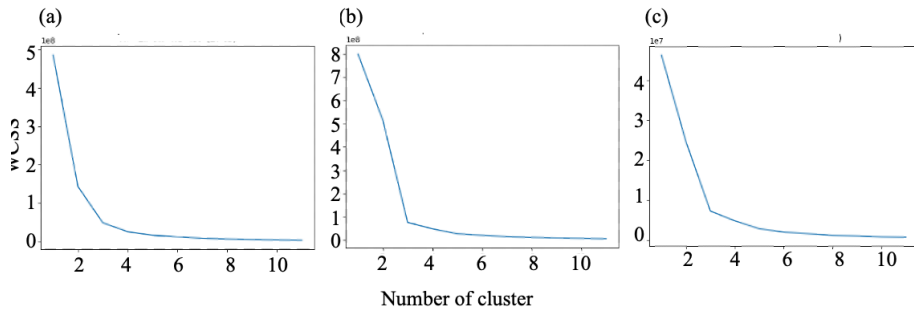


Hình 2: Kết quả các ảnh chỉ số. (a) Ảnh tổ hợp màu tự nhiên, (b) Ảnh NDWI, (c) ảnh MDWI và (d) ảnh tỷ số tích hợp.

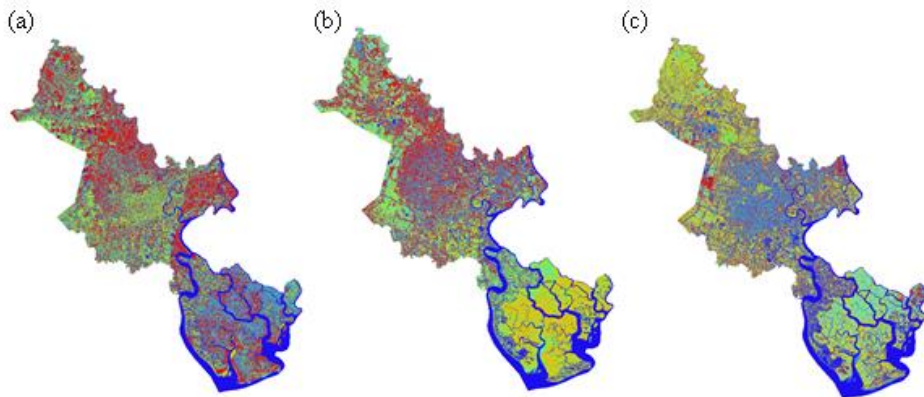
Kết quả cho thấy các ảnh cấp độ xám được tạo thành từ các giá trị chỉ số NDWI, MNDWI cung cấp sự thể hiện nổi bật đối với các vùng nước bề mặt. Cụ thể, vùng mặt nước sẽ có giá trị tối màu hơn các vùng khác. Đối với trường hợp ảnh tỷ số, yếu tố đường mép nước được

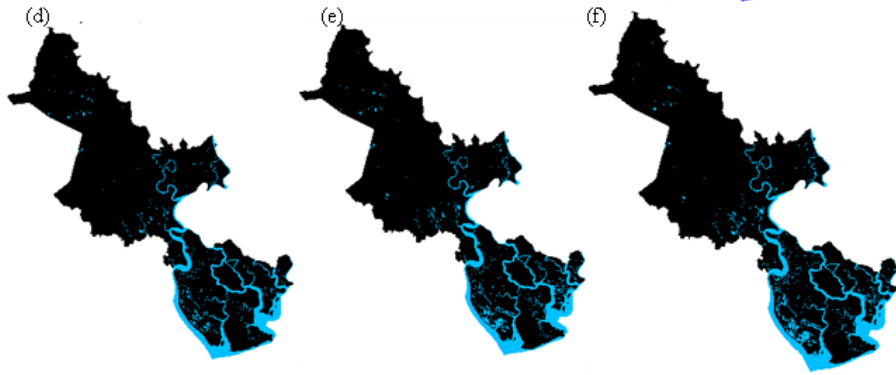
ưu tiên thể hiện rõ nét hơn. Vùng đường mép nước thể hiện sự chuyển tiếp giữa mặt nước và vùng đất liền nên các giá trị thể hiện đường biên sẽ có giá trị nhỏ hơn 1.

Dựa trên kết quả phân tích lược đồ elbow cho khu vực thành phố Hồ Chí Minh thấy số lượng nhóm thích hợp cho thuật toán K-mean trong cả 3 trường hợp là 6 nhóm (Hình 3). Kết quả phân nhóm được thể hiện trong Hình 4. Chương trình phát triển sẽ gán màu cho các phân lớp với màu của lớp nước là xanh biển. Từ kết quả phân loại có thể nhận thấy các chỉ số được đề xuất chỉ thể hiện sự nhạy cảm đối với đối tượng mặt nước. Đối với các đối tượng khác kết quả phân nhóm của của ba chỉ số là khác nhau (Hình 4 a-c). Vùng bề mặt nước được thể hiện rõ nét trong cả ba trường hợp (Hình 4 e-f). Trong trường hợp ảnh tỷ số sự khác biệt không thể phân biệt bằng mắt thường nhưng kết quả phân loại vẫn thành công (Hình 2d và Hình 4c).



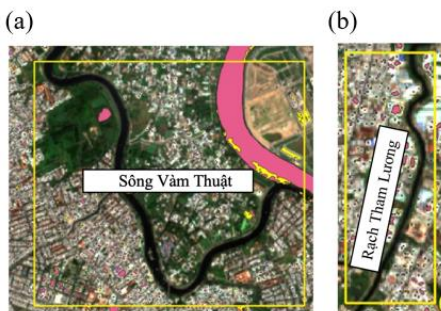
Hình 3: Biểu đồ Elbow. (a) NDWI, (b) MNDWI và (c) tỷ số tích hợp





Hình 4: Kết quả tách bề mặt nước đô thị từ giá trị (a,d) NDWI, (b, e) MNDWI và (c,f) tỷ số tích hợp. Các hàng thể hiện (a-c) kết quả phân nhóm theo k-means và vùng nước được tách thành công (d-f)

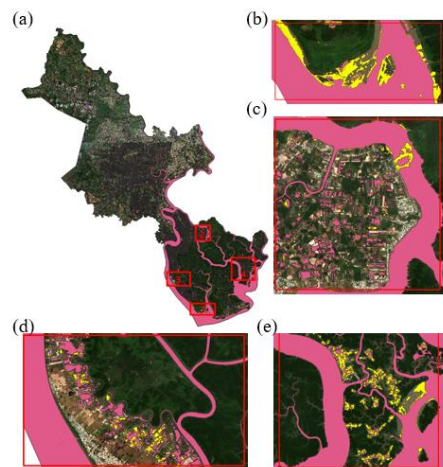
Để tiến hành đánh giá độ chính xác, vùng nước mặt được tách dựa trên 3 chỉ số khác nhau được chồng lớp lên dữ liệu ảnh tổ hợp màu tự nhiên của ảnh vệ tinh để tiến hành đánh giá khách quan. Kết quả cho thấy, các tuyến sông lớn và sông chính được trích lọc thành công. Tuy nhiên, không phải tất cả các vùng nước đều được trích xuất thành công. Khu vực sông Vàm Thuật bị ô nhiễm nghiêm trọng và chuyển màu đen thì không thể được nhận biết trong cả ba trường hợp được đề cập trong nghiên cứu. Kết quả này có thể được ứng dụng trong quan trắc ô nhiễm nguồn nước. Bên cạnh đó, một số tuyến kênh nhỏ cũng không được nhận diện do hạn chế về độ phân giải không gian của ảnh sử dụng (Hình 5).



Hình 5: Vùng nước không thể trích xuất do chất lượng nước và kích thước dòng kênh

Để tiến hành so sánh kết quả của ba phương pháp, vùng nước mặt được trích xuất từ ba phương pháp được chồng lớp với nhau

(Hình 6a). Nhằm hiển thị trực quan các màu riêng biệt được gán cho từng lớp giá trị. Cụ thể, màu hồng là kết quả trích xuất mặt nước từ NDWI, màu vàng là kết quả của ảnh tỷ số, và màu xanh lá cây là kết quả từ MNDWI. Kết quả cho thấy, chỉ số NDWI và MNDWI cho kết quả tương đồng khi chồng lớp vì không nhận thấy sự khác biệt. Sự khác biệt gây ra bởi ảnh tỷ số tại một số khu vực có bãi bồi là do đặc tính của ảnh tỷ số khi tính từ kênh SWIR1 nhạy cảm với các vùng đất trống có độ ẩm cao như các bãi bồi hoặc khu vực đất canh tác nông nghiệp ngập nước (Hình 6b-e).



Hình 6: Kết quả chồng lớp của 3 phương pháp. màu hồng là kết quả trích xuất mặt nước từ NDWI, màu vàng là kết quả của ảnh tỷ số, và màu xanh lá cây là kết quả từ MNDWI

Nhìn chung, việc sử dụng hình ảnh vệ tinh có thể trích lọc thông tin bề mặt đất nhanh chóng. Kỹ thuật này đã được ứng dụng rộng rãi trên thế giới trong vài thập kỷ qua. Bản đồ bề mặt nước thành phố được tạo lập từ ảnh với độ phân giải cao có thể giúp các nhà nghiên cứu và người làm chính sách trong việc ra quyết định tốt hơn và nhanh chóng hơn trong mục đích giám sát hệ sinh thái thành thị, và đưa ra các chính sách giảm tác động của các đảo nhiệt đô thị. Trong nghiên cứu này độ phân giải của ảnh vệ tinh là 20 m nên nhiều vùng kênh nhỏ không thể được nhận biết. Bên cạnh đó chất lượng màu nước cũng ảnh hưởng đến kết quả. Tuy nhiên, dựa trên các chỉ số đặc trưng được tính cho từng loại lớp phủ, việc trích xuất thông tin từ ảnh sẽ dễ dàng hơn trong phạm vi rộng lớn. Việc này đóng góp rất lớn cho công tác quản lý cũng như giám sát sự biến động của nguồn nước mặt.

### **5. Kết luận**

Các nguồn dữ liệu vệ tinh miễn phí đóng một vai trò quan trọng trong việc thu thập thông tin về tình hình mặt đất tổng quan và đặc biệt là dữ liệu về mặt nước. Trong nghiên cứu được trình bày ở đây, chúng tôi sử dụng dữ liệu vệ tinh Sentinel-2 để trích xuất thông tin về mặt nước tại Thành phố Hồ Chí Minh. Quá trình trích xuất thông tin này được thực hiện tự động dựa trên thuật toán phân loại phi giám định K-means từ các giá trị chỉ số NDWI, MNDWI, và tỷ số tích hợp. Kết quả cho thấy rằng số lớp tối ưu để phân loại cho khu vực này là 6. Tuy nhiên, quá trình trích xuất thông tin từ dữ liệu ảnh vệ tinh vẫn phụ thuộc vào chất lượng của nước. Khu vực nước ô nhiễm và có màu đen như sông Vàm Thuật không thể được nhận diện từ ảnh dựa trên bất kỳ chỉ số nào được sử dụng trong nghiên cứu này. Nhìn chung, kết quả trích xuất thông tin về mặt

nước từ cả ba chỉ số đều cho thấy sự tương đồng đáng kể. Tuy nhiên, chỉ số tỷ số tích hợp tương quan cả với vùng đất có độ ẩm cao. Điều này dẫn đến sự khác biệt đáng kể trong kết quả của chỉ số tỷ số tích hợp ở khu vực cửa sông so với hai chỉ số còn lại. Nói chung, phương pháp phân loại phi giám định có nhiều lợi thế trong việc tự động trích xuất thông tin về mặt nước. Tuy nhiên, cần phải xác định các chỉ số phù hợp cho quá trình phân loại. Đối với các ảnh tỷ số khó nhận diện bằng mắt thường, phương pháp K-means vẫn cho kết quả phân loại thành công. Việc lựa chọn chỉ số để tiến hành phân loại nên phụ thuộc vào đối tượng cần trích xuất. Việc tích hợp các thông tin khác như bản đồ sử dụng đất, cao độ địa hình và các yếu tố địa lý khác có thể cải thiện độ chính xác của quá trình phân loại và trích xuất thông tin về mặt nước. Các yếu tố này có thể cung cấp thông tin bổ sung về môi trường đô thị và hệ thống nước bề mặt. ○

### **Lời cảm ơn**

Chúng tôi xin cảm ơn Trường Đại học Bách khoa, ĐHQG-HCM đã hỗ trợ cho nghiên cứu này.

### **Tài liệu tham khảo**

- [1]. Xiang X, Li Q, Khan S, Khalaf OI. 2021. Urban water resource management for sustainable environment planning using artificial intelligence techniques. *Environ Impact Assess Rev.* 86:106515.
- [2]. Nwakaire CM, Onn CC, Yap SP, Yuen CW, Onodagu PD. 2020. Urban heat island studies with emphasis on urban pavements: a review. *Sustain Cities Soc.* 63:102476.
- [3]. Yang L, Driscoll J, Sarigai S, Wu Q, Lippitt CD, Morgan M. *Towards Synoptic Water Monitoring Systems: A Review of AI Methods for Automating Water Body*

Detection and Water Quality Monitoring Using Remote Sensing. *Sensors*. 2022; 22(6):2416.

<https://doi.org/10.3390/s22062416>

[4]. Wangchuk S, Bolch T. 2020. Mapping of glacial lakes using Sentinel-1 and Sentinel-2 data and a random forest classifier: strengths and challenges. *Sci Remote Sens*. 2:100008

[5]. Liu Q, Huang C, Shi Z, Zhang S. 2020. Probabilistic river water mapping from Landsat-8 using the support vector machine method. *Remote Sens*. 12(9):1374.

[6]. Hartigan JA, Wong MA. 1979. Algorithm AS 136: a k-means clustering algorithm. *J Royal Statist Soc. Series c (Applied Statistics)*. 28(1):100–108.

### **Summary**

### **Extracting surface water in urban areas in Ho Chi Minh City from sentinel-2 images using the k-means clustering algorithm**

*Phan Thi Anh Thu, Nguyen Tuan Kiet, Tran Thanh Long*

*Civil Engineering Faculty, Ho Chi Minh City University of Technology*

*Viet Nam National University of Ho Chi Minh City*

The existence of surface water reservoirs, such as rivers, lakes, and other aquatic bodies, plays a pivotal role in temperature regulation and the preservation of ecological equilibrium within urban settings. Precise monitoring and evaluation of the spatial distribution of these surface water features in urban areas, particularly with high precision, have become indispensable for effective urban environmental management. In the face of escalating climate change, cities grapple with various challenges, including the urban heat island effect. The presence of surface water reservoirs can ameliorate this phenomenon by creating cooler zones and maintaining temperature equilibrium. This not only benefits the environment but also contributes to the future sustainability and safety of urban regions. This study introduces an automated algorithm designed to extract information about urban surface water. The proposed methodology has been assessed using Sentinel-2 data at a spatial resolution of 20 meters and applies to large areas. It leverages the K-means clustering algorithm to autonomously categorise images based on NDWI, MNDWI, and ratio indices. The research was conducted in Ho Chi Minh City, and the optimal number of clusters for K-means in the proposed approach is determined to be 6. The results of the accuracy assessments affirm that the proposed method is well-suited for rapid and precise extraction of data concerning surface water reservoirs.○

Keywords: Sentinel -2, K-means, urban surface water.

[7]. McFeeters SK. 1996. The use of the Normalized Difference Water Index (NDWI) in the delineation of open water features. *Int J Remote Sens*. 17(7):1425–1432.

[8]. Feyisa GL, Meilby H, Fensholt R, Proud SR. 2014. Automated water extraction index: a new technique for surface water mapping using Landsat imagery. *Remote Sens. Environ*. 140:23–35.

[9]. Xu H. 2006. Modification of Normalized Difference Water Index (NDWI) to Enhance Open Water Features in Remotely Sensed Imagery. *International Journal of Remote Sensing* 27(14):3025–3033. DOI: 10.1080/0143116060058917.○