

Phát triển mô hình trí tuệ nhân tạo kết hợp DC-CSFT trong thành lập bản đồ dự báo không gian sạt lở đất tại Quốc lộ 6, tỉnh Hòa Bình, Việt Nam

■ **TS. PHẠM THÁI BÌNH; TS. NGÔ QUỐC TRINH**
THS. NGUYỄN ĐỨC ĐẰM; TS. BÙI THỊ QUỲNH ANH
 Trường Đại học Công nghệ Giao thông vận tải

TÓM TẮT: Ở Việt Nam, sạt lở đất là một trong những hiểm họa thiên tai gây thiệt hại lớn về con người và tài sản, đặc biệt tại các khu vực miền núi. Việc phát triển các mô hình dự báo mới và xây dựng các bản đồ dự báo không gian sạt lở đất tin cậy là cần thiết và hữu ích để giảm nguy cơ sạt lở đất và lập kế hoạch phát triển các khu vực đồi núi. Trong nghiên cứu này, mô hình trí tuệ nhân tạo kết hợp DC-CSFT đã được phát triển và sử dụng để thành lập bản đồ dự báo không gian sạt lở đất tại QL6, tỉnh Hòa Bình, Việt Nam. Hiệu suất của mô hình được đánh giá thông qua các chỉ số thống kê khác nhau bao gồm diện tích dưới đường cong ROC (AUC). Kết quả cho thấy rằng, mô hình DC-CSFT có hiệu suất cao trong việc dự báo không gian sạt lở đất tại khu vực nghiên cứu (AUC = 0,831). Cách tiếp cận này cũng có thể được áp dụng ở các khu vực đồi núi dễ bị sạt lở khác tại Việt Nam để mang lại hiệu quả cao trong công tác phòng ngừa và quản lý sạt lở đất.

TỪ KHÓA: Trí tuệ nhân tạo, sạt lở đất, dự báo không gian, DC-CSFT, Quốc lộ 6, Hòa Bình.

ABSTRACT: In Vietnam, landslide is one of the most serious natural hazards which caused a huge loss of deaths and properties, especially in mountainous areas. Development of new prediction models and construction of reliable landslide spatial prediction maps are essential and useful for reduction of landslide risks and making better land use planning in the affected areas. In this study, a hybrid artificial intelligence model namely DC-CSFT was developed and applied in landslide spatial prediction mapping at national road NO6, Hoa Binh province, Vietnam. Prediction capability of the model was validated using various statistical indexes including area under the ROC curve (AUC). The results show that the DC-CSFT model has a high performance in landslide spatial prediction at the study area (AUC = 0.831). This approach can be also applied in other landslide prone areas of Vietnam for better landslide hazard management.

KEYWORDS: Artificial intelligence, landslide, spatial prediction, DC-CSFT, National road 6, Hoa Binh province.

1. ĐẶT VẤN ĐỀ

Sạt lở đất là một trong những thảm họa thiên nhiên nguy hiểm với tính mạng con người, tài sản, cơ sở hạ tầng. Đây là hiện tượng đất đá chuyển dịch từ phía đỉnh dốc về chân dốc theo các cơ chế và tốc độ khác nhau. Nguyên nhân gây sạt lở đất có thể do yếu tố tự nhiên (mưa, điều kiện địa chất, thảm thực vật...) hay do yếu tố con người (phá rừng, xây dựng công trình...). Sạt lở đất tập trung nhiều ở vùng núi, trung du dốc và thường xảy ra sau các trận mưa kéo dài hoặc bão lũ, hoặc động đất. Các bản đồ dự báo không gian sạt lở đất ở các khu vực nhạy cảm dễ xảy ra sạt lở đất đóng vai trò quan trọng đồng thời đưa ra các phương án quy hoạch, xây dựng phát triển cơ sở hạ tầng hay giải pháp phòng chống phù hợp nhằm giảm thiểu thiệt hại do hiện tượng này gây ra.

Các phương pháp thống kê truyền thống, chẳng hạn như mô hình phân cấp thứ bậc [1], mô hình tỷ số tần suất [2] thường được sử dụng trong phân tích dự báo không gian sạt lở đất. Ngày nay, với sự tiến bộ của khoa học công nghệ các phương pháp trí tuệ nhân tạo trong đó có các mô hình học máy, chẳng hạn như mô hình cây hồi quy, rừng ngẫu nhiên, hồi quy logistics [3], mạng nơ-ron nhân tạo (ANN) [4] đã và đang được sử dụng trong các nghiên cứu dự báo không gian sạt lở đất cho thấy sự hiệu quả với độ chính xác cao. Ngoài ra, các mô hình trí tuệ nhân tạo lai - các mô hình kết hợp giữa các kỹ thuật tối ưu hóa và kỹ thuật phân loại cũng đã được phát triển và thể hiện tính ưu việt hơn so với các mô hình đơn trong dự báo không gian sạt lở đất [5, 6].

Trong nghiên cứu này, mục tiêu chính của nghiên cứu là phát triển mô hình trí tuệ nhân tạo kết hợp DC-CSFT để thành lập bản đồ dự báo không gian sạt lở đất dọc theo QL6 đi qua địa hình đồi núi của tỉnh Hòa Bình, Việt Nam. Mô hình kết hợp DC-CSFT được kết hợp bởi hai kỹ thuật bao gồm kỹ thuật tối ưu hóa Decorate (DC) và kỹ thuật phân loại chi phí nhạy cảm (CSFT). Dữ liệu của nghiên cứu được xử lý trên ứng dụng ArcGIS và quá trình mô hình hóa được thực hiện sử dụng công cụ Weka.

2. KHU VỰC NGHIÊN CỨU

QL6 (QL ND-6), tỉnh Hòa Bình, Việt Nam được lựa chọn làm khu vực nghiên cứu, nằm giữa các vĩ độ của 20°19' 21°08' N và kinh độ 104°48' - 105°40' E. (Hình 3.1). QL NH-6 dài 504 km nối Hà Nội với các tỉnh miền núi Tây Bắc của Việt Nam. Đoạn tuyến QL6 đi qua tỉnh Hòa Bình dài khoảng 115 km từ km38 - km153 qua 6 khu vực: Lương Sơn, Kỳ Sơn, Tân Hòa Bình, Cao Phong, Tân Lạc, Mai Châu. Khu vực dọc theo

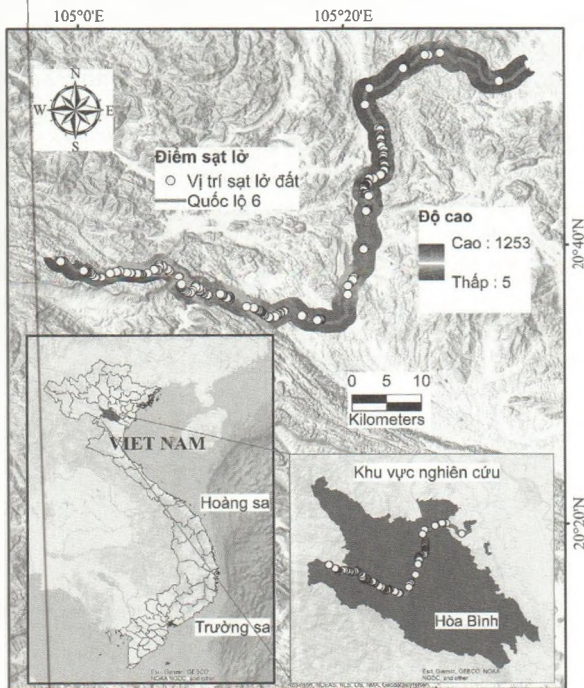
và xung quanh đường này, tại tỉnh Hòa Bình, nơi bị ảnh hưởng chủ yếu bởi sạt lở đất, đã được lựa chọn là khu vực nghiên cứu cho nghiên cứu hiện tại. Địa hình của khu vực là đồi núi với các thung lũng hẹp. Độ cao từ 0 đến 1.163 m. Đoạn km38 - km90 chạy qua núi và đồi thấp ở phía Đông Nam ở độ cao từ 0 - 500 m; từ km91 - km 153 chạy qua vùng núi cao phía Tây Bắc ở độ cao từ 600 đến 1.163 m (Hình 3.1).

3. THU THẬP VÀ PHÂN TÍCH DỮ LIỆU

Dữ liệu được sử dụng trong nghiên cứu dự báo không gian sạt lở đất bao gồm có hiện trạng sạt lở đất và các tham số ảnh hưởng tới sạt lở đất.

3.1. Hiện trạng sạt lở đất khu vực nghiên cứu

Hiện trạng sạt lở đất khu vực nghiên cứu được xây dựng từ dữ liệu khảo sát hiện trường và giải đoán ảnh Google Earth [6]. Có tổng cộng 235 vụ sạt lở đất trong quá khứ đã được thu thập trong giai đoạn từ năm 2016 đến 2018 dọc theo tuyến Quốc lộ NH-6 (Hình 3.1). Trong đó, 219 vụ sạt lở đất (93,2%) được phát hiện nằm trong phạm vi 100 m từ tìm đường và 16 vụ sạt lở đất (16,8%) nằm trong phạm vi 200 - 300 m từ tìm đường. Hầu hết các vụ sạt lở đất tại khu vực khảo sát là các vị trí mất ổn định nông và thường xảy ra vào mùa mưa. Bên cạnh hiện trạng sạt lở đất, các vị trí không sạt lở cũng được nhận diện dựa vào phân tích địa hình địa mạo khu vực nghiên cứu, số lượng các vị trí không sạt lở được xác định bằng với số lượng vị trí sạt lở được nhận diện. Để ứng dụng mô hình dự báo không gian sạt lở đất, cơ sở dữ liệu bao gồm dữ liệu đào tạo (70% dữ liệu hiện trạng) và dữ liệu kiểm chứng (30% dữ liệu hiện trạng) được xây dựng. Trong đó, các lớp dữ liệu hiện trạng sạt lở đất được mã hóa thành "1" và lớp không sạt lở được mã hóa thành "0".



Hình 3.1: Vị trí khu vực nghiên cứu và hiện trạng sạt lở đất

3.2. Các yếu tố ảnh hưởng đến sạt lở đất

Các yếu tố ảnh hưởng đến sạt lở đất được xác định dựa trên giả thuyết rằng các vụ sạt lở đất trong tương lai sẽ

xảy ra trong điều kiện tương tự như các vụ sạt lở đất trong quá khứ [6]. Trong nghiên cứu này, 12 yếu tố ảnh hưởng đến sạt lở đất được sử dụng cho nghiên cứu mô hình dự báo [6], cụ thể là: (1) Góc mái dốc, (2) Hướng mái dốc, (3) Hình dáng bề mặt địa hình, (4) Độ cao địa hình, (5) Bao phủ thực vật (NDVI), (6) Địa chất, (7) Khoảng cách đứt gãy, (8) Tích lũy dòng chảy, (9) Sức mạnh dòng chảy (SPI), (10) Độ ẩm địa hình (TWI), (11) Khoảng cách đến sông suối, (12) Khoảng cách đến đường giao thông. Trong đó, các yếu tố địa hình địa mạo được lấy từ Mô hình Độ cao kỹ thuật số (DEM) của khu vực nghiên cứu được tạo từ dữ liệu trái đất ALOS PALSAR (<https://search.asf.alaska.edu/>) ở độ phân giải không gian 12,5 m. Các yếu tố ảnh hưởng khác như TWI, SPI, STI đã được tạo ra từ DEM. Các dữ liệu còn lại về địa chất, khoảng cách đứt gãy... được thu thập từ bản đồ liên quan của Bộ Tài nguyên và Môi trường. Các lớp thông tin của các bản đồ thành phần được mã hóa thành các số tương ứng với số lớp của mỗi tham số trong cơ sở dữ liệu dùng cho mô hình dự báo.

4. PHƯƠNG PHÁP NGHIÊN CỨU

Để dự báo không gian sạt lở đất, mô hình lai kết hợp có tên là D-CSF được sử dụng trong nghiên cứu này. Mô hình DC-CSFC là mô hình kết hợp từ hai phương pháp là phương pháp phân loại chi phí nhạy cảm (a Cost-Sensitive Filter Tree - CSFT) và phương pháp tập hợp Decorate (DC). Trong đó, kỹ thuật Decorate là kỹ thuật tối ưu hóa dữ liệu đầu vào để tạo ra được tập dữ liệu đào tạo tối ưu được sử dụng để phân loại các lớp sạt lở và các lớp không sạt lở sử dụng kỹ thuật phân loại CSFT. Các thông tin mô tả chi tiết về hai phương pháp DC và CSFT được mô tả trong các nghiên cứu đã công bố [7, 8].

Để đánh giá hiệu suất của mô hình dự báo, các kỹ thuật đánh giá định lượng như diện tích dưới đường cong ROC (AUC) [9] và các chỉ số thống kê được sử dụng bao gồm: giá trị dự đoán dương (PPV), giá trị dự đoán âm (NPV), độ nhạy (SST), độ đặc hiệu (SPF), độ chính xác (ACC), chỉ số Kappa (K), lỗi bình phương trung bình gốc (RMSE), lỗi tuyệt đối (MAE) [3, 4]. Giá trị AUC thay đổi giữa "0,5 - 1". Giá trị AUC càng gần 1 độ chính xác của thuật toán càng cao, trong khi gần 0,5 độ chính xác của thuật toán càng thấp hơn [5]. Chỉ báo K là một biện pháp thống kê hiệu quả giúp đo lường sự đồng thuận ngẫu nhiên giữa các yếu tố phân loại. K thay đổi giữa 1 và 0. Nếu các giá trị K gần gũi với 1, nó cho thấy độ tin cậy cao và độ tin cậy của thuật toán trong việc dự đoán sự nhạy cảm sạt lở đất. Tiêu chí ACC ước tính tỷ lệ hoặc dự báo chính xác để dự báo toàn bộ sạt lở [10, 11]. RMSE cho biết sự khác biệt giữa dữ liệu được quan sát và dữ liệu ước tính. MAE là một phạm vi lỗi giữa các quan sát nhị phân. Các giá trị cao hơn của SPF, PPV, NPV, ACC, SST, K và các giá trị thấp hơn của RMSE và MAE cho biết hiệu suất cao hơn của mô hình trong việc dự đoán sự nhạy cảm sạt lở đất.

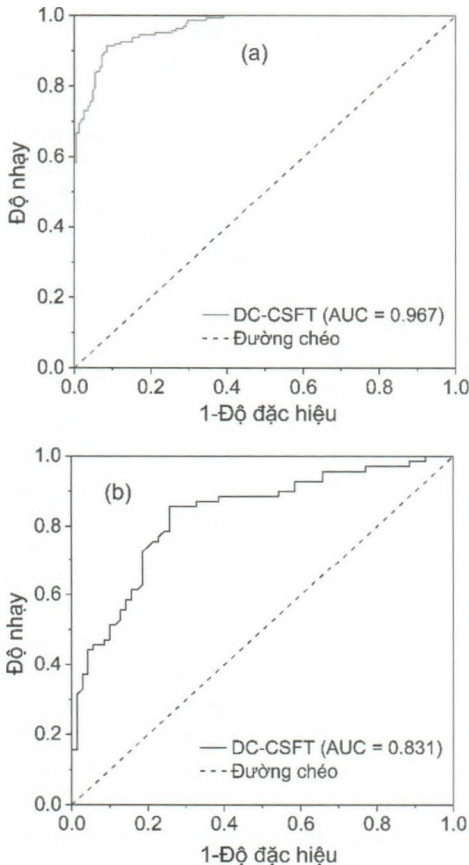
5. KẾT QUẢ VÀ THẢO LUẬN

5.1. Đánh giá hiệu suất của mô hình dự báo

Kết quả đánh giá hiệu suất của mô hình dự báo DC-CSFT trên hai tập dữ liệu đào tạo và dữ liệu kiểm chứng

được thể hiện trên Hình 5.1, Hình 5.2 và Bảng 5.1. Kết quả đánh giá mô hình sử dụng kỹ thuật đường cong ROC cho thấy, diện tích dưới đường cong khi mô hình sử dụng dữ liệu đào tạo là $AUC = 0,967$ và dữ liệu kiểm chứng là $AUC = 0,831$, kết quả này cho thấy hiệu suất dự báo của mô hình DC-CSFT là tốt. Bảng 5.1 và Hình 5.2 thể hiện các giá trị của các chỉ số thống kê đánh giá hiệu suất của mô hình dự báo, cho thấy các giá trị của các chỉ số trên cả hai tập dữ liệu đào tạo và kiểm chứng đều thể hiện rằng năng lực dự báo của mô hình DC-CSFT là tốt.

Kết quả đánh giá hiệu suất của mô hình cho thấy mô hình lai kết hợp DC-CSFT có năng lực dự báo không gian sạt lở đất tốt, có độ chính xác cao và có thể sử dụng trong việc xây dựng bản đồ dự báo không gian sạt lở đất khu vực nghiên cứu. Kết quả của nghiên cứu này cũng khẳng định và phù hợp với kết quả của các nghiên cứu trước đó [6], [3, 4] rằng mô hình trí tuệ nhân tạo, đặc biệt là các mô hình lai kết hợp là các mô hình dự báo mới, có hiệu quả trong dự báo không gian sạt lở đất.

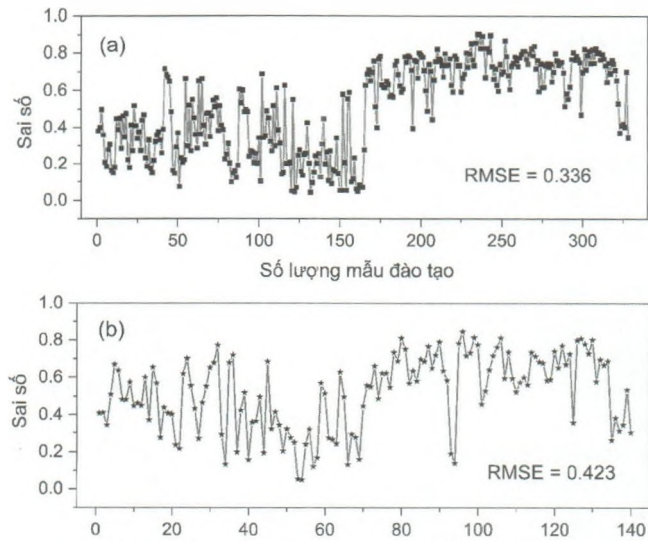


Hình 5.1: Giá trị AUC của mô hình DC-CSFT sử dụng: a) - Dữ liệu đào tạo, b) - Dữ liệu kiểm chứng

Bảng 5.1. Hiệu suất của mô hình DC-CSFT sử dụng các chỉ số thống kê

STT	Tham số	Dữ liệu đào tạo	Dữ liệu kiểm chứng
1	PPV (%)	84,85	70,00
2	NPV (%)	93,87	85,71
3	SST (%)	93,33	83,05
4	SPF (%)	85,96	74,07

STT	Tham số	Dữ liệu đào tạo	Dữ liệu kiểm chứng
5	ACC (%)	89,33	77,86
6	K	0,79	0,56
7	MAE	0,30	0,39



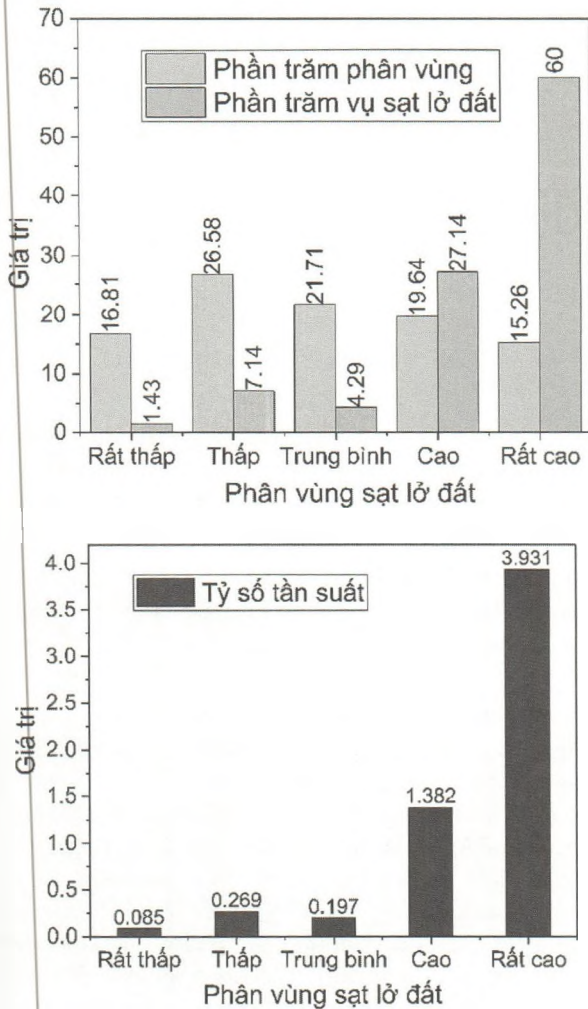
Hình 5.2: Giá trị lỗi bình phương trung bình gốc (RMSE) của mô hình DC-CSFT sử dụng dữ liệu đào tạo (a) và dữ liệu kiểm chứng (b)

5.2. Xây dựng bản đồ dự báo không gian sạt lở đất

Bản đồ dự báo không gian sạt lở đất khu vực nghiên cứu được thành lập dựa trên kết quả đào tạo mô hình DC-CSFT được thể hiện trên Hình 5.3. Bản đồ được xây dựng với 5 cấp độ về sắc xuất xảy ra sạt lở đất: rất cao, cao, trung bình, thấp và rất thấp tương ứng với các giá trị xác suất xảy ra sạt lở đất được xác định từ mô hình dự báo. Trong đó, phương pháp phân loại “điểm nghi tự nhiên” trong phần mềm ArcGIS 10,8 đã được sử dụng để phân chia các lớp [12]. Để đánh giá độ tin cậy của bản đồ dự báo không gian sạt lở đất khu vực nghiên cứu, 30% các vụ sạt lở đất chưa được sử dụng để chồng lán lên các lớp bản đồ và xác định tỷ số tần suất xuất hiện (Hình 5.4). Kết quả đánh giá cho thấy, hầu hết các vụ sạt lở đất trong quá khứ xảy ra tại khu vực có xác suất xảy ra sạt lở đất rất cao (3,931) và cao (1,382) và có rất ít các vụ sạt lở đất xảy ra trên các lớp thấp (0,269) và lớp rất thấp (0,085). Kết quả này khẳng định, bản đồ dự báo không gian sạt lở đất được xây dựng từ mô hình DC-CSFT có độ tin cậy rất cao.



Hình 5.3: Bản đồ dự báo không gian sạt lở đất khu vực nghiên cứu sử dụng mô hình DC-CSFT



Hình 5.4: Phân trăm các vụ sạt lở đất và tỷ số tần suất của các lớp nhạy cảm sạt lở đất

6. KẾT LUẬN

Trong nghiên cứu này, mô hình trí tuệ nhân tạo lai kết hợp DC-CSFT đã được phát triển và ứng dụng trong dự báo không gian sạt lở đất dọc theo tuyến QL6, tỉnh Hòa Bình, Việt Nam. Trong đó, mô hình DC-CSFT là mô hình kết hợp giữa hai kỹ thuật CSFT và DC là các kỹ thuật trí tuệ nhân tạo tiên tiến. Cơ sở dữ liệu dùng cho mô hình dự báo bao gồm 235 vụ sạt lở đất trong quá khứ và 12 tham số ảnh hưởng tới sạt lở đất khu vực nghiên cứu. Kết quả nghiên cứu cho thấy rằng, mô hình DC-CSFT có năng lực dự báo tốt ($AUC = 0,83$). Bản đồ dự báo không gian sạt lở đất đã được xây dựng có độ tin cậy cao thể hiện rằng có 27,14% và 60% khu vực nghiên cứu nằm trong khu vực có xác suất xảy ra sạt lở đất rất cao và cao. Kết quả nghiên cứu cũng khẳng định rằng, mô hình trí tuệ nhân tạo trong đó có mô hình DC-CSFT là các công cụ hữu ích có độ chính xác cao, có thể được sử dụng trong dự báo không gian sạt lở đất tại các khu vực chịu ảnh hưởng bởi thiên tai sạt lở đất.

Lời cảm ơn: Nhóm tác giả trân trọng cảm ơn sự chia sẻ dữ liệu từ TS. Hà Thị Hằng, Trường Đại học Xây dựng Hà Nội để phục vụ cho nghiên cứu này.

Tài liệu tham khảo

- [1]. Đỗ, M.N., M.Đ.J.V.J.o.S.E. Đỗ and E. Sciences (2016), *Ứng dụng GIS và phương pháp phân tích thứ bậc (AHP) thành lập bản đồ nguy cơ trượt lở huyện Xín Mần, tỉnh Hà Giang, Việt Nam*, 32(2S).
- [2]. Thanh, D.Q., et al. (2020.), *GIS based frequency ratio method for landslide susceptibility mapping at Da Lat City, Lam Dong province, Vietnam*, 42(1), pp.55-66.
- [3]. Nhu, V.-H., et al. (2020), *Landslide Detection and Susceptibility Modeling on Cameron Highlands (Malaysia): A Comparison between Random Forest, Logistic Regression and Logistic Model Tree Algorithms*, 11(8): p.830.
- [4]. Lee, D.-H., Y.-T. Kim and S.-R.J.R.S. Lee (2020), *Shallow Landslide Susceptibility Models Based on Artificial Neural Networks Considering the Factor Selection Method and Various Non-Linear Activation Functions*, 12(7): p.1194.
- [5]. Chen, W., et al. (2019), *Spatial Prediction of Landslide Susceptibility Using GIS-Based Data Mining Techniques of ANFIS with Whale Optimization Algorithm (WOA) and Grey Wolf Optimizer (GWO)*, Applied Sciences, 9(18), p.3755.
- [6]. Hang, H.T., et al. (2021), *Spatial prediction of landslides along National Highway-6, Hoa Binh province, Vietnam using novel hybrid models*, pp.1-26.
- [7]. Siers, M.J. and M.Z. Islam (2015), *Cost-Sensitive Decision Forest: CSForest*.
- [8]. Sun, B., H. Chen and J.J.K.-B.S. Wang (2015), *An empirical margin explanation for the effectiveness of DECORATE ensemble learning algorithm*, 78, pp.1-12.
- [9]. Avand, M., et al. (2020), *A tree-based intelligence ensemble approach for spatial prediction of potential groundwater*, International Journal of Digital Earth, 13(12): p.1408-1429.
- [10]. Prävälje, R. and R. Costache (2014), *The analysis of the susceptibility of the flash-floods' genesis in the area of the hydrographical basin of Bâsca Chiojduului river*, Forum geografic. XIII(1), pp.39-49.
- [11]. De Rosa, P., A. Fredduzzi and C. Cencetti (2019), *Stream Power Determination in GIS: An Index to Evaluate the Most 'Sensitive' Points of a River*, Water, 11(6), p.1145.
- [12]. Roy, S., et al. (2021), *Coastal erosion risk assessment in the dynamic estuary: The Meghna estuary case of Bangladesh coast*, International Journal of Disaster Risk Reduction, 61, p.102364.

Ngày nhận bài: 17/5/2022

Ngày chấp nhận đăng: 20/6/2022

Người phản biện: TS. Lý Hải Bằng

TS. Đỗ Minh Ngọc