

# Ứng dụng học máy nhận dạng tướng đất đá khu vực bồn trũng An Châu

Doãn Ngọc San<sup>1\*</sup>, Nguyễn Văn Thắng<sup>2</sup>, Hoàng Hữu Hiệp<sup>2</sup>, Nguyễn Anh Hào<sup>2</sup>, Trần Thị Oanh<sup>1</sup>, Nguyễn Thị Minh Hồng<sup>3</sup>

<sup>1</sup>Trường Đại học Dầu khí Việt Nam, 762 Cách Mạng Tháng Tám, phường Long Toàn, TP Bà Rịa, tỉnh Bà Rịa - Vũng Tàu, Việt Nam

<sup>2</sup>Công ty TNHH MTV Dầu khí Sông Hồng, 117 Trần Duy Hưng, phường Trung Hòa, quận Cầu Giấy, Hà Nội, Việt Nam

<sup>3</sup>Trường Đại học Mỏ - Địa chất, 18 phố Viên, phường Đức Thắng, quận Bắc Từ Liêm, Hà Nội, Việt Nam

Ngày nhận bài 22/2/2023; ngày chuyển phân biện 24/2/2023; ngày nhận phân biện 7/3/2023; ngày chấp nhận đăng 11/3/2023

## Tóm tắt:

Bồn trầm tích An Châu (bồn An Châu) là một cấu trúc địa chất kéo dài theo phương Tây Bắc - Đông Nam, phân bố ở vùng Đông Bắc Việt Nam. Bồn An Châu được cho là bồn có quy mô lớn và có tiềm năng dầu khí. Thực tế này cho thấy, nhiều khả năng sẽ phát hiện được các tích tụ dầu khí trong bồn An Châu thuộc địa phận của Việt Nam. Mặc dù tiềm năng dầu khí của bồn này được nhận định từ rất sớm, nhưng vì nhiều lý do mà cho đến nay công tác điều tra, khảo sát địa chất và thăm dò ở khu vực bồn An Châu còn rất sơ sài, chưa đáp ứng được các yêu cầu của công tác tìm kiếm - thăm dò dầu khí. Việc ứng dụng mô hình học máy (machine learning - ML) vào nhận dạng tướng đá là một phương pháp mới giúp giảm thiểu thời gian xử lý, tổng hợp cơ sở dữ liệu về cả số lượng và định dạng, phát hiện các mối quan hệ ẩn sâu giữa các lớp thông tin nhận dạng. Mục tiêu chính của nghiên cứu này là nhận dạng tướng đất đá khu vực bồn An Châu từ phần số liệu đầy đủ đã được huấn luyện bởi cấu trúc mạng cây quyết định (DT) kết hợp thuật toán gradient boosting (XGB) để đánh giá cấu trúc và xác định tiềm năng dầu khí khu vực này. Điều kiện tiên quyết để nâng cao độ chính xác của học máy là phải làm giàu cơ sở dữ liệu thông qua tích hợp số liệu địa chất - địa chấn và tính toán thêm các thuộc tính để xây dựng “mô hình học” - huấn luyện ML và sử dụng kết quả huấn luyện đó để nhận dạng tướng đất đá trong khu vực bồn An Châu.

**Từ khóa:** cây quyết định, cơ sở dữ liệu địa chất - địa vật lý, học máy, trí tuệ nhân tạo.

**Chỉ số phân loại:** 1.2, 1.8, 2.1

## Application of machine learning for facies recognition in An Chau basin

Ngoc San Doan<sup>1\*</sup>, Van Thang Nguyen<sup>2</sup>, Huu Hiep Hoang<sup>2</sup>, Anh Hao Nguyen<sup>2</sup>, Thi Oanh Tran<sup>1</sup>, Thi Minh Hong Nguyen<sup>3</sup>

<sup>1</sup>Petroleum Vietnam University (PVU), 762 Cach Mang Thang Tam Street, Long Toan Ward, Ba Ria City, Ba Ria - Vung Tau Province, Vietnam

<sup>2</sup>PVEP Song Hong, 117 Tran Duy Hung Street, Trung Hoa Ward, Cau Giay District, Hanoi, Vietnam

<sup>3</sup>Hanoi University of Mining and Geology, 18 Vien Street, Duc Thang Ward, Bac Tu Liem District, Hanoi, Vietnam

Received 22 February 2023; revised 7 March 2023; accepted 11 March 2023

## Abstract:

The An Chau basin (An Chau basin) is a geological structure extending in a northwest-southeast direction, located in the northeastern region of Vietnam. The An Chau basin is considered a large-scale basin with significant petroleum potential. This fact suggests a high likelihood of discovering hydrocarbon accumulations within the An Chau basin within Vietnam's territory. Although the petroleum potential of this basin was identified early, for various reasons, geological surveys and exploration activities in the An Chau basin remain rudimentary and have not yet met the requirements for effective petroleum exploration and production. The application of machine learning (ML) models to lithofacies identification offers a novel approach to reducing processing time, consolidating large and diverse datasets, and uncovering hidden relationships among layers of identification information. The primary objective of this study is to identify lithofacies in the An Chau basin using comprehensive datasets trained with a decision tree (DT) structure combined with a gradient boosting algorithm (XGB) to evaluate the structure and assess the petroleum potential of this area. A prerequisite for improving the accuracy of machine learning is enriching the database through the integration of geological and seismic data, along with the calculation of additional attributes to build the “learning model” training the ML system and applying the trained model to identify lithofacies in the An Chau basin.

**Keywords:** artificial intelligence, decision tree, geophysical database, machine learning.

**Classification numbers:** 1.2, 1.8, 2.1

\*Tác giả liên hệ: Email: doannngocsan@gmail.com

## 1. Đặt vấn đề

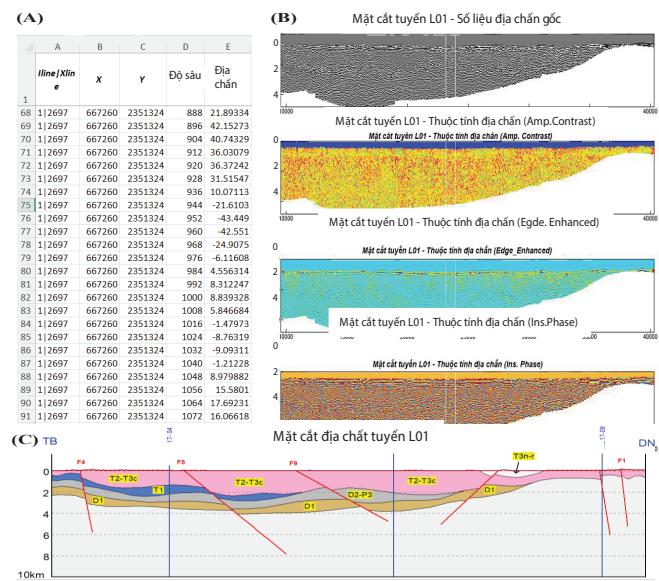
Đòi hỏi đổi mới công nghệ tìm kiếm thăm dò dầu khí (khảo sát và xử lý minh giải) để tìm ra các mỏ dầu khí mới là rất cấp bách trong điều kiện trữ lượng dầu khí ngày càng cạn kiệt, môi trường địa chất rất phức tạp của các vùng tiềm năng còn lại thì việc đổi mới công nghệ tìm kiếm thăm dò dầu khí (khảo sát và xử lý minh giải) để tìm ra các mỏ mới là rất cấp bách. Công tác xử lý số liệu yêu cầu đẩy mạnh việc phát triển về thuật toán, đòi hỏi việc sử dụng và tối ưu hóa các hệ thống phần cứng và đáp ứng việc xử lý Big data [1]. Với việc các bài toán cơ bản như xử lý hình ảnh và xử lý ngôn ngữ tự nhiên ngày càng tiếp cận với trình độ cơ bản của con người, trí tuệ nhân tạo (artificial intelligence - AI) sẽ là một giải pháp mang tính đột phá để khai phá lượng dữ liệu địa chất - địa vật lý để tìm ra các mối quan hệ ẩn giữa các lớp thông tin địa chất - địa vật lý và dầu khí, qua đó tìm ra các bộ thuộc tính đặc trưng để đánh giá triển vọng dầu khí [2]. Khu vực trung An Châu hiện đang được đánh giá là một trong những vị trí có tiềm năng dầu khí cao trên đất liền ở Việt Nam [3].

## 2. Nội dung nghiên cứu

### 2.1. Tích hợp cơ sở dữ liệu phục vụ đầu vào cho hệ thống học máy

Sự phát triển mạnh mẽ của công nghệ thông tin, sự nâng cấp về thiết bị và các chương trình ứng dụng cho phép chiết xuất được khối lượng thông tin ngày càng nhiều từ tài liệu địa chấn, làm tăng số lượng tham số có thể sử dụng được trong quá trình minh giải. Thuộc tính địa chấn (Seismic Attributes) được hiểu là những đặc trưng động lực học như: pha, tần số, biên độ hay các thông số phụ trợ khác như thuộc tính đa mạch (coherency...). Các thuộc tính địa chấn có thể được tính theo mặt cắt, theo bề mặt hoặc theo khối [4]. Các thuộc tính địa chấn bao gồm cả các đặc điểm hình học động (thời gian, tốc độ...) và đặc điểm động lực (pha, biên độ, tần số, độ suy giảm năng lượng). Thuộc tính địa chấn còn có thể được xác định theo đơn mạch (được tính cho từng mạch địa chấn), liên kết giữa các mạch hoặc thông qua tần số, biên độ, pha, tần số tức thời, pha tức thời, cường độ phản xạ... Các thuộc tính đa mạch được tính trên cơ sở hàm tương quan liên kết theo một nhóm mạch địa chấn và trên nhiều xung địa chấn (coherency).

Trong các thuộc tính địa chấn, thuộc tính được quan tâm nhiều hiện nay là biên độ. Từ các dị thường biên độ như các “điểm sáng”, “điểm tối”... có thể phát hiện các dấu hiệu liên quan đến dầu khí. Dị thường biên độ thường được sử dụng để thành lập các bản đồ phản ánh sự biến đổi tương đất đá và tham số vỉa của tầng chứa. Sự biến đổi biên độ theo chiều ngang có thể là cơ sở phân biệt sự khác nhau của các loại tương đất đá [5].



**Hình 1. Tích hợp số liệu địa chất - địa chấn. (A)** Số liệu địa chấn, **(B)** Các thuộc tính địa chấn, **(C)** Mặt cắt địa chất.

Hình 1 là các thông tin địa chấn và kết quả minh giải theo phương pháp truyền thông tại tuyến L01. Theo đó, đã minh giải được 5 tầng trầm tích (biểu hiện theo các màu) theo các mặt phản xạ địa chấn, có các tuổi khác nhau từ dưới lên như sau (hình 1C): Địa tầng tuổi Devon sớm (D<sub>1</sub>); Địa tầng tuổi Devon giữa - Pecmi muộn (D<sub>2</sub>-P<sub>3</sub>); Địa tầng tuổi Trias sớm (T<sub>1</sub>); Trias giữa - muộn (T<sub>2</sub>-T<sub>3c</sub>) và Trias muộn (T<sub>3n-r</sub>).

Hiện nay, có trên 50 thuộc tính địa chấn khác nhau có thể chiết xuất được từ tài liệu địa chấn và được sử dụng để minh giải các cấu trúc địa chất, địa tầng, tính chất chất lưu, cơ lý đá [6]. Có thể nói, sự phát triển của phân tích thuộc tính địa chấn gắn liền với phát triển trong lĩnh vực máy tính - điện toán. Phân tích thuộc tính địa chấn trong tìm kiếm dầu khí là quá trình sử dụng các kỹ thuật phân tích dữ liệu địa chất và địa chấn để đánh giá tiềm năng của các khu vực địa chất có khả năng chứa dầu khí. Trong quá trình phân tích thuộc tính địa chấn, các chuyên gia sử dụng các phương pháp phân tích dữ liệu như phân tích tương quan, phân tích phân cụm, phân tích thành phần chính và các phương pháp khác để phân tích và tìm kiếm các mô hình địa chấn. Các mô hình này cung cấp thông tin về các tính chất địa chất và địa chấn của khu vực đó, bao gồm các đặc điểm về cấu trúc, độ sâu, hướng và hình dạng của các lớp đất và đá.

Phân tích thuộc tính địa chấn rất quan trọng trong tìm kiếm dầu khí. Nó giúp xác định các khu vực có tiềm năng khai thác dầu khí và giảm thiểu các rủi ro và chi phí không cần thiết trong quá trình tìm kiếm. Nó cũng giúp các chuyên gia địa chất và địa chấn tìm hiểu sâu hơn về đặc điểm địa chất và địa chấn của các khu vực khai thác dầu khí, để có thể phát triển các phương pháp khai thác hiệu quả và bền vững hơn. Khả năng quản trị dữ liệu lớn cũng là một ưu điểm nổi trội của hệ thống AI. Độ chính xác “nhận dạng” đối tượng của ML càng tăng lên khi chúng ta sử dụng càng nhiều “features” đặc trưng cho

đối tượng cần nghiên cứu [7]. Chính vì vậy, nhằm mục đích nâng cao độ chính xác của thuật toán nhận dạng bằng ML, bằng phần mềm PETREL đã tiến hành phân tích các thuộc tính địa chấn như: <amp\_contrast>, <cur\_contour>, <Cur\_max>, <Curvature\_Az\_min>, 'Dominant\_Freq', 'Edge\_Enhanced', 'GLCM\_X', 'ins\_fre', 'ins\_phase', 'Polarity', 'RAI', 'trace\_gradient'... (hình 1).

Tuy nhiên, việc sử dụng số lượng feature một cách hợp lý, loại bỏ những features quá "liên kết chặt" với nhau sẽ làm tăng tính khách quan của tập hợp số liệu đầu vào. Để loại bỏ các thuộc tính có hệ số tương quan cao, tác giả đã sử dụng thuật toán 'Pearson' (linear correlation). Trong hình 2, các thuộc tính có  $R > |0,7|$  đã được loại bỏ khỏi đầu vào của hệ thống ML.

(A)	org	amp_contrast	chaos	cur_contour	Cur_max	Cur_max_std	Cur_max_diff	Cur_max_iderv	Curvature_Az_min	Dominant_Freq	Edge_Enhanced	envelope	GLCM_X	gradient_mag	ins_fre	ins_phase	Polarity	RAI	rms	sweetness	trace_gradient		
org	1.00	-0.01	0.02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
amp_contrast	-0.01	1.00	-0.66	-0.71	-0.09	-0.11	0.02	0.00	-0.02	0.01	-0.23	0.43	0.01	-0.11	0.41	-0.12	-0.43	0.01	-0.05	-0.01	-0.11	-0.11	-0.01
chaos	-0.02	-0.66	1.00	0.80	-0.05	0.06	0.09	0.00	-0.01	0.11	0.18	-0.45	-0.10	-0.11	0.41	-0.12	-0.43	0.01	-0.05	-0.01	-0.11	-0.11	-0.01
cur_contour	0.00	-0.71	0.80	1.00	0.01	0.00	-0.10	0.08	0.10	-0.08	0.10	-0.51	-0.08	-0.08	0.38	-0.08	-0.48	0.00	0.00	0.00	-0.09	-0.08	0.00
Cur_max	0.00	-0.09	0.05	0.01	1.00	0.94	0.44	0.23	0.23	-0.19	0.45	0.29	-0.14	-0.06	-0.06	0.06	-0.07	-0.13	-0.04	-0.02	0.00	-0.06	0.00
Cur_max_std	0.00	-0.11	0.06	0.00	0.94	1.00	0.47	0.21	-0.19	0.43	0.33	-0.15	-0.06	-0.07	0.07	-0.08	-0.14	-0.04	-0.03	0.00	-0.07	-0.07	0.00
Cur_max_diff	0.00	0.02	0.09	-0.10	0.44	0.47	1.00	-0.01	-0.18	0.86	0.11	-0.12	-0.05	-0.06	0.07	-0.06	-0.11	0.06	0.09	0.00	-0.06	-0.06	0.00
Curvature_Az_min	0.00	0.00	0.08	0.23	0.01	-0.01	1.00	0.29	0.07	-0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00	-0.01	0.00	0.00	0.00	0.00	0.00
Dominant_Freq	0.00	-0.03	-0.01	0.10	-0.19	-0.19	-0.18	0.29	1.00	-0.06	-0.03	0.01	0.01	0.01	-0.01	-0.01	0.00	0.02	0.01	-0.01	-0.01	-0.01	0.00
Edge_Enhanced	0.00	0.01	0.11	0.08	0.45	0.43	0.86	0.07	-0.06	1.00	0.10	-0.16	-0.07	-0.07	0.08	-0.07	-0.15	-0.05	0.11	0.00	-0.07	-0.07	0.00
envelope	0.00	0.23	0.18	0.10	0.29	0.33	0.11	0.03	-0.03	0.10	1.00	-0.05	0.00	0.01	0.09	0.00	0.02	0.00	0.06	-0.01	0.01	0.01	0.00
GLCM_X	0.00	0.43	-0.45	-0.51	-0.14	-0.15	-0.12	0.00	-0.01	-0.16	-0.05	1.00	0.18	0.20	-0.17	0.22	0.95	0.03	0.04	-0.02	-0.21	0.20	0.01
gradient_mag	0.00	0.01	-0.10	-0.08	-0.06	-0.06	-0.05	0.00	-0.01	-0.07	0.00	0.18	1.00	0.91	-0.09	0.81	0.20	0.03	0.05	0.14	0.90	0.91	-0.15
ins_fre	0.04	0.00	-0.11	-0.08	-0.06	-0.07	-0.06	0.00	-0.01	-0.07	0.01	0.20	0.91	1.00	-0.01	0.93	0.22	0.00	0.03	-0.06	1.00	1.00	0.03
ins_phase	-0.13	-0.37	0.41	0.38	0.06	0.07	0.07	0.00	-0.01	0.08	0.09	-0.17	-0.09	-0.01	1.00	-0.01	-0.15	-0.51	-0.03	-0.27	-0.02	-0.01	0.17
Polarity	0.01	0.02	-0.12	-0.08	-0.07	-0.08	-0.06	0.00	0.00	-0.01	-0.07	0.00	0.22	0.81	0.93	-0.01	1.00	0.23	0.01	0.02	-0.10	0.95	0.93
RAI	0.00	0.39	0.43	0.48	-0.13	-0.14	-0.11	0.00	0.00	-0.15	-0.02	0.95	0.20	0.21	-0.15	-0.23	1.00	0.04	0.04	-0.07	0.23	0.23	0.01
rms	-0.01	0.05	-0.01	0.00	-0.04	0.04	-0.04	0.00	0.02	0.05	0.00	0.03	0.03	0.00	-0.51	0.01	0.04	1.00	-0.01	0.21	0.01	0.00	0.18
sweetness	0.06	0.05	0.05	0.00	-0.02	-0.03	0.09	-0.01	0.01	0.11	0.06	0.04	0.05	0.03	-0.03	0.02	0.04	-0.01	1.00	-0.02	0.03	0.03	0.03
trace_gradient	0.20	0.00	-0.01	0.00	0.00	0.00	0.00	-0.01	0.00	-0.01	-0.02	0.14	-0.06	-0.22	-0.10	-0.02	0.21	-0.02	1.00	-0.03	-0.06	-0.73	

(B)	org	amp_contrast	chaos	Cur_max	Curvature	Dominant	GLCM_X	ins_phase	Polarity	sweetness	trace_grad
org	1										
amp_contrast		1									
chaos			1								
Cur_max				1							
Curvature_Az_min					1						
Dominant_Freq						1					
GLCM_X							1				
ins_phase								1			
Polarity									1		
sweetness										1	
trace_gradient											1

Hình 2. Ma trận PC tương quan các thuộc tính địa chấn. (A) gốc và (B) đã loại bỏ các thuộc tính có R cao.

A	B	C	D	E	K	Q	R	X	AD	AI	AP	AV	AW	BC	BD	BI	CB	CC	CL
line	X	Y	Độ sâu	Địa chấn	amp_contrast	chaos	cur_contour	Cur_max	Curvature_Az_min	Dominant_Freq	Edge_Enhanced	envelope	GLCM_X	gradient_mag	ins_fre	ins_phase	rms	sweetness	trace_grad
12897	66720	235124	538	0	0.594957	0.254632	-0.00209	-0.00169	179.6415	0.007704	0.02411	0.013539	1	0.005065	0.007697	-90	0.012725	0.013539	72.79
12897	66720	235124	536	0.011713	0.594531	0.156196	-0.00317	-0.00225	179.6236	0.02202	0.02802	0.028038	1	0.013854	0.012015	-85.3262	0.017216	0.028038	72.79
12897	66720	235124	544	0.029282	0.566624	0.101117	-0.00454	-0.00356	179.6031	0.009596	0.032784	0.03324	1	0.010612	0.009955	-28.5633	0.020256	0.03324	72.79
12897	66720	235124	552	0.021474	0.552518	0.066664	-0.0056	-0.00492	179.4501	0.008031	0.025922	0.022315	1	0.01824	0.06604	-15.7877	0.025602	0.022315	72.79
12897	66720	235124	560	0.025378	0.554084	0.030467	-0.0082	-0.00771	179.4748	0.001187	0.029481	0.029095	1	0.012172	0.010111	-18.8002	0.029677	0.029095	72.79
12897	66720	235124	568	0.025378	0.567525	0.017783	-0.00496	-0.00352	179.4121	0.003272	0.025973	0.029342	1	0.022239	0.003266	-26.4369	0.029362	0.029342	72.79
12897	66720	235124	576	0.031817	0.585511	0.030177	-0.01283	-0.01069	182.1763	0.005576	0.026945	0.024601	1	0.016807	0.005022	-16.4401	0.028833	0.024601	72.79
12897	66720	235124	584	0.021474	0.594069	0.03501	-0.000167	-0.00252	180.0059	0.0058831	0.022416	1	0.015439	-0.00591	-23.5005	0.026735	0.022416	72.79	
12897	66720	235124	592	0.017561	0.579659	0.068888	-0.01447	-0.00938	180.0052	0.028866	0.025803	1	0.015336	-0.0052	-47.0861	0.025244	0.025803	72.79	
12897	66720	235124	600	0.019521	0.542862	0.288312	-0.04124	-0.02104	180.0027	0.114743	0.03902	1	0.019344	0.020267	-99.9894	0.047782	0.03902	72.79	
12897	66720	235124	608	0.040995	0.494713	0.840548	-0.1164	-0.074601	180.00179	0.11264	0.07503	1	0.027076	0.007338	-66.8773	0.068821	0.07503	72.79	
12897	66720	235124	616	0.10913	0.551711	0.25401	-0.06776	-0.02964	180.001739	0.125139	0.14899	1	0.019776	0.009761	-35.9014	0.017703	0.14899	72.79	
12897	66720	235124	624	0.228401	0.43908	0.12036	0.044716	0.019915	0.002391	0.141831	0.220288	0.056693	1	0.049945	0.011812	-6.0079	0.207994	0.220288	72.79
12897	66720	235124	632	0.304535	0.480075	0.112245	-0.0875	-0.00329	0.285315	0.037369	0.177672	0.357926	1	0.115881	0.013258	-31.6975	0.218936	0.357926	72.79
12897	66720	235124	640	0.173741	0.434094	0.138959	-0.0225	-0.0048	0.38874	0.014495	0.143804	0.860584	1	0.013746	0.013746	-70.2047	0.348063	0.143804	72.79
12897	66720	235124	648	0.22059	0.426784	0.161671	-0.01667	-0.0088	0.490614	0.014641	0.05913	0.671997	1	0.031244	0.273893	0.013764	199.1653	0.161671	0.671997
12897	66720	235124	656	-0.88918	0.423711	0.131713	-0.0074	-0.00204	179.4571	1.42142	0.063958	0.9564	1	0.134232	0.013915	-148.9973	0.384981	0.063958	72.79
12897	66720	235124	664	-0.8482	0.425314	0.254574	-0.00165	-0.00095	179.4659	1.45412	0.06238	0.87751	1	0.012123	0.452978	0.014444	-170.237	0.384981	0.87751
12897	66720	235124	672	-0.51361	0.438838	0.22006	-0.06405	-0.03701	1.61142	0.286163	0.844889	0.887704	1	0.47559	0.015945	-127.423	0.614227	0.844889	
12897	66720	235124	680	-0.162038	0.435781	0.424444	-0.02951	-0.00909	0.76977	1.98142	0.493876	0.709715	1	0.843988	0.341219	-0.019296	-76.802	0.51683	0.709715
12897	66720	235124	688	0.002214	0.465737	0.4317	-0.00162	-0.0009	0.804621	1.156142	0.663879	0.819819	1	0.540887	0.021255	-41.0021	0.41279	0.819819	
12897	66720	235124	696	0.134531	0.440461	0.378074	-0.02457	-0.013963	1.88142	0.662674	0.6343								

Cây quyết định được xây dựng bằng cách chia tập dữ liệu gốc, bắt đầu từ nút gốc của cây, sau đó tạo ra các tập con kế tiếp. Việc chia tách này dựa trên các quy tắc phân chia được thiết lập dựa trên các đặc điểm phân loại. Quá trình này được lặp đi lặp lại một cách đệ quy, được gọi là phân vùng đệ quy. Quá trình đệ quy kết thúc khi tập hợp con tại một nút có tất cả các giá trị giống nhau của biến mục tiêu hoặc khi việc phân tách không mang lại cải thiện cho dự đoán. Phương pháp suy diễn từ trên xuống của DT là một ví dụ về thuật toán tham lam và cho đến nay vẫn là chiến lược phổ biến nhất để xây dựng cây quyết định từ dữ liệu.

Ưu điểm của DT là:

- Dễ hiểu và diễn giải: DT có cấu trúc rõ ràng, dễ hiểu và diễn giải, do đó dễ dàng giải thích cho những người không có kinh nghiệm về ML.

- Khả năng xử lý các dữ liệu lớn: DT có khả năng xử lý các tập dữ liệu lớn, có thể xử lý các tập dữ liệu chứa hàng triệu bản ghi.

- Tính linh hoạt: DT có tính linh hoạt cao trong việc chọn các tính năng và tham số để tối ưu hóa kết quả.

- Hiệu quả tính toán: DT có thể được xây dựng và tối ưu hóa bằng các thuật toán hiệu quả, do đó, thời gian tính toán và tài nguyên cần thiết để xây dựng một DT là tương đối thấp.

- Có thể xử lý dữ liệu không đầy đủ: DT có khả năng xử lý các tập dữ liệu không đầy đủ hoặc bị thiếu thông tin.

Nhược điểm của DT:

- Dễ bị quá khớp: DT có thể trở nên quá phức tạp và dễ bị quá khớp với dữ liệu huấn luyện nếu không được tối ưu hóa một cách chính xác.

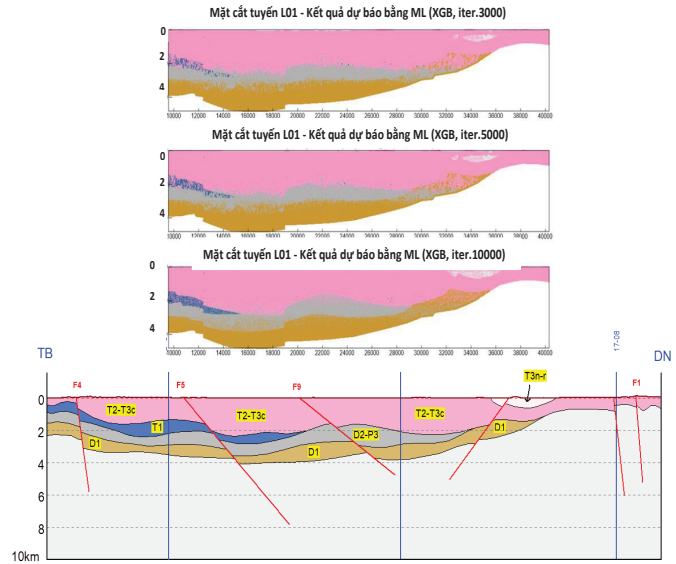
- Nhạy cảm với nhiễu và các giá trị ngoại lai: DT có thể dễ bị ảnh hưởng bởi các giá trị ngoại lai hoặc nhiễu trong dữ liệu.

- Không hiệu quả với các dữ liệu dạng liên tục: DT không hiệu quả trong việc xử lý các tập dữ liệu dạng liên tục.

- Khó khăn trong việc xử lý các mối quan hệ phức tạp giữa các tính năng: DT có thể không hiệu quả trong việc xử lý các mối quan hệ phức tạp giữa các trường thông tin.

Thuật toán XGBoost [8] là viết tắt của eXtreme Gradient Boosting. “eXtreme” nói đến các cải tiến về tốc độ như tính toán song song và nhận biết bộ nhớ cache giúp XGBoost nhanh hơn khoảng 10 lần so với tăng cường Gradient Boosting thông thường. Ngoài ra, XGBoost còn bao gồm một thuật toán tìm kiếm phân tách độc đáo để tối ưu hóa cây, cùng với tính năng chính quy tích hợp giúp giảm việc trang bị quá mức. Nói chung, XGBoost là phiên bản Gradient Boost nhanh hơn, chính xác hơn. Vì XGBoost là phiên bản nâng cao của Gradient Boosting và kết quả của nó là không song song nên nó có thể được cho là ML tốt nhất chúng ta có hiện nay.

### 3. Kết quả và bàn luận



Hình 5. Kết quả nhận dạng tương đất đá theo DT-XGBoost.

Nghiên cứu đã sử dụng các lát cắt địa chất truyền thống đã xây dựng để làm “mẫu luyện mạng” (hình 5). Số lượng mẫu học chỉ chiếm 60% và hoàn toàn ngẫu nhiên về vị trí XY. Cấu trúc mạng DT-XGBoost đã được sử dụng thử nghiệm với các số iteration lần lượt là 3.000, 5.000 và 10.000. Mục đích sử dụng DT-XGBoost là xác định mối liên hệ giữa mặt cắt địa chất đã biết (mẫu học) và dữ liệu địa chấn (số liệu gốc và các thuộc tính địa chấn). Kết quả nhận dạng được trình bày trong hình 6.



Hình 6. Sai số nhận dạng tương đất đá với số lần lặp. (A) >3.000, (B) >5.000, (C) >10.000.

Các tham số thống kê đánh giá sai số được trình bày trong hình 6 cho thấy, sai số nhận dạng phân loại với 3000 lần lặp còn thấp: nhận dạng sai 16%, đặc biệt đất đá  $T_3n-r$  và T1 có sai số cao (75-80%). Khi phân loại với số lần lặp trên 10.000 lần thì độ chính xác đã được cải thiện đáng kể (hình 6C).

Các tham số thống kê của số liệu địa chấn và 11 thuộc tính địa chấn cho thấy nếu chỉ sử dụng các tham số thống kê với các thuật toán thống kê truyền thống để phân loại đất đá thì rất khó khăn, vì các nhóm tham số có giá trị rất gần nhau và cấu trúc số liệu không hoàn toàn phản ánh cấu trúc địa chất của các tuyến (hình 1). Các thuật toán ML đã chứng tỏ được khả năng vượt trội của mình so với mình giải theo các phương pháp truyền thống trong công tác nhận dạng/phân loại đối tượng theo các thuộc tính. Đó là khả năng quản trị dữ liệu lớn: Bảng số liệu tích hợp các tuyến địa chấn có dung lượng khoảng xấp xỉ 10TB (hàng triệu dòng và gần 80 cột); khả năng khai thác các mối quan hệ ẩn sâu giữa các cột thuộc tính từ đó nhận dạng các tương đá với độ chính xác cao.

Các tham số thống kê mức độ chính xác nhận dạng phân loại đất đá bằng ML và truyền thống cho thấy (hình 7 và 8) [9]:

Hệ tầng Nà Khuất, Mẫu Sơn ( $T_2-T_{3c}$ ) gồm sét vôi, bột kết, cát kết, đá phiến sét, đá vôi; sạn kết, cát kết, sét vôi, cát kết màu đỏ, bột kết xen cát kết màu đỏ; Phân hệ tầng dưới: cát kết, cát kết dạng quazit, thấu kính cuội kết, bột kết màu đỏ được nhận dạng chính xác (~90% đúng).

Hệ tầng Tam Hoa, Bắc Sơn, Đồng Đăng ( $D_2-P_3$ ) với thành phần đất đá gồm cuội kết, sạn kết, cát kết arko, đá phiến sét, đá vôi, đá vôi dolomit, bauxit, bột kết, đá vôi, sét silic được nhận dạng với độ chính xác trung bình cao (70-90% đúng).

Ngược lại, hệ tầng Văn Lãng ( $T_3n-r$ ) gồm cát kết, bột kết, sét - bột kết màu đen, xám đen, đá vôi - silic, sét than lại có độ chính xác nhận dạng trung bình thấp nhất (~45-50%).

Tuyến L09						Tuyến L07					
Total	Corr. forecast	Aver. ident				Total	Corr. forecast	Aver. ident			
1484940	67.65	57.29				1097313	84.52	76.93			
Formation	Class_No	Sample_count	Percent	Corr. forecast	Aver. ident	Formation	Class_No	Sample_count	Percent	Corr. forecast	Aver. ident
D1	1	260466	17.54	66.61	60.46	D1	1	154252	14.06	77.35	68.05
D2-P3	2	402326	27.09	76.3	55.57	D2-P3	2	221537	20.19	77.81	62.62
J-K	4	12266	0.83	45.68	63.41	J-Kr	5	6894	0.63	56.7	79.34
Rhyolit T2	8	104279	7.02	23.95	50.74	T1	9	98852	9.01	60.83	60.24
T1	9	248162	16.71	61.07	51.19	T2-T3c	11	581010	52.95	95.11	87.45
T2-T3c	11	359093	24.18	82.72	63.46	T3n-r	12	34768	3.17	54.97	78.56
T3n-r	12	98348	6.62	45.62	54.97						

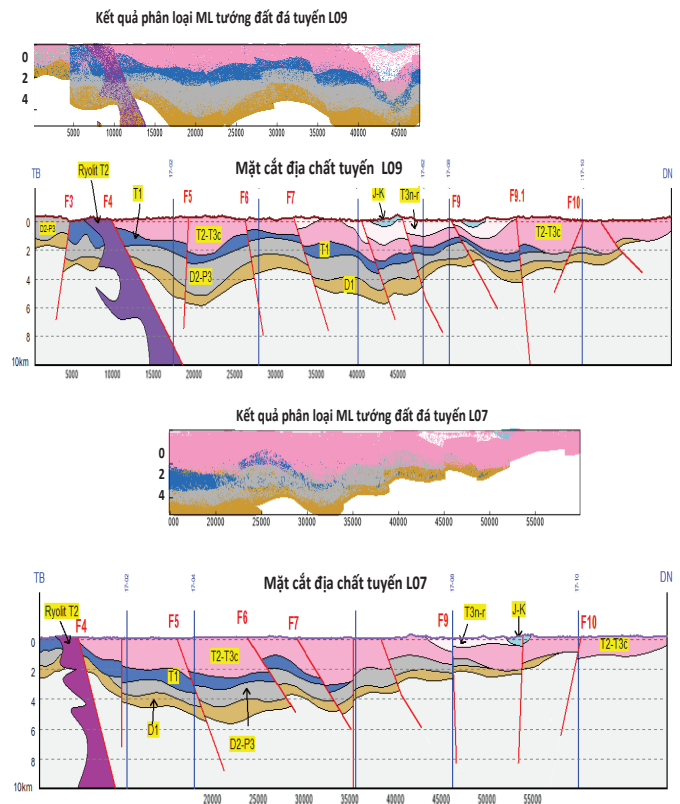
  

class No	No sample	Incorrect	Percent
-1	553631	50728	9.16 %
11	307904	27815	9.04 %
2	299948	33989	11.33 %
1	177529	15931	8.98 %
9	93678	50126	53.51 %
12	33208	20266	61.01 %
8	14170	5	0.04 %
4	4872	3835	78.73 %
Total:	1484940	480452	32.36 %

class No	No sample	Incorrect	Percent
11	589010	178156	30.23 %
2	178156	139989	78.58 %
-1	139989	115931	82.83 %
1	115931	9	0.01 %
9	50126	12	0.02 %
12	20266	5	0.02 %
5	3835	169866	44.30 %
Total:	1097313	169866	15.48 %

Hình 8. Sai số nhận dạng tương đất đá trên các mặt cắt L09 và L07.



Hình 9. Mặt cắt địa chất và kết quả nhận dạng tương đất đá trên các mặt cắt L09 và L07.

Theo kết quả xây dựng mặt cắt địa chấn - địa chấn trên cơ sở nhận dạng theo dữ liệu địa chấn (hình 9 và 10) có thể thấy: ML đã khẳng định có mối liên quan chặt chẽ giữa các hệ tầng trên mặt cắt địa chất do các chuyên gia truyền thống xây dựng và dữ liệu địa chấn, đặc biệt là dữ liệu thuộc tính địa chấn. Đa số các đoạn tuyến đều có tính phù hợp cao giữa mặt cắt địa chất đã biết và kết quả dự báo bằng ML.

Tuy nhiên, DT (XGB) cũng phát hiện một vài chỗ không phù hợp giữa kết quả của truyền thống và liên kết địa chấn [10]:

Tuyến L06.1						Tuyến L06.2						Tuyến L08					
Total	Corr. forecast	Aver. ident				Total	Corr. forecast	Aver. ident				Total	Corr. forecast	Aver. ident			
3811064	79.43	71.71				4688156	82.85	75.71				1423018	77.01	64.97			
Formation	Class_No	Sample_count	Percent	Corr. forecast	Aver. ident	Formation	Class_No	Sample_count	Percent	Corr. forecast	Aver. ident	Formation	Class_No	Sample_count	Percent	Corr. forecast	Aver. ident
D1	1	620314	16.28	84.47	63.74	D1	1	249639	5.31	57.7	60.32	D1r	0	340749	23.95	79.59	63.94
D2-P3	2	745819	19.56	69.27	60.11	J-K	4	1153303	23.79	92.55	81.05	D1	1	94954	6.67	64.27	61.55
T1	9	147762	3.88	50.89	54.72	T1	9	241578	5.15	70.53	70.39	D2-P3	2	133205	9.36	57.36	57.13
T2-T3c	11	194062	5.15	89.98	78.95	T2-P3	10	240178	5.12	62.38	62.95	J-K	4	300398	13.38	78.57	72.25
T3n-r	12	348097	9.13	45.19	69.86	T2-T3c	11	857358	18.29	81.67	74.72	T1	9	144266	10.14	68.73	61.6
						T3n-r	12	1084080	42.31	85.05	76.3	T2-T3c	11	519446	36.51	84.4	65.99

class No	No sample	Incorrect	Percent
11	199556	1	0.00 %
1	548040	4	0.00 %
-1	519944	13446	2.58 %
2	517937	2	0.00 %
12	173739	9	0.01 %
9	55448	1	0.00 %
Total:	3811064	794075	20.83 %

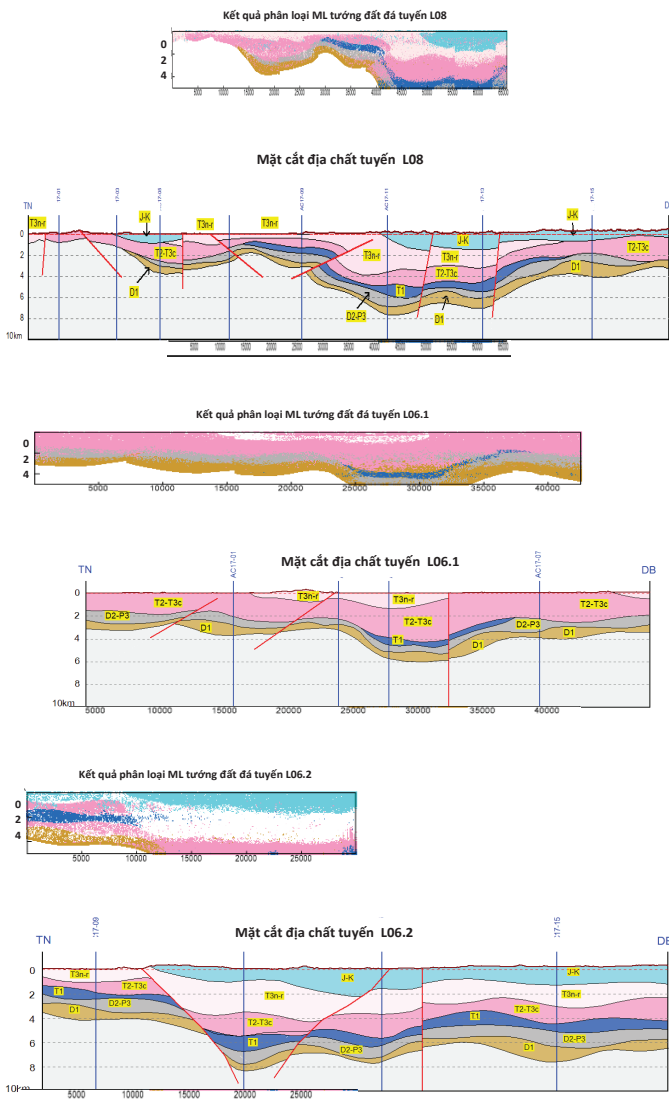
  

class No	No sample	Incorrect	Percent
12	1088057	4	0.00 %
4	1193541	11	0.00 %
11	754810	1	0.00 %
-1	431653	843	0.19 %
9	177899	10	0.01 %
10	143492	1	0.00 %
1	118704	1	0.00 %
Total:	4688156	803940	17.15 %

class No	No sample	Incorrect	Percent
11	431528	1	0.00 %
-1	141341	2310	1.63 %
0	277879	4	0.00 %
4	144711	9	0.01 %
9	95761	2	0.00 %
2	58909	1	0.00 %
1	52889	1	0.00 %
Total:	1423018	327159	23.00 %

Hình 7. Sai số nhận dạng tương đất đá trên các mặt cắt L06.1, L06.2 và L08.



Hình 10. Mặt cắt địa chất và kết quả nhận dạng tương đất đá trên các mặt cắt L08, L06.1 và L06.2.

Trên tuyến L09, đoạn 40.000-47.000 m, độ sâu 2.000-3.000 m. Khả năng hệ tầng Lạng Sơn ( $T_1$ , cát kết, bột kết, đá phiến sét, đá vôi sét, sét vôi) có thể nằm ngang, không tạo nên đứt gãy.

Trên tuyến L08, đoạn 40.000-45.000 m, hệ tầng Lạng Sơn không có biểu hiện rõ và vị trí đứt gãy được xác định ở độ sâu nông hơn.

Trên dữ liệu địa chấn, các khối rhyorit thể hiện khá mờ nhạt và không rõ ràng với độ nhận dạng không cao (29%).

#### 4. Kết luận

Kết quả cho thấy, các phần mềm được phát triển trong phạm vi nghiên cứu cũng như hệ thống máy tính công suất cao (High Performance Computer - HPC) đáp ứng được khả năng xử lý cơ sở dữ liệu khổng lồ bằng các thuật toán tiên tiến. Tuy nhiên, để đáp ứng được công tác xử lý một lượng

lớn số liệu địa chất - địa vật lý rất cần các thuật toán mới, phương pháp mới, hướng đi mới để khắc phục được các yếu điểm của hệ thống truyền thống hiện nay. Đó là khả năng tổng hợp cơ sở dữ liệu nhiều cả về số lượng và định dạng; xác định các mối liên kết giữa dầu khí và thông tin địa chất - địa vật lý.

Với thử nghiệm mô hình cấu trúc ML dạng DT và thuật toán Gradient Boosting (GB) để phân loại có kiểm định (có mẫu học) trên cơ sở tài liệu địa chất - địa chấn bước đầu đã thu được một số kết quả nhất định:

Thử nghiệm thành công tích hợp dữ liệu địa chất - địa chấn thành cơ sở dữ liệu lớn theo nguyên tắc GIS.

Phát hiện các mối quan hệ ẩn sâu giữa các lớp thông tin qua đó nhận dạng với độ chính xác chấp nhận được các đối tượng địa chất theo mẫu đã biết.

Với DT - GB số lần lặp phải trên 10.000 lần thì kết quả nhận dạng mới chấp nhận được. Thời gian tính toán còn quá lâu (~60-90 giờ chạy máy) để đạt được kết quả cao.

#### LỜI CẢM ƠN

Nghiên cứu này là kết quả đề tài “Nghiên cứu xây dựng hệ thống trí tuệ nhân tạo tích hợp dữ liệu địa chất dầu khí đánh giá triển vọng dầu khí” (mã số KC-4.0-01/19-25) thuộc Chương trình trọng điểm quốc gia “Hỗ trợ nghiên cứu, phát triển và ứng dụng công nghệ của công nghiệp 4.0”. Các tác giả xin trân trọng cảm ơn.

#### TÀI LIỆU THAM KHẢO

- [1] D.N. San (2003), *Conducting and Completing Technology for Synthesizing and Analyzing Data using Geological Geographic Information Systems and Geological Expert Systems, Applied for in the Cho Don Area*, Ministerial level research project (in Vietnamese).
- [2] D.N. San, N.D. Nuong (2019), *Research and Development of an Artificial Intelligence System Integrated with Petroleum Geological Data for Petroleum Prospect Evaluation*, Government level research project (in Vietnamese).
- [3] N.T. San (1975), “How should the petroleum prospects in the An Chau Basin be assessed”, *Geological Journal*, **122(11-12)**, 6pp (in Vietnamese).
- [4] Schlumberger (2015), *3D Seismic Facies Generated From Independent Seismic Attributes Using The Neural Network Technique*, 56pp.
- [5] M.T. Tan (2011), *Seismic Exploration*, Transport Publishing House, 45pp (in Vietnamese).
- [6] S. Chopra, K.J. Marfurt (2005), *Seismic Attribute Mapping of Structure and Stratigraphy*, EAGE, 30pp.
- [7] L. Zhang, C. Zhan (2017), “Machine learning in rock facies classification: An application of XGBoost”, *International Geophysical Conference*, Qingdao, China, Society of Exploration Geophysicists and Chinese Petroleum Society, DOI: 10.1190/IGC2017-351.
- [8] Boosting Algorithm (2020), *XGBoost*, Towards Data Science, 2017-05-14.
- [9] PVEP Song Hong (2019), *Phase I Report of The An Chau Basin Research Project* (in Vietnamese).
- [10] PVEP Song Hong (2013), *Petrographic Analysis Report of The An Chau Basin Region* (in Vietnamese).