

# Dự báo chất lượng không khí tại Việt Nam bằng phương pháp chuỗi thời gian

1<sup>st</sup> Huỳnh Lê Phong  
Khoa Hệ thống thông tin  
Trường Đại học Công nghệ Thông tin  
21520086@gm.uit.edu.vn

2<sup>nd</sup> Nguyễn Quốc Trọng  
Khoa Hệ thống thông tin  
Trường Đại học Công nghệ Thông tin  
21521556@gm.uit.edu.vn

3<sup>rd</sup> Nguyễn Triệu Vy  
Khoa Hệ thống thông tin  
Trường Đại học Công nghệ Thông tin  
21522812@gm.uit.edu.vn

4<sup>th</sup> Vũ Thị Phương Linh  
Khoa Hệ thống thông tin  
Trường Đại học Công nghệ Thông tin  
20521541@gm.uit.edu.vn

5<sup>th</sup> Đỗ Đình Đăng Khoa  
Khoa Hệ thống thông tin  
Trường Đại học Công nghệ Thông tin  
21522218@gm.uit.edu.vn

**Tóm tắt nội dung**—Nghiên cứu này tập trung vào phân tích chất lượng không khí ở ba thành phố lớn tại Việt Nam: Hà Nội, Hà Long và Việt Trì. Mục tiêu của nghiên cứu là sử dụng các kỹ thuật phân tích chuỗi thời gian tiên tiến để dự báo và hiểu các chỉ số chất lượng không khí như PM2.5, PM10, O3, NO2, SO2 và CO. Phương pháp phân tích bao gồm cả các thuật toán truyền thống và tiên tiến như Hồi quy Tuyến tính, Mô hình ARIMA, Mạng Nơ-ron Hồi quy (RNN), Đơn vị Hồi quy Cổng (GRU), LSTM, VAR, Rừng Ngẫu nhiên, Mô hình Mã hóa Dày Chuỗi Thời gian (TiDE), Autoformer, Mạng Nơ-ron Đa Tầng (MLP). Chúng tôi sử dụng hai tỷ lệ khác nhau giữa các tập dữ liệu huấn luyện: testing, cụ thể là 7:3, 8:2 và 9:1, để đánh giá hiệu suất của các mô hình dưới các điều kiện khác nhau. Sau đó, nó so sánh hiệu suất của các mô hình khác nhau dựa trên ba chỉ số: Mean Absolute Percentage Error (MAPE), Root Mean Square Error (RMSE), và Mean Logarithmic Square Error (MLSE). Cuối cùng, hai mô hình có hiệu suất tốt nhất được sử dụng để dự báo chất lượng không khí cho 30 ngày tiếp theo, thể hiện hiệu quả của chúng trong việc dự báo chất lượng không khí. Bằng cách sử dụng các tập dữ liệu thời gian thực, nghiên cứu này nhằm dự báo, khám phá các mẫu, xu hướng và mối quan hệ giữa các chỉ số chất lượng không khí. Ngoài ra, nghiên cứu cũng đóng góp vào việc hiểu rõ hơn về việc giải quyết các vấn đề ô nhiễm không khí, thúc đẩy môi trường đô thị khỏe mạnh và bền vững hơn.

**Index Terms**—Chất lượng không khí, ô nhiễm không khí, phân tích chuỗi thời gian, dự báo, Hồi quy Tuyến tính, ARIMA, RNN, GRU, LSTM, VAR, Random Forest, TiDE, Autoformer, MLP, Hà Nội, Hà Long, Việt Trì, Việt Nam.

## I. GIỚI THIỆU

Trong những thập kỷ gần đây, sự chú ý toàn cầu ngày càng tập trung vào suy thoái môi trường, đặc biệt là ô nhiễm không khí, do tác động sâu rộng của nó đối với sức khỏe con người, hệ sinh thái và sự ổn định của khí hậu. Hạt phấn mịn (PM2.5 và PM10) cùng với các khí như ozone (O3), dioxide nitơ (NO2), dioxide lưu huỳnh (SO2), và carbon monoxide (CO) đóng vai trò quan trọng trong sự suy giảm chất lượng không khí. Trong

bối cảnh này, nghiên cứu của chúng tôi nhằm khám phá động lực chất lượng không khí trong môi trường đô thị, đặc biệt là ở ba thành phố lớn của Việt Nam: Hà Nội, Hà Long và Việt Trì. Mỗi thành phố đại diện cho một ngữ cảnh kinh tế-xã hội và địa lý khác biệt, mang lại cái nhìn sâu sắc về các thách thức và cơ hội của quản lý ô nhiễm không khí. Mục tiêu tổng thể là sử dụng các kỹ thuật phân tích chuỗi thời gian tiên tiến để dự báo và phân tích các tham số về chất lượng không khí. Chúng tôi sử dụng một loạt các thuật toán, từ các phương pháp truyền thống như Hồi quy Tuyến tính và Mô hình Trung bình Chuyển động Tích hợp Tự động (ARIMA) đến các kiến trúc học sâu tiên tiến như Mạng Nơ-ron Hồi quy (RNN), Đơn vị Hồi quy Cổng (GRU), Bộ nhớ Ngắn hạn Dài - Ngắn hạn (LSTM), Hồi quy Vector (VAR), Rừng Ngẫu nhiên, Mô hình Mã hóa Dày Chuỗi Thời gian (TiDE), Autoformer, và Mạng Nơ-ron Đa Tầng (MLP). Bằng cách tận dụng những kỹ thuật này, chúng tôi nhằm mục đích khám phá các mẫu, xu hướng và mối liên hệ cơ bản trong dữ liệu chất lượng không khí. Các tập dữ liệu của chúng tôi bao gồm các số liệu đo thời gian thực về các tham số chất lượng không khí chính, bao gồm PM2.5, PM10, O3, NO2, SO2 và CO, được thu thập trong một khoảng thời gian dài. Những tập dữ liệu này cho phép khám phá các biến thể thời gian, sự đa dạng không gian và tương tác phức tạp giữa các chất gây ô nhiễm. Qua phân tích toàn diện và mô hình hóa dự đoán, nghiên cứu của chúng tôi nhằm cung cấp cái nhìn hành động cho các nhà quyết định chính sách, quy hoạch đô thị và các cơ quan môi trường được giao nhiệm vụ lập kế hoạch các chiến lược quản lý chất lượng không khí hiệu quả. Hơn nữa, các kết quả của chúng tôi đóng góp vào sự hiểu biết rộng lớn hơn về việc giải quyết các thách thức đa mặt do ô nhiễm không khí gây ra, thúc đẩy môi trường đô thị khỏe mạnh và bền vững hơn. Ngoài ra, điều quan trọng là nhấn mạnh về vấn đề ô nhiễm không khí cấp bách mà Việt Nam đang phải đối mặt, đặc biệt là ở các thành phố như Hà Nội, nơi chất lượng không khí đã nằm trong số tồi tệ nhất ở Đông Nam Á. Các yếu tố như mức độ cao của PM2.5 và các

Identify applicable funding agency here. If none, delete this.

chất gây ô nhiễm khác đã gây ra những lo ngại sức khỏe đáng kể và tác động tiêu cực đến GDP của đất nước, nhấn mạnh tính cấp bách của các chiến lược quản lý ô nhiễm không khí hiệu quả. Những nỗ lực của chính phủ để giải quyết những thách thức này, bao gồm các sáng kiến để thiết lập các tiêu chuẩn và quy định, cũng như tổ chức các hệ thống để đối phó với ô nhiễm không khí, là các yếu tố quan trọng trong bối cảnh rộng lớn của nghiên cứu của chúng tôi.

## II. NGHIÊN CỨU LIÊN QUAN

Trong những năm gần đây, đã có một sự gia tăng đáng kể trong nghiên cứu nhằm dự báo chất lượng không khí bằng cách sử dụng một loạt các kỹ thuật máy học, học sâu và thống kê.

Evgeniy Marinov, Dessislava Petrova-Antonova và Simeon Malinov, trong một nghiên cứu [1], tập trung vào việc cải thiện độ chính xác của việc dự báo chất lượng không khí thông qua việc áp dụng phương pháp ARIMA. Nghiên cứu của họ, dựa trên dữ liệu từ các trạm giám sát chất lượng không khí tại Thành phố Sofia, Bulgaria, từ ngày 1 tháng 1 năm 2015 đến ngày 31 tháng 12 năm 2019, đã cho thấy hiệu quả của các mô hình ARIMA trong việc dự báo nồng độ ô nhiễm như CO, NO<sub>2</sub>, O<sub>3</sub> và PM<sub>2.5</sub>.

S. H. Khaerun Nisa, Irfan Irfani và Utriweni Mukhaiyar, trong nghiên cứu của họ [2], đã khám phá việc dự báo mức độ ô nhiễm không khí tại Jakarta bằng phương pháp phân tích Vector Autoregressive (VAR). Bằng cách phân tích dữ liệu chuỗi thời gian về chỉ số chất lượng không khí (AQI) và nồng độ PM<sub>2.5</sub> thu thập từ ngày 16 tháng 8 đến ngày 25 tháng 9 năm 2023, để huấn luyện, và từ ngày 25 tháng 9 đến ngày 1 tháng 10 năm 2023, để kiểm tra, họ đã xác định mô hình VAR(2) tối ưu cho việc dự báo ô nhiễm không khí chính xác.

Khawaja Hassan Waseem và cộng sự đã khám phá tác động của các yếu tố khí hậu đối với nồng độ PM<sub>2.5</sub> và phát triển mô hình dự báo trong một nghiên cứu gần đây [3]. Nghiên cứu của họ, kéo dài trong 30 ngày và 72 giờ cho các dự đoán hàng ngày và hàng giờ tương ứng, đã kết hợp dữ liệu chất lượng không khí, chất gây ô nhiễm và điều kiện khí hậu từ nhiều thành phố ở Pakistan. Sử dụng các mô hình máy học và học sâu bao gồm FbProphet và LSTM, họ đã phát hiện mô hình mã hóa-giải mã LSTM vượt trội so với các mô hình khác, đạt được cải thiện đáng kể về độ chính xác dự báo.

Hai tác giả, Marwa Winis Misbah Esager và Kamil Demirberk Ünlü [4], đã áp dụng mô hình LSTM (Long Short-Term Memory) để dự báo nồng độ hàng giờ của hạt phần mịn PM<sub>2.5</sub> tại Tripoli, Libya. Họ đã sử dụng 100 epochs để huấn luyện mô hình, và kết quả tốt nhất đã được đạt được khi số nút được đặt là 20. Kết quả RMSE trên tập kiểm tra cho thấy mức độ lỗi thấp, khoảng 0.0146, chứng tỏ mô hình dự báo có độ chính xác khá cao.

Ngoài ra, các nghiên cứu gần đây cũng đã khám phá việc sử dụng các kỹ thuật mô hình hóa tiên tiến để cải thiện dự báo chất lượng không khí. Ví dụ, Abhimanyu Das và các đồng nghiệp đã giới thiệu TiDE (Time-series Dense Encoder) trong công việc của họ [5]. TiDE, một mô hình mới được tùy chỉnh cho dự báo chuỗi thời gian dài hạn, cho thấy khả năng hứa hẹn trong việc xử lý các biến đổi phi tuyến tính trong dữ liệu chuỗi thời gian.

Các đánh giá thực nghiệm của TiDE đã cho thấy sự vượt trội hoặc tương đương với các phương pháp hiện tại trên các chỉ số dự báo dài hạn phổ biến trong thế giới thực, đồng thời tự hào về tốc độ suy luận và huấn luyện nhanh hơn đáng kể.

Ngược lại, nghiên cứu trước đó của Ong và đồng nghiệp [6] đã khám phá việc sử dụng các phương pháp dựa trên RNN để dự báo mức độ chất lượng không khí, trình bày một kỹ thuật động để tiền huấn luyện mô hình tập trung vào dự báo chuỗi thời gian nhiều bước trước. Tương tự, Lim và đồng nghiệp [7] đã sử dụng RNN để dự báo các chất gây ô nhiễm không khí khác nhau nhưng không tìm thấy sự khác biệt hoặc cải thiện đáng kể so với các mô hình truyền thống.

Bốn tác giả, Haixu Wu, Jiehui Xu, Jianmin Wang và Ming-sheng Long, đã phát triển một mô hình gọi là Autoformer [8]. Mô hình này được thiết kế để giải quyết vấn đề dự báo chuỗi thời gian dài hạn, một thách thức đáng kể trong các ứng dụng thực tế như dự báo thời tiết cực đoan và lập kế hoạch tiêu thụ năng lượng dài hạn. Autoformer vượt trội so với các mô hình dựa trên Transformer trước đó bằng cách tích hợp các khối Phân rã và Mô hình Tương quan Tự động dựa trên tính chu kỳ của chuỗi thời gian, cho phép mô hình khám phá và tổng hợp thông tin ở mức con-chuỗi. Kết quả thực nghiệm trên sáu thử nghiệm khác nhau, bao gồm năm ứng dụng thực tế từ năng lượng đến dịch bệnh, cho thấy rằng Autoformer đạt được độ chính xác hàng đầu với cải thiện tương đối 38% so với các phương pháp hiện tại.

Những phương pháp và kết quả đa dạng này nhấn mạnh những nỗ lực tiếp tục để nâng cao các kỹ thuật dự báo chất lượng không khí, tận dụng cả các phương pháp thống kê truyền thống và các mô hình máy học tiên tiến.

## III. TÀI NGUYÊN

### A. Bộ dữ liệu

Bài viết sử dụng 3 bộ dữ liệu lấy từ dữ liệu chất lượng không khí trên trang web aqicn.org từ 01/03/2019 - 01/03/2024 bao gồm 3 thành phố Hà Nội, Hạ Long và Việt Trì. Bộ dữ liệu bao gồm các thuộc tính cụ thể sau:

Bảng I  
MÔ TẢ THUỘC TÍNH

Thuộc tính	Mô tả
Ngày	Thời gian (YYYY-MM-DD)
Pm25	Bụi mịn có đường kính từ 2.5 đến 10 micron
Pm10	Bụi mịn có đường kính nhỏ hơn 2.5 micron
O3	Ozon
NO2	Điôxit nitơ
SO2	Điôxit lưu huỳnh
CO	Carbon monoxit

Statistic	HaLong	HaNoi	VietTri
Count	1828	1828	1828
Mean	39.89	62.03	42.13
Std Dev	23.21	40.14	32.02
Min	5	2	1
Q1	22	31	19
Q2 (Median)	37	53	34
Q3	54	86	59
Max	163	217	178
Mode	5	24	1
Variance	538.861	1611.963	1025.602
Kurtosis	0.398	0.411	1.934
Skewness	0.705	0.943	1.263
CV	0.581	0.647	0.76

Bảng II  
STATISTICAL SUMMARY FOR HALONG, HANOI, AND VIETTRI

## B. Công cụ

Trong quá trình nghiên cứu và phân tích dữ liệu, chúng tôi đã tận dụng một loạt các công cụ phân tích thống kê trong Python để khám phá sâu hơn các mẫu dữ liệu và rút ra những kết luận có ý nghĩa. Các công cụ chính bao gồm: numpy, pandas, sklearn, matplotlib.pyplot,... Sử dụng các công cụ phân tích thống kê này đã giúp chúng tôi hiểu sâu hơn về dữ liệu và đưa ra những phát hiện quan trọng. Chi tiết các kết quả có thể được tìm thấy trong bảng mô tả và biểu đồ đi kèm.

## C. Tỷ lệ phân chia dữ liệu

Trong phân tích dữ liệu chuỗi thời gian của chúng tôi, chúng tôi đã chia tập dữ liệu thành hai phần: tập huấn luyện và tập kiểm tra, với các tỷ lệ khác nhau như 70% cho huấn luyện và 30% cho kiểm tra, 80% cho huấn luyện và 20% cho kiểm tra, và 90% cho huấn luyện và 10% cho kiểm tra. Những tỷ lệ này giúp chúng tôi đánh giá cách mà chúng ảnh hưởng đến hiệu suất của mô hình bằng cách xem xét sự phân phối của dữ liệu trong mỗi tập. Tỷ lệ phổ biến nhất là 7:3, chia 70% dữ liệu cho huấn luyện và 30% cho kiểm tra, tạo ra sự cân bằng giữa việc cung cấp đủ dữ liệu huấn luyện và đảm bảo các tập riêng biệt để điều chỉnh và đánh giá. Một lựa chọn khác là tỷ lệ 8:2, ưa chuộng việc sử dụng 80% dữ liệu cho huấn luyện, có ích cho các mô hình phức tạp yêu cầu tập dữ liệu huấn luyện lớn hơn. Trong một số trường hợp cụ thể, một cách tiếp cận thận trọng như tỷ lệ 9:1 có thể được ưa chuộng, đặc biệt khi xử lý tập dữ liệu lớn và một mô hình đơn giản. Tỷ lệ này đảm bảo có đủ dữ liệu huấn luyện trong khi vẫn giữ được một tập kiểm tra đáng kể để đánh giá hiệu suất.

## D. Đánh giá mô hình

Trong việc đánh giá hiệu suất của các mô hình, chúng tôi sử dụng ba chỉ số là Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), và Root Mean Squared Error (RMSE). Thuật toán có giá trị thấp nhất trong ba chỉ số này sẽ cho thấy mức độ chính xác tốt nhất. Dưới đây là các công thức để tính MAE, MAPE và RMSE.

Với các tham số:

- $n$ : Số lượng điểm dữ liệu.
- $y_i$ : Giá trị thực tế.
- $\hat{y}_i$ : Giá trị dự đoán.

**MAE (Mean Absolute Error):**

MAE là trung bình của các giá trị tuyệt đối của sai số giữa giá trị dự đoán và giá trị thực tế. Nó đo lường độ lớn của sai số trung bình và không quan tâm đến hướng của sai số. MAE càng nhỏ, mô hình càng chính xác.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (1)$$

**RMSE (Root Mean Square Error):**

RMSE là căn bậc hai của trung bình của các bình phương của sai số giữa giá trị dự đoán và giá trị thực tế. Nó đo lường độ lớn của sai số trung bình và cung cấp một con số tương đối về mức độ sai lệch giữa dự đoán và giá trị thực tế. RMSE càng nhỏ, mô hình càng chính xác.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (2)$$

**MAPE (Mean Absolute Percentage Error):**

MAPE là trung bình của tỉ lệ phần trăm của sai số tuyệt đối so với giá trị thực tế. Nó thường được sử dụng để đo lường tỷ lệ phần trăm trung bình của sai số so với giá trị thực tế, cung cấp cái nhìn về mức độ chính xác của mô hình dự đoán. MAPE càng nhỏ, mô hình càng chính xác.

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100 \quad (3)$$

## IV. PHƯƠNG PHÁP

### A. LINEAR REGRESSION

Hồi quy tuyến tính là một kỹ thuật quan trọng trong thống kê, cho phép chúng ta khám phá và mô hình hóa mối quan hệ giữa các biến. Giúp chúng ta hiểu rõ hơn về cách các yếu tố độc lập ảnh hưởng đến một biến phụ thuộc. Trong hồi quy tuyến tính đa biến, chúng ta có thể xem xét tác động của nhiều biến độc lập đến một biến phụ thuộc, giúp tạo ra các dự đoán hoặc giải thích phức tạp hơn về thực tế. Một mô hình hồi quy tuyến tính đa biến có dạng:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon$$

Trong đó:

- $Y$ : Biến phụ thuộc.;
- $X_1, X_2, \dots, X_k$ : Các biến độc lập.;
- $\beta_0$ : Hệ số chặn;
- $\beta_1, \beta_2, \dots, \beta_k$ : Các hệ số hồi quy.;
- $\epsilon$ : Thành phần sai số.

## B. ARIMA

Phương pháp ARIMA (Autoregressive Integrated Moving Average) sử dụng một mô hình AR (Autoregressive) kết hợp với một mô hình MA (Moving Average) để thực hiện dự báo chuỗi thời gian. Các tham số chính cần xem xét bao gồm:

- Số lượng quan sát trước đó ( $p$ );
- Độ chênh lệch ( $d$ );
- Kích thước của trung bình chuyển động ( $q$ ).

Mô hình AR hiển thị sự phụ thuộc của một quan sát vào một giai đoạn thời gian trước đó. Mô hình AR thu được  $p$  quan sát trước đó như sau:

$$y_t = \alpha + \sum_{i=1}^p \phi_i y_{t-i} + e_t \quad (5)$$

trong đó  $y_t$  là biến dự đoán cho thời điểm  $t$  từ phân phối chuẩn và  $y_{t-i}$  xác định  $p$  quan sát trước của cùng một chuỗi thời gian.  $\phi_i$  biểu thị các hệ số hồi quy,  $\alpha$  là một hằng số và  $e_t$  là thuật ngữ lỗi ngẫu nhiên. Thứ tự  $p$  cho mô hình AR( $p$ ) được lựa chọn dựa trên các đỉnh quan trọng của PACF (Partial Autocorrelation Function). Một chỉ báo bổ sung là sự giảm chậm chạp của ACF (Autocorrelation Function).

MA thực hiện dự báo dựa trên các trung bình chuyển động của các thuật ngữ lỗi ngẫu nhiên trước đó như sau:

$$y_t = \mu + \sum_{i=1}^q \theta_i e_{t-i}$$

trong đó  $\theta_i$  đại diện cho các hệ số hồi quy,  $q$  là thứ tự của trung bình chuyển động, và  $\mu$  là một hằng số. Thứ tự  $q$  cho mô hình MA( $q$ ) được lấy từ ACF, nếu nó có một đoạn cắt sắc sau lags  $q$ . PACF giảm chậm trong trường hợp này.

Mô hình ARIMA có thể được xác định như sau:

$$\phi_p(B)(1 - B)^d y_t = \theta_q(B)e_t \quad (7)$$

trong đó  $B$  là toán tử backshift,  $p$  là thứ tự hồi quy,  $d$  là thứ tự chênh lệch và  $q$  là thứ tự của trung bình chuyển động.

## C. VAR

Mô hình Vector Autoregression (VAR) là một công cụ mạnh mẽ trong phân tích dữ liệu thời gian đa biến. Thay vì chỉ tập trung vào một biến duy nhất như các mô hình hồi quy tuyến tính thông thường, VAR cho phép chúng ta đánh giá các tương tác phức tạp giữa nhiều biến số cùng một lúc. Điều này giúp chúng ta hiểu rõ hơn về cách các biến ảnh hưởng lẫn nhau qua thời gian, và tạo ra dự báo linh hoạt dựa trên các kịch bản khác nhau của tương lai. VAR đã trở thành một công cụ không thể thiếu trong lĩnh vực kinh tế lượng, tài chính và các lĩnh vực khác đòi hỏi sự hiểu biết sâu sắc về mối quan hệ giữa các biến số thời gian. Công thức chung của mô hình VAR với  $p$  lags (trễ) được biểu diễn như sau:

$$Y_t = c + A_1 Y_{t-1} + A_2 Y_{t-2} + \dots + A_p Y_{t-p} + \epsilon_t$$

Trong đó:

- $Y_t$ : là một vector chứa các biến phụ thuộc tại thời điểm  $t$ ;
- $c$ : là một vector hằng số;
- $A_1, A_2, \dots, A_p$ : là các ma trận hệ số tương ứng với các lags;
- $\epsilon_t$ : là vector của thành phần sai số tại thời điểm  $t$ .

## D. Random Forest

### Định nghĩa

Random Forest là phương pháp xây dựng một tập hợp (ensemble) cây quyết định (decision tree) cho các bài toán phân lớp (classification) và hồi quy (regression). Cây quyết định thường có độ lệch thấp và phương sai cao, và Random Forest khai thác điều này bằng cách lấy trung bình các cây để cải thiện hiệu suất.

### Thách thức

Thách thức chính của Random Forest là tạo ra tính ngẫu nhiên để giảm mối tương quan giữa các cây ( $\rho(x)$ ) và duy trì độ lệch thấp.

### Thuật toán

Thuật toán “Bagging” (Bootstrap Aggregating) là cơ chế chính của Random Forest. Mỗi cây quyết định được xây dựng từ một tập dữ liệu huấn luyện khác nhau và không bị cắt tỉa, tăng tính đa dạng và giảm lỗi. *Dự đoán hồi quy*: Lấy giá trị trung bình của dự đoán từ tất cả các cây để giảm phương sai và cải thiện hiệu suất. *Dự đoán phân lớp*: Lấy đa số phiếu bầu cho nhãn lớp từ tất cả các cây, giúp khắc phục sai sót của từng cây và tăng độ chính xác.

## E. Time-series Dense Encoder

**TiDE** (Time-series Dense Encoder) là một mô hình dự báo chuỗi thời gian. Mô hình hoạt động mã hóa chuỗi thời gian quá khứ cùng với hiệp phương sai bằng cách sử dụng “dense MLP” (Multi-Layer Perceptron). Sau đó, mô hình giải mã (decode) chuỗi thời gian được mã hóa (encode) cùng với các hiệp phương sai trong tương lai.

Kiến trúc tổng quan được trình bày ở Hình 1. Đầu vào (input) của mô hình là dữ liệu quá khứ và phương sai của một chuỗi thời gian tại một thời điểm  $(y_{1:L}^{(i)}, x_{1:L}^{(i)}, a^{(i)})$  và ánh xạ tới dự đoán của chuỗi thời gian  $\hat{y}_{L+1:L+H}^{(i)}$ . Thành phần chính của mô hình là residual block MLP.

**Residual Block**: Là một thành phần quan trọng của kiến trúc TiDE vì nó cho phép mô hình nắm bắt các tính chất phi tuyến tính vốn có trong dữ liệu chuỗi thời gian, đồng thời duy trì các mối quan hệ tuyến tính nhằm giúp cải thiện hiệu suất dự báo dài hạn.

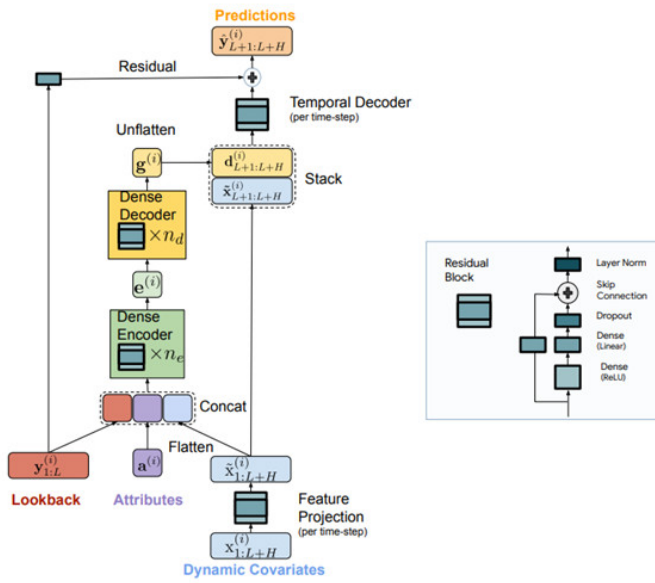
Mô hình được chia thành hai phần: mã hóa (encoding) và giải mã (decoding):

**Mã hóa (Encoding)**: Nhiệm vụ của bước mã hóa là ánh xạ dữ liệu quá khứ và phương sai của chuỗi thời gian thành một biểu diễn dày đặc (dense). Quá trình mã hóa có hai bước chính

**Feature Projection**: Sử dụng residual block để ánh xạ  $x_t^{(i)}$  tại mỗi time-step, hoạt động này được mô tả như sau

$$\tilde{x}_t^{(i)} = \text{ResidualBlock}(x_t^{(i)}) \quad (4)$$

**Giải mã (Decoding)**: Việc giải mã trong mô hình ánh xạ các biểu diễn ẩn được mã hóa thành các dự đoán trong tương lai của chuỗi thời gian. Nó cũng bao gồm hai hoạt động, bộ giải mã dày đặc (dense decoder) và bộ giải mã tạm thời (temporal decoder).



Hình 1. Tổng quan về kiến trúc TiDE

## F. RNN

Mạng nơ-ron hồi tiếp (RNN) là một loại mạng có khả năng xử lý dữ liệu tuần tự. Điểm đặc biệt của RNN là các nơ-ron trong mạng có thể kết nối với nhau theo chu kỳ. Nhờ vậy, RNN có thể học được mối quan hệ giữa các giá trị trong một chuỗi thời gian.

Công thức biểu diễn của RNN như sau:

$$h_t = f(x_t, h_{t-1}) \quad (5)$$

Trong đó:

- $h_t$ : trạng thái của RNN tại thời điểm  $t$
- $x_t$ : đầu vào của RNN tại thời điểm  $t$
- $h_{t-1}$ : trạng thái của RNN tại thời điểm  $t - 1$
- $f$ : hàm kích hoạt của RNN

## V. KẾT QUẢ THÍ NGHIỆM

### A. Cài đặt mô hình

- 1) LINEAR REGRESSION:
- 2) ARIMA:

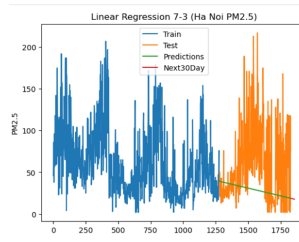
- 3) VAR:

### B. Đánh giá

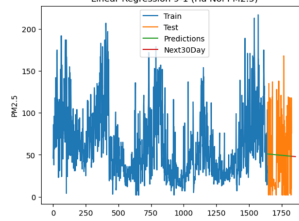
## VI. KẾT LUẬN

### TÀI LIỆU

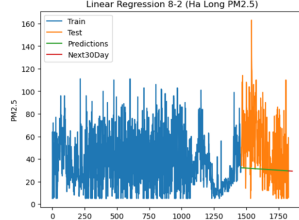
- [1] E. Marinov, D. Petrova-Antonova, and S. Malinov, "Time Series Forecasting of Air Quality: A Case Study of Sofia City," *Atmosphere*,



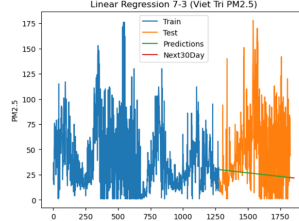
Hình 2. Linear Regression (7:3) PM2.5 Hà Nội



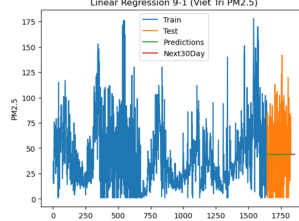
Hình 4. Linear Regression (9:1) PM2.5 Hà Nội



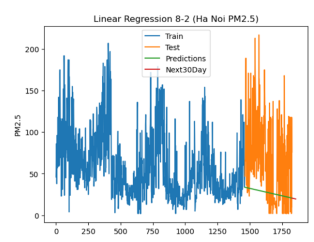
Hình 6. Linear Regression (8:2) PM2.5 Hạ Long



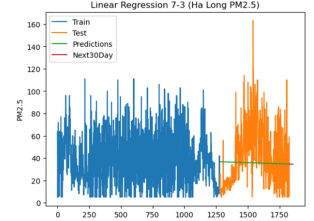
Hình 8. Linear Regression (7:3) PM2.5 Việt Trì



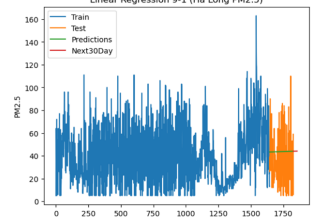
Hình 10. Linear Regression (9:1) PM2.5 Việt Trì



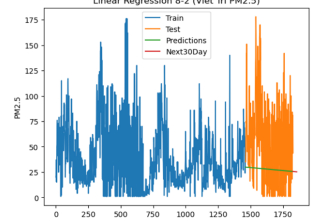
Hình 3. Linear Regression (8:2) PM2.5 Hà Nội



Hình 5. Linear Regression (7:3) PM2.5 Hạ Long



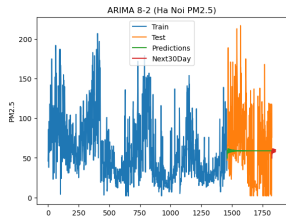
Hình 7. Linear Regression (9:1) PM2.5 Hạ Long



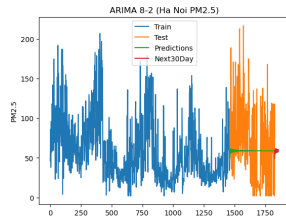
Hình 9. Linear Regression (8:2) PM2.5 Việt Trì

vol. 13, p. 788, 2022. [Online]. Available:<https://www.mdpi.com/2073-4433/13/5/788>

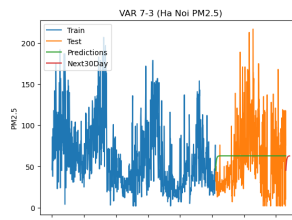
- [2] K. N. Sh, I. Irfani, and U. Mukhaiyar, "Predicting Air Pollution Levels in Jakarta Using Vector Autoregressive Analysis," *Proceedings of the 5th International Conference on Statistics, Mathematics, Teaching, and Research 2023 (ICSMTTR 2023)*, pp. 14-22, Atlantis Press, 2023. [Online]. Available:[https://doi.org/10.2991/978-94-6463-332-0\\_3](https://doi.org/10.2991/978-94-6463-332-0_3)



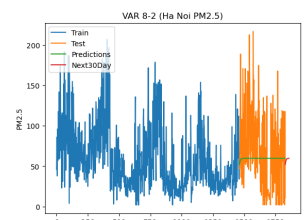
Hình 11. ARIMA (7:3) PM2.5 Hà Nội



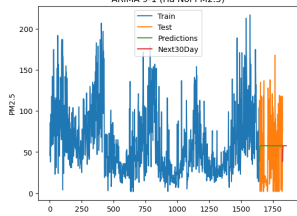
Hình 12. ARIMA (8:2) PM2.5 Hà Nội



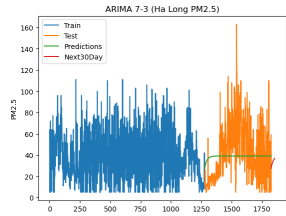
Hình 20. VAR (7:3) PM2.5 Hà Nội



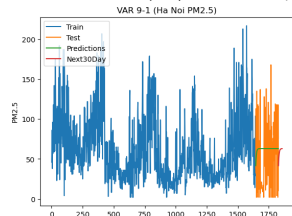
Hình 21. VAR (8:2) PM2.5 Hà Nội



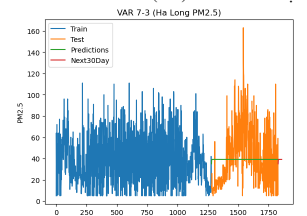
Hình 13. ARIMA (9:1) PM2.5 Hà Nội



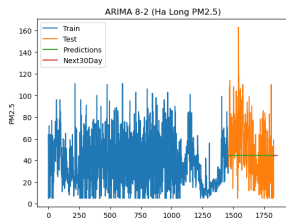
Hình 14. ARIMA (7:3) PM2.5 Hà Nội



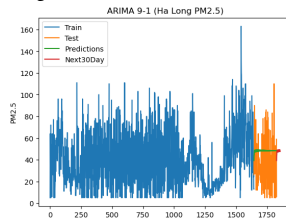
Hình 22. VAR (9:1) PM2.5 Hà Nội



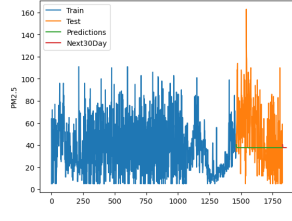
Hình 23. VAR (7:3) PM2.5 Hà Nội



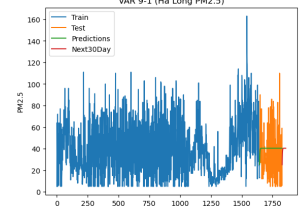
Hình 15. ARIMA (8:2) PM2.5 Hà Nội



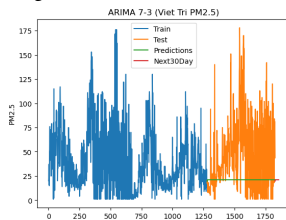
Hình 16. ARIMA (9:1) PM2.5 Hà Nội



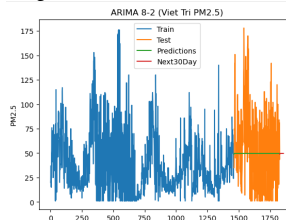
Hình 24. VAR (8:2) PM2.5 Hà Nội



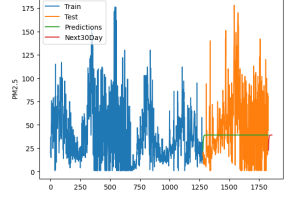
Hình 25. VAR (9:1) PM2.5 Hà Nội



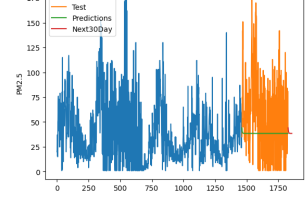
Hình 17. ARIMA (7:3) PM2.5 Việt Trì



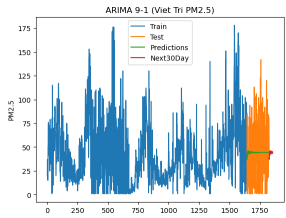
Hình 18. ARIMA (8:2) PM2.5 Việt Trì



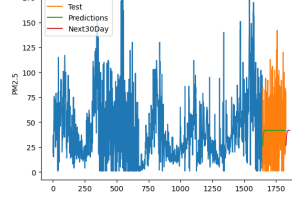
Hình 26. VAR (7:3) PM2.5 Việt Trì



Hình 27. VAR (8:2) PM2.5 Việt Trì



Hình 19. ARIMA (9:1) PM2.5 Việt Trì



Hình 28. VAR (9:1) PM2.5 Việt Trì

- [3] K. H. Waseem, H. Mushtaq, F. Abid, A. M. Abu-Mahfouz, A. Shaikh, M. Turan, and J. Rasheed, "Forecasting of Air Quality Using an Optimized Recurrent Neural Network," *Processes*, vol. 10, no. 10, p. 2117, 2022. [Online]. Available: <https://doi.org/10.3390/pr1010211>
- [4] Citation: Esager, M.W.M.; Ünlü, K.D. "Forecasting Air Quality in Tripoli: An Evaluation of Deep Learning Models for Hourly PM2.5 Surface Mass Concentrations. *Atmosphere* 2023, 14, 478. [Online]. Available: <https://doi.org/10.3390/atmos14030478>
- [5] A. Das, W. Kong, A. Leach, S. Mathur, R. Sen, and R. Yu, "Long-term Forecasting with TiDE: Time-series Dense Encoder," *arXiv:2304.08424 [stat.ML]*, April 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2304.08424>

- [6] B. T. Ong, K. Sugiura, and K. Zettsu, "Dynamically pre-trained deep recurrent neural networks using environmental monitoring data for predicting PM2.5," *Neural Comput Appl*, vol. 27, no. 6, pp. 1553–1566, 2016. [Online]. Available: <https://doi.org/10.1007/s00521-015-1955-3>
- [7] Y. B. Lim, I. Aliyu, and C. G. Lim, "Air pollution matter prediction using recurrent neural networks with sequential data," In: *Proceedings of the 2019 3rd International Conference on Intelligent Systems, Metaheuristics & Swarm Intelligence. ISMSI 2019*, pp. 40–44, Association for Computing Machinery, New York, NY, USA, 2019. [Online]. Available: <https://doi.org/10.1145/3325773.3325788>
- [8] H. Wu, J. Xu, J. Wang, & M. Long, "Autoformer: Decomposition Transformers with Auto-Correlation for Long-Term Series Forecasting," 2021.