

THU THẬP VÀ PHÂN TÍCH DỮ LIỆU VÉ MÁY BAY ỨNG DỤNG HỌC MÁY TRONG VIỆC RA QUYẾT ĐỊNH MUA VÉ MÁY BAY

Môn: Tiền xử lý và xây dựng bộ dữ liệu

Lớp: DS108.P21

GVHD: TS. Nguyễn Gia Tuấn Anh

CN. Trần Quốc Khánh

SVTH: Đinh Bảo Thy - 23521563

Võ Ngọc Anh Thy - 23521565

Nguyễn Vũ Thùy Trâm - 23521617



#### DS108 - Tiền xử lý và xây dựng bộ dữ liệu

# Nội dung báo cáo

Tổng quan	01
Thu thập và tiền xử lý bộ dữ liệu	02
Xây dựng mô hình	03
Thực nghiệm và đánh giá	04
Kết luận và hướng phát triển	05

### Tổng quan

#### Mục tiêu ----

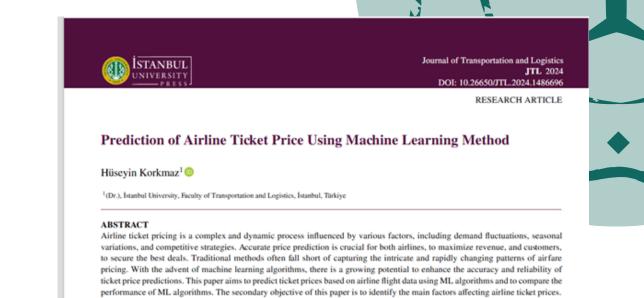
- Hệ thống dự đoán giá vé theo thời gian
- Khuyến nghị về thời điểm mua vé (mua ngay hoặc chờ)

#### 

- Tính phù hợp
- Tính cập nhật
- Tính chính xác

#### Khác biệt ———

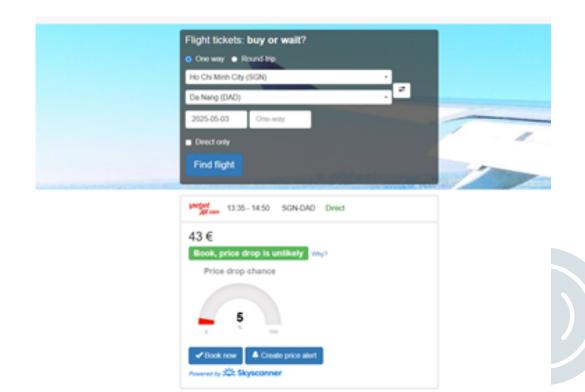
- Bao gồm: dự đoán + hệ khuyến nghị
- Theo nhu cầu khách hàng



The flight and ticket price datasets of THY and PGS that were obtained from open-access sources are used in this paper. The final dataset consists of 962 records for three months from June 1st, 2022 to August 30th, 2022 and includes 19 different variables. Statistical tests and ML algorithms were applied to the final dataset. This paper compares various ML models to predict airline ticket prices, considering performance metrics such as MAE, MSE, RMSE, and R2 during training and test phases. According to the model training and test results, the best algorithm is GPR with R2: 0.86 (training) and R2: 0.90 (test). The findings are consistent

with existing literature, further validating the superior efficacy of certain models in specific contexts and demonstrating significant progress in the field. This paper contributes to the literature by comparing the effectiveness of various machine learning algorithms in predicting airline ticket prices, providing new and valuable insights into model performance and key price-determining factors.

Keywords: Price Prediction, Ticket Price, Airfare Price, Machine Learning, Intelligent Transportation Systems



#### Bộ dữ liệu

• Website: traveloka

• Crawl dữ liệu: Selenium WebDriver

- Địa điểm: 6 thành phố lớn, du lịch ở Việt Nam
- Tính đa dạng: is\_holiday
  (ngày lễ cận lễ cuối tuần ngày thường)
- Bronze layer: 46,826 samples với 12 thuộc tính



brand	id	price	start_time	start_day	end_day	end_time	trip_time	destination	hand_luggage	checked_baggage	crawl_date
Bamboo Airways	PQC0002	1.245.267 VND/khách	16:30	02 thg 5	02 thg 5	17:30	1h 0m	Phú Quốc (PQC)\nSân bay Phú Quốc	Hành lý xách tay 7 kg	Hành lý 0 kg	2025-05- 01
Vietravel Airlines	PQC0323	1.224.930 VND/khách	17:25	02 thg 5	02 thg 5	18:25	1h 0m	Phú Quốc (PQC)\nSân bay Phú Quốc	Hành lý xách tay 7 kg	Hành lý 0 kg	2025-05- 01
Vietnam Airlines	PQC0292	1.535.463 VND/khách	17:25	02 thg 5	02 thg 5	18:25	1h 0m	Phú Quốc (PQC)\nSân bay Phú Quốc	Hành lý xách tay 1 x 12 kg	Hành lý 23 kg	2025-05- 01
VietJet Air	PQC0109	1.567.000 VND/khách	15:40	02 thg 5	02 thg 5	16:40	1h 0m	Phú Quốc (PQC)\nSân bay Phú Quốc	Hành lý xách tay 7 kg	Hành lý 0 kg	2025-05- 01
VietJet Air	PQC0132	1.567.000 VND/khách	17:10	02 thg 5	02 thg 5	18:15	1h 5m	Phú Quốc (PQC)\nSân bay Phú Quốc	Hành lý xách tay 7 kg	Hành lý 0 kg	2025-05- 01

# Tiền xử lý dữ liệu

Xử lý tổng thể

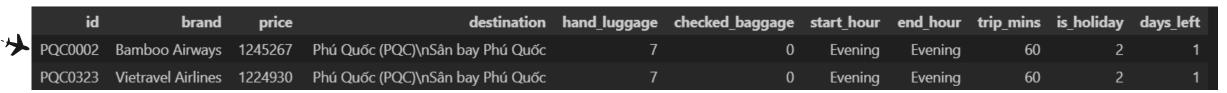
Chuẩn hóa kiểu dữ liệu

Trích xuất đặc trưng

- Loại bỏ các dòng trùng lặp
- Xử lý missing value

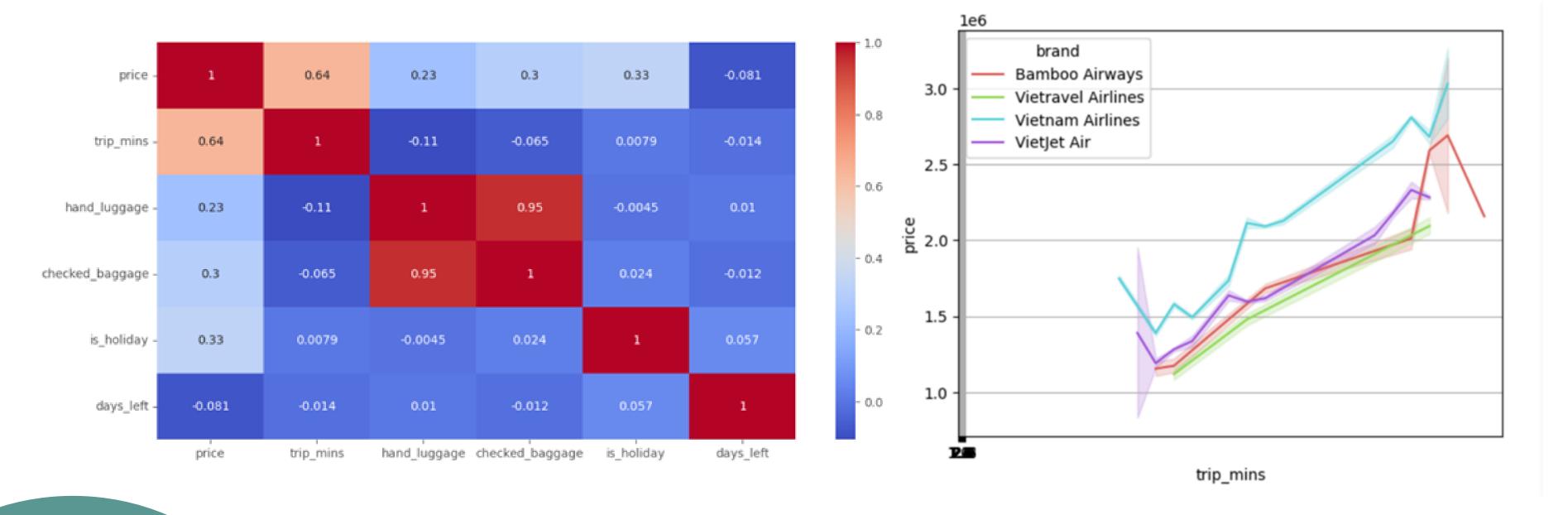
- Chuyển đổi kiểu dữ liệu phù hợp, đồng nhất định dạng cho các đặc trưng
- Tách và lấy các thông tin quan trong
- Phân loại thời gian theo các nhóm
- Thêm đặc trưng
- Xác định và gán nhãn các ngày đặc biệt

brand	id	price	start_time	start_day	end_day	end_time	trip_time	destination	hand_luggage	checked_baggage	crawl_date
Bamboo Airways	PQC0002	1.245.267 VND/khách	16:30	02 thg 5	02 thg 5	17:30	1h 0m	Phú Quốc (PQC)\nSân bay Phú Quốc	Hành lý xách tay 7 kg	Hành lý 0 kg	2025-05- 01
Vietravel Airlines	PQC0323	1.224.930 VND/khách	17:25	02 thg 5	02 thg 5	18:25	1h 0m	Phú Quốc (PQC)\nSân bay Phú Quốc	Hành lý xách tay 7 kg	Hành lý 0 kg	2025-05- 01



→ Silver layer: 46,549 samples với 11 thuộc tính

• Độ tương quan giữa các dữ liệu dạng số

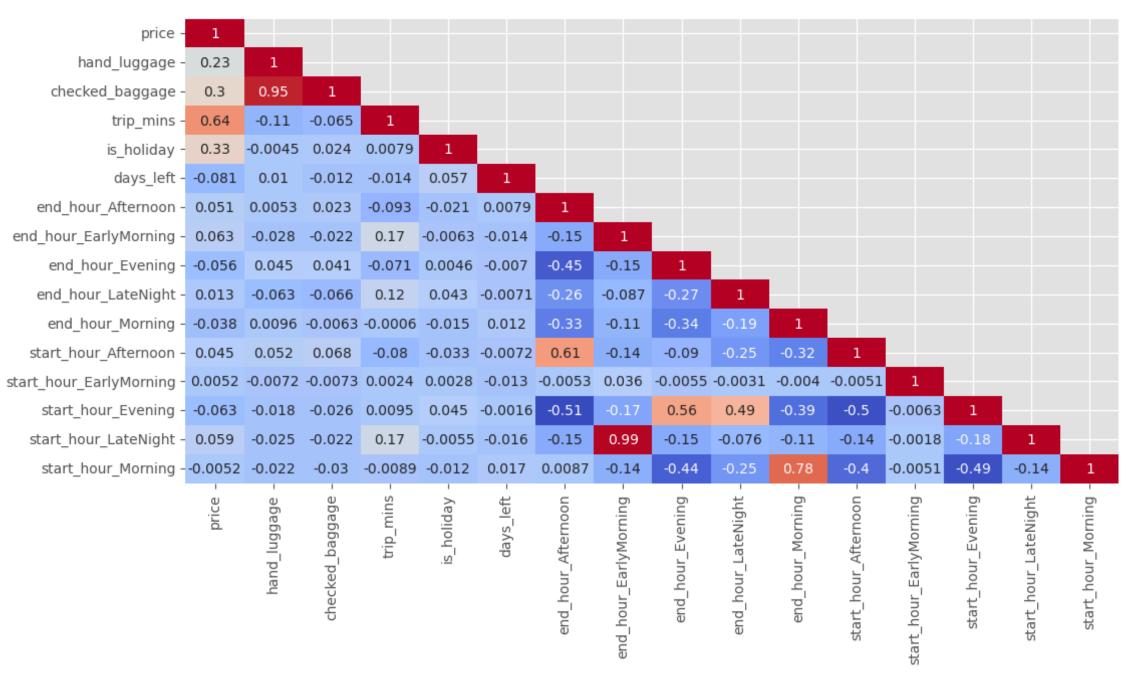


→ Giữa 2 biến **price** và **trip\_mins** có sự tương quan đồng biến





• Độ tương quan giữa các thuộc tính







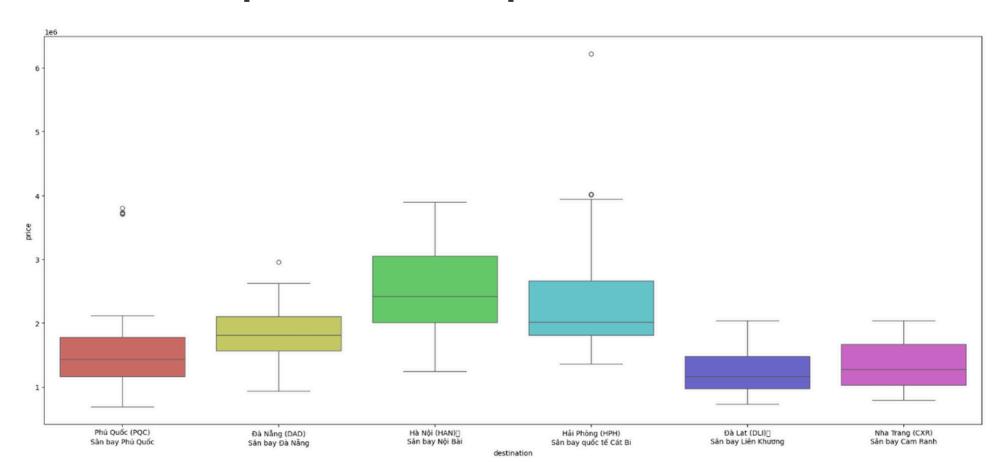
- 0.8

- 0.6

- 0.4

- 0.2

- 0.0

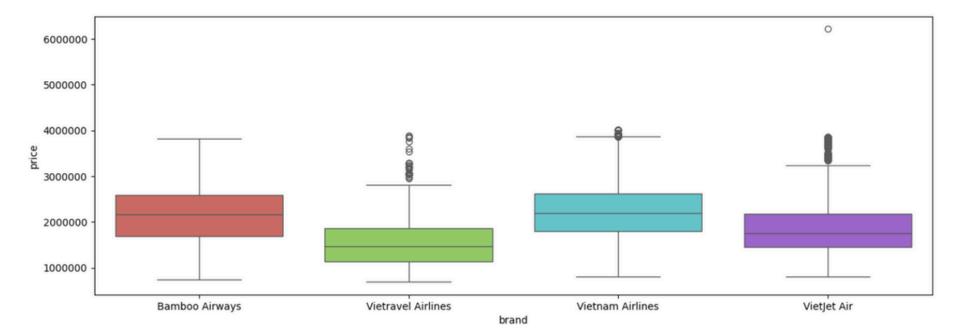


Phân bố giá vé theo điểm đến:

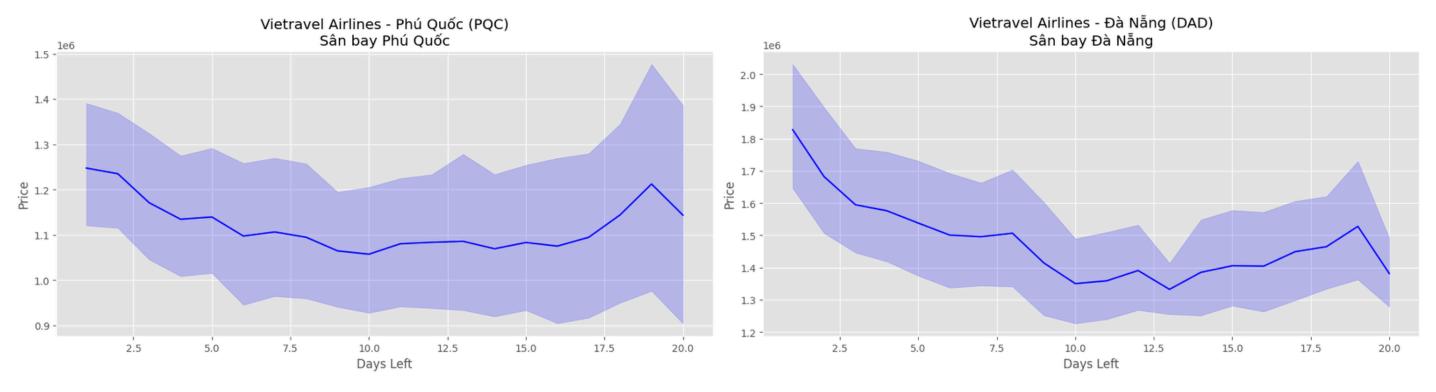
• Có sự khác biệt giữa các sân bay lớn và nhỏ

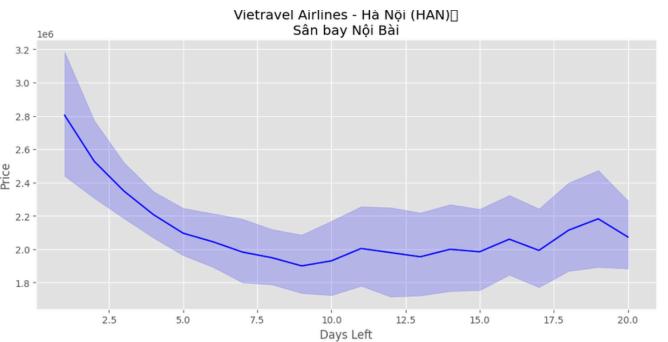
Phân bố giá vé theo hãng bay:

• Phân bố có sự tương đồng

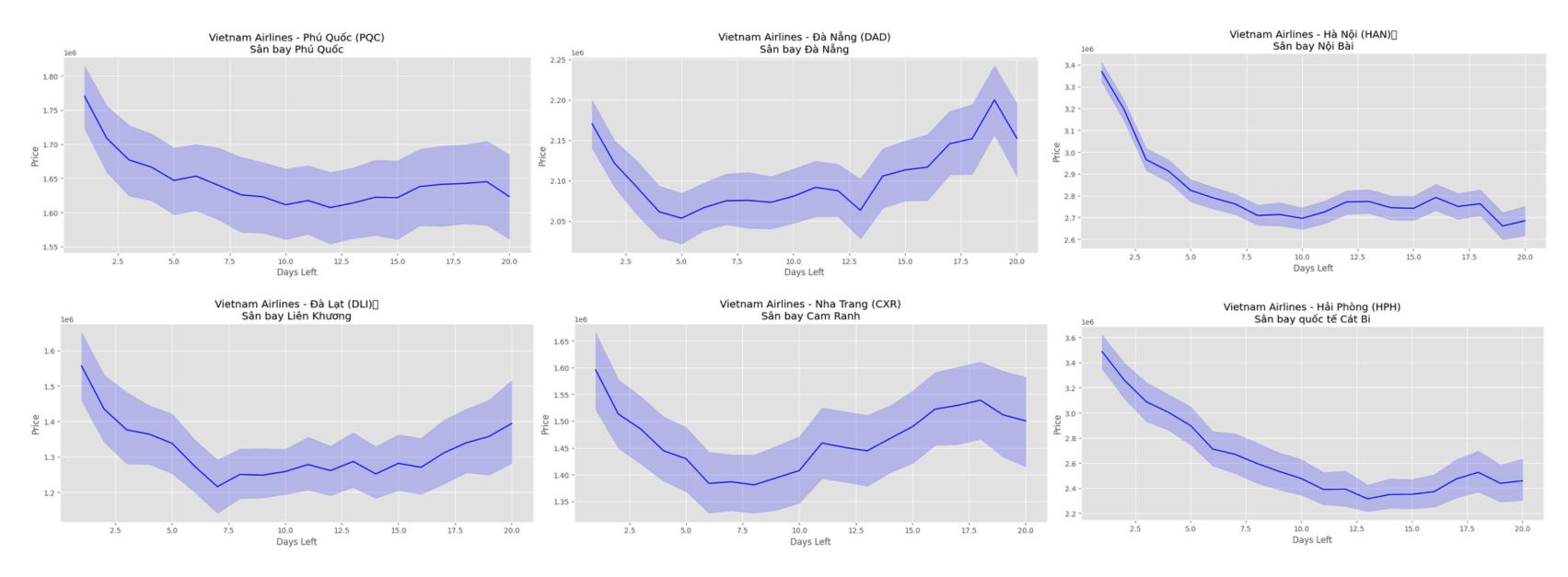


Biểu đồ phân bố price và days\_left của brand theo destination





Biểu đồ phân bố price và days\_left của brand theo destination



- Theo brand, giá vé tăng, giảm theo quy luật
- → Huấn luyện mô hình theo brand có khả năng cho kết quả tốt hơn theo từng destination.

## Tách tập train/test



Test\_data:

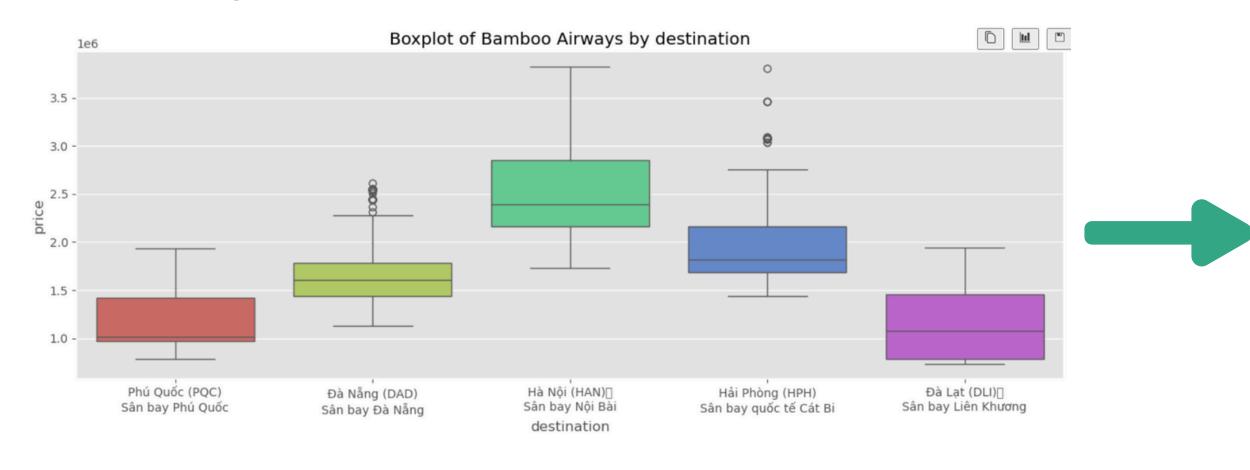


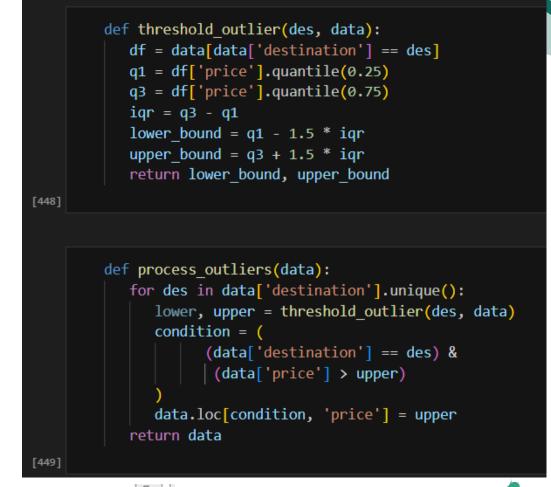


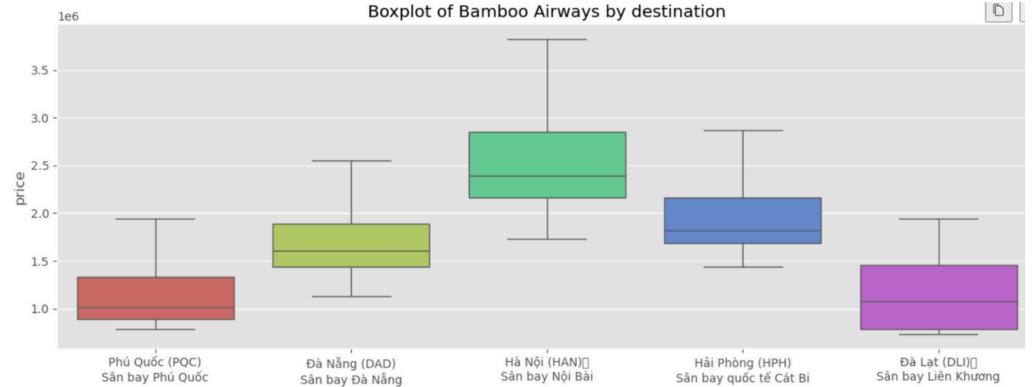


Train= bộ dữ liệu - Test

### Xử lý outliers:







destination

# Xây dựng mô hình



	id	price	destination	hand_luggage	checked_baggage	trip_mins	is_holiday	days_left
(	) PQC0002	1245267.0	Phú Quốc (PQC)\nSân bay Phú Quốc	7	0	60	2	1
	PQC0003	869953.0	Phú Quốc (PQC)\nSân bay Phú Quốc	7	0	60	2	2
2	PQC0004	1013756.0	Phú Quốc (PQC)\nSân bay Phú Quốc	7	0	60	2	3

→ Gold layer: 4 tập dữ liệu theo từng brand với 8 thuộc tính

CHUẨN HÓA DỮ LIỆU

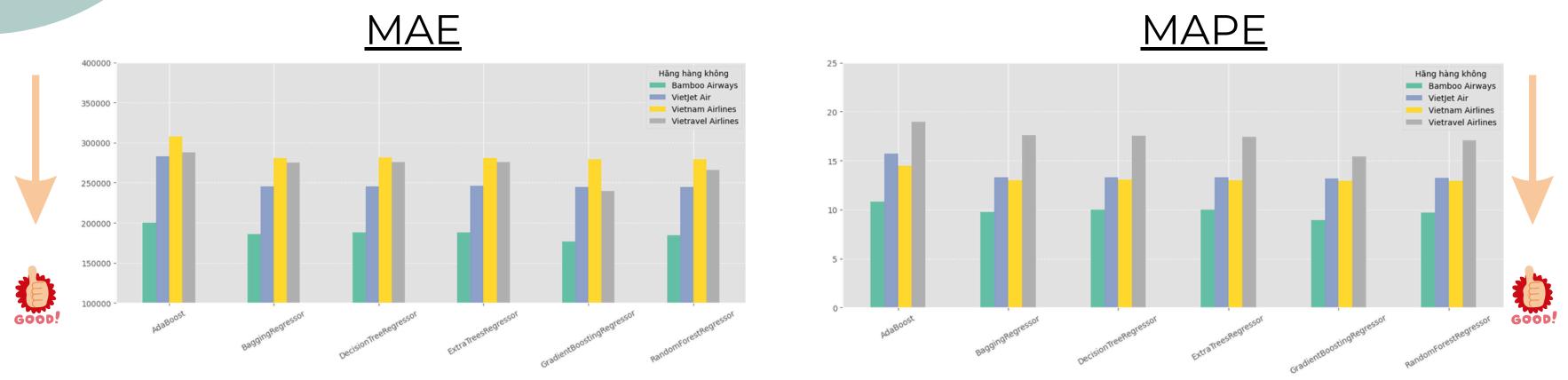
**Standard Scaler** 

ONE-HOT ENCODING

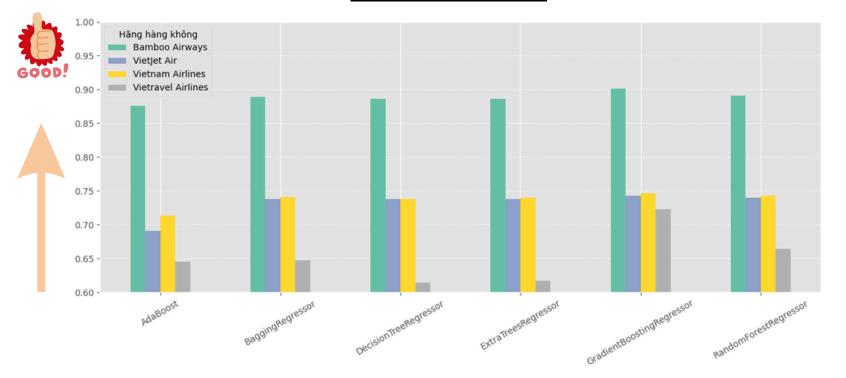
destination



# Xây dựng mô hình





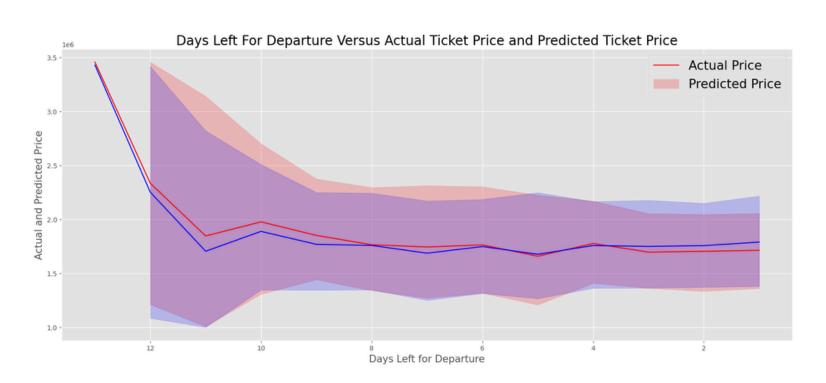




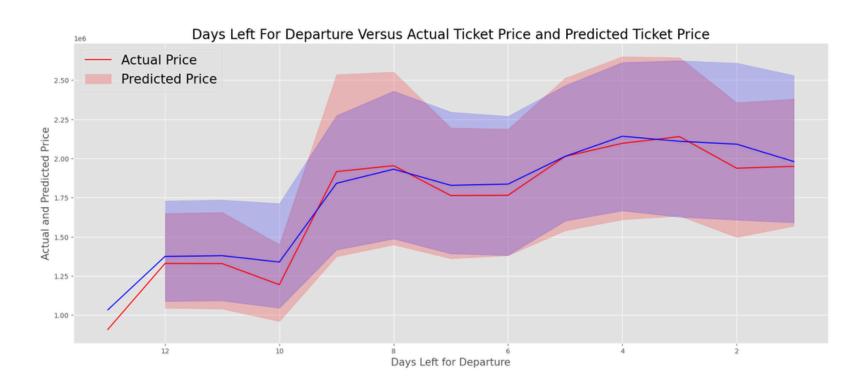




#### **Bamboo Airways**

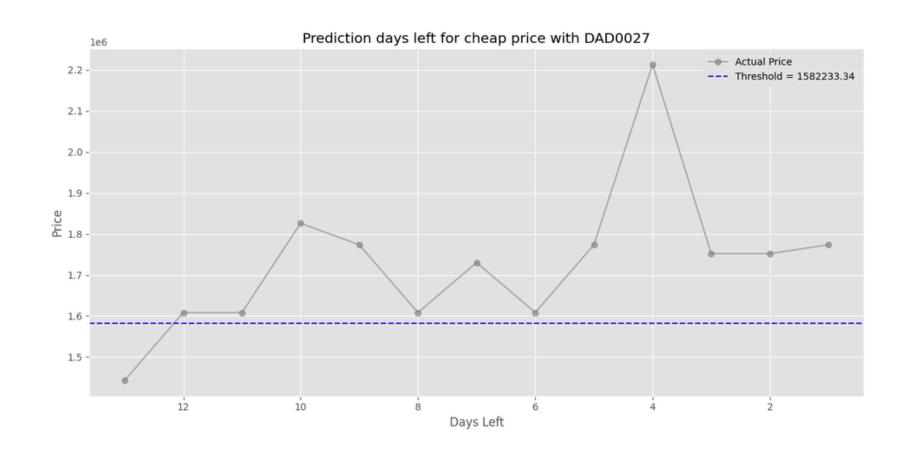


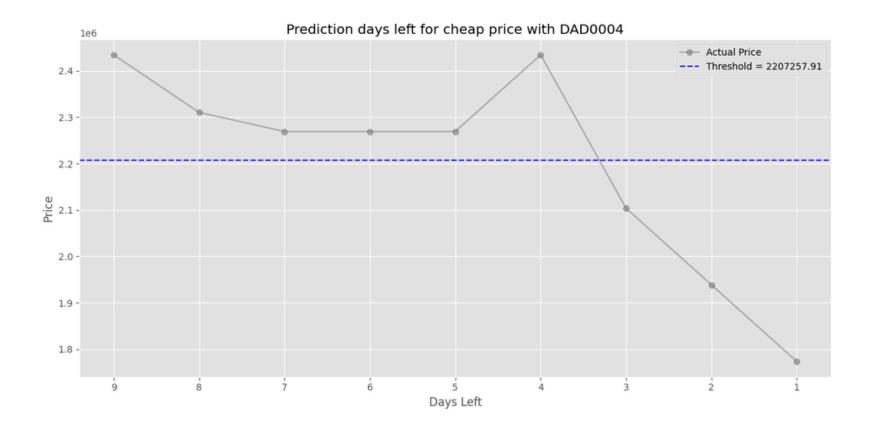
#### **Vietnam Airlines**

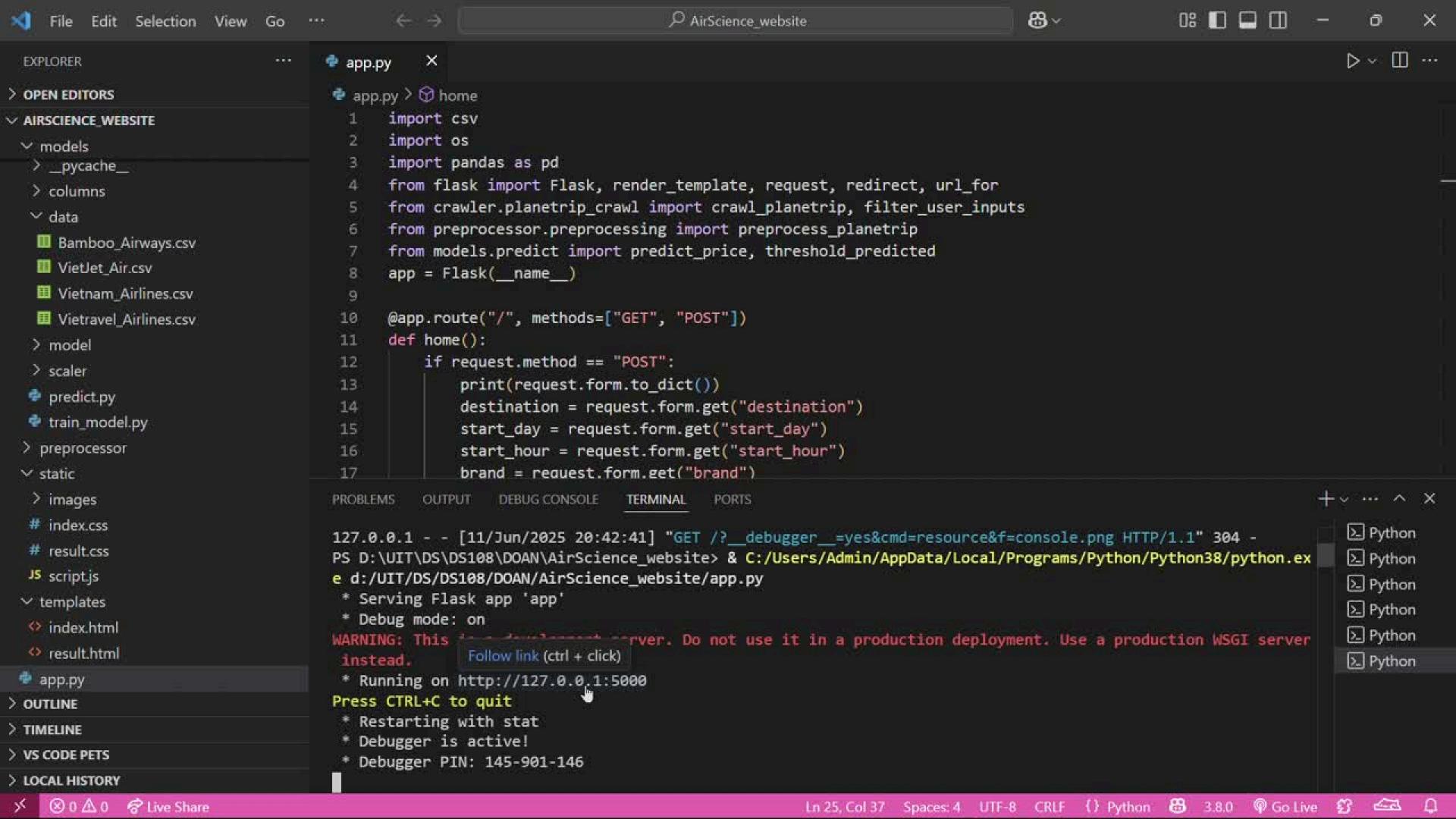


Mô hình Gradient Boosting Regressor

# - Thực nghiệm & Đánh giá













- Các chuyển bay giới hạn trong 20 ngày trước ngày bay
- Chuyến bay chỉ xuất phát từ HCM → 6 thành phố khác
- Không tự động hoàn toàn do reCAPTCHA



#### • HƯỚNG PHÁT TRIỂN:

- Thời gian đợi từ nhập dữ liệu → dự đoán
- Học tăng cường



