



ĐỒ ÁN TỐT NGHIỆP DATA SCIENCE (DL07_308T37_ON)

DỰ ĐOÁN VÀ XÁC ĐỊNH BẤT THƯỜNG GIÁ XE MÁY TRÊN NỀN TẢNG CHỢ TỐT

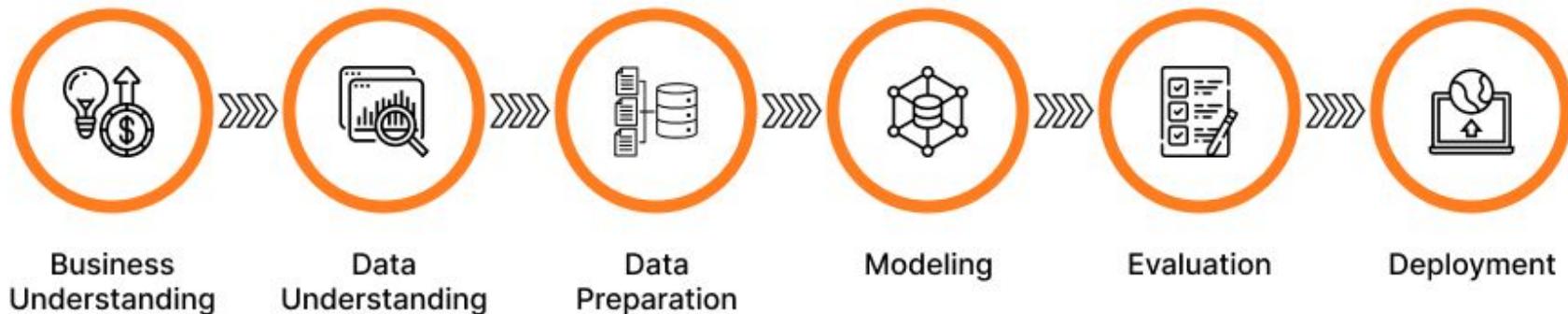
GVHD: Th.S. KHUẤT THÙY PHƯƠNG

HVTH: NGUYỄN MAI XUÂN BÁCH
VÕ THỊ HOÀNG ANH

Tp.HCM, 29/11/2025

CONTENT

DATA SCIENCE PROCESS

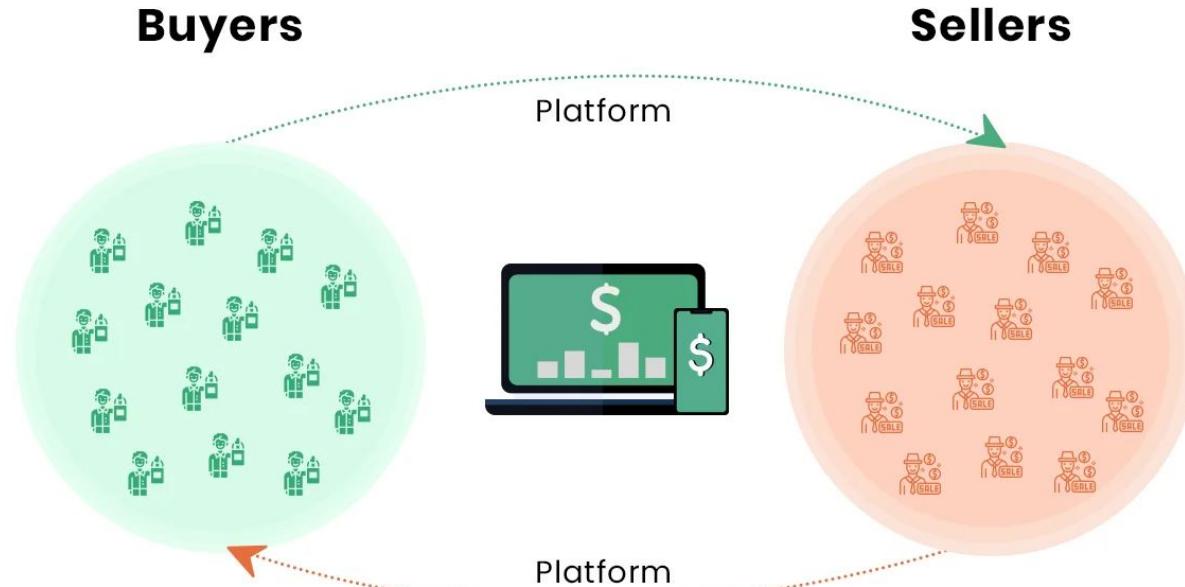




Chợ Tốt – Nền tảng mua bán trực tuyến lớn tại Việt Nam



DATA – Thông tin về xe máy cũ được đăng bán trên nền tảng Chợ Tốt.



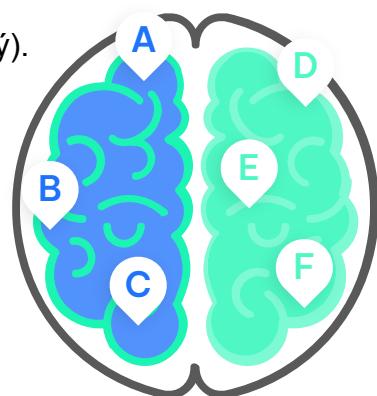
BÀI TOÁN & MỤC TIÊU

⌚ Bối cảnh & Vấn đề

- ❖ Giá xe máy cũ trên thị trường biến động lớn theo hãng, dòng xe, năm sản xuất, tình trạng và số km đã dùng.
- ❖ Người bán thường định giá theo cảm tính → giá thiếu nhất quán, chênh lệch so với giá trị thực.
- ❖ Xuất hiện nhiều tin đăng bất thường: giá quá rẻ (nguy cơ lừa đảo) hoặc quá cao (không hợp lý).

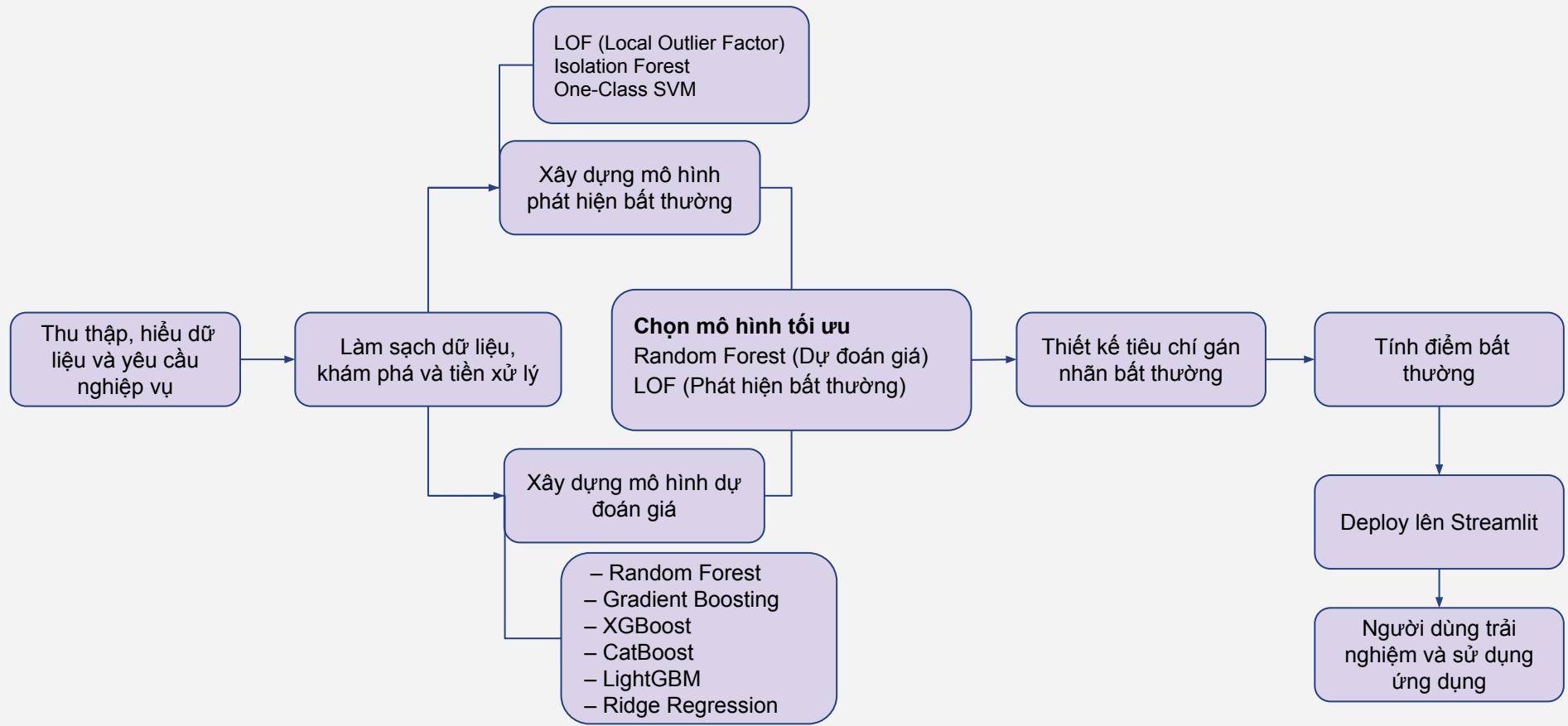
⌚ Mục tiêu của dự án

- Xây dựng mô hình dự đoán mức giá hợp lý cho từng tin đăng dựa trên dữ liệu thực.
- Áp dụng phát hiện bất thường (Anomaly Detection) để nhận diện các tin đăng có giá sai lệch.
- Giúp người mua tham khảo giá chính xác hơn và nền tảng có cơ sở kiểm duyệt các tin bất thường.



QUY TRÌNH THỰC HIỆN

Từ tiền xử lý dữ liệu tới ứng dụng cho người dùng cuối



DATA UNDERSTANDING

id	Tiêu đề	Giá	Khoảng giá min	Khoảng giá max	Địa chỉ	Mô tả chi tiết	Thương hiệu	Dòng xe	Năm đăng ký	Số Km đã đi	Tình trạng	Loại xe	Dung tích xe	Xuất xứ	Chính sách bảo hành	Trọng lượng	Href
0 1	Bán Vespa Sprint 125cc 2024 xanh dương, xe đẹp...	66.000.000 đ	72.53 tr	85.14 tr	Phường Bến Thành, Quận 1, Tp Hồ Chí Minh	Bán xe #Vespa Sprint 125cc. Mua mới tại #Topco...	Piaggio	Vespa	2024	14000	Đã sử dụng	Tay ga	100 - 175 cc	Đang cập nhật	Bảo hành	> 50 kg	https://xe.chotot.com/mua-ban-xe-may-quan-1-tp...
1 2	🔥 SH 150i Tháng ABS 2019 Bstp Chính Chủ	79.500.000 đ	62.76 tr	73.68 tr	Phường Tân Định, Quận 1, Tp Hồ Chí Minh	Bán SH 150i Tháng ABS 2019 Xám Bạc, Ủy Team X...	Honda	SH	2019	28000	Đã sử dụng	Tay ga	100 - 175 cc	Đang cập nhật	Bảo hành	> 50 kg	https://xe.chotot.com/mua-ban-xe-may-quan-1-tp...
2 3	CC Vision Thể Thao 2023 Đen+bộ đèn Demi audi A7	37.000.000 đ	28 tr	32.86 tr	Phường Cầu Kho, Quận 1, Tp Hồ Chí Minh	Chính chủ bán Vision phiên bản Thể Thao 2023 Đ...	Honda	Vision	2023	12000	Đã sử dụng	Tay ga	100 - 175 cc	Đang cập nhật	Bảo hành	> 50 kg	https://xe.chotot.com/mua-ban-xe-may-quan-1-tp...

- Dataset có 7208 dòng và 18 cột.
- Một vài cột chính:
- `tiêu đề`, `giá`, `khoảng giá min/max`, `địa chỉ`, `thương hiệu`, `dòng xe`, `năm đăng ký`, `tình trạng`, `loại xe`, `dung tích xe`, `xuất xứ`, `href` ...

Nhận xét: Dữ liệu dạng văn bản mô tả (Mô tả chi tiết, Địa chỉ, Tiêu đề) có thể chứa hashtag, ký tự đặc biệt, emoji 🔥 ... → Cần làm sạch để sử dụng trong phân tích NLP hoặc trích xuất thông tin.

Dữ liệu số (Giá, Km đã đi, Khoảng giá) có thể cần chuyển đổi về dạng số thực (float) để phân tích, vì hiện tại nhiều giá trị ở dạng chuỗi có ký tự "đ", "tr".



EDA – PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

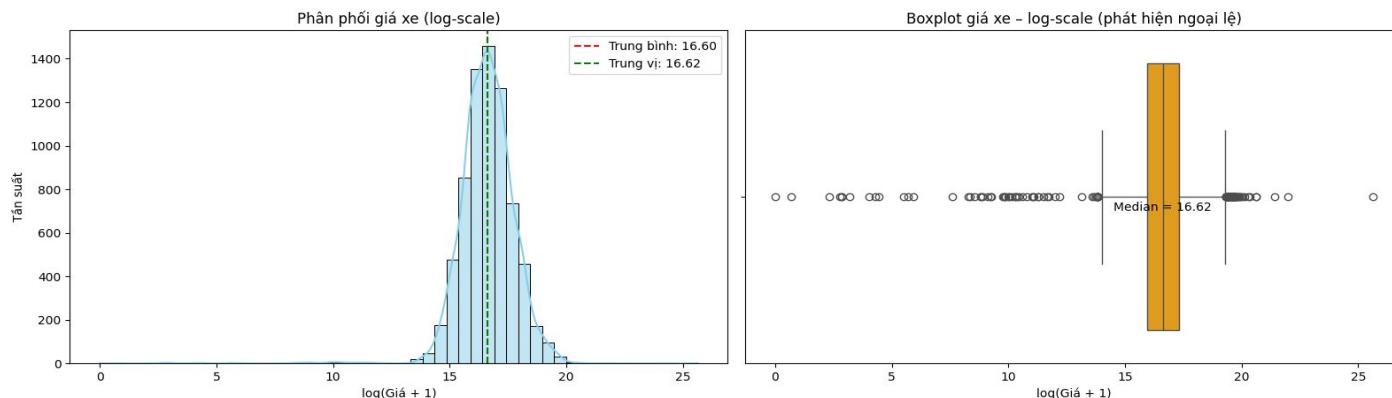
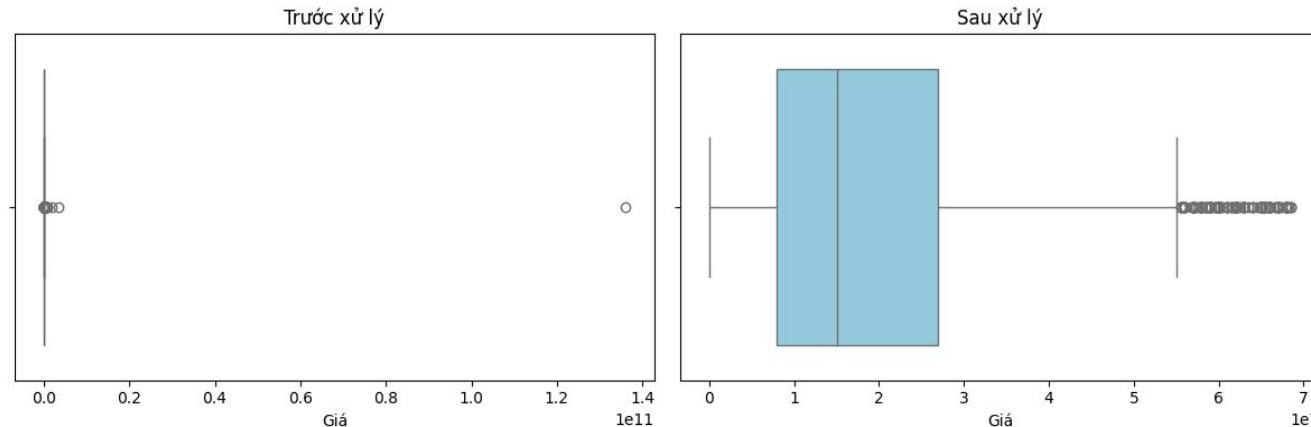
Thống kê về giá

THỐNG KÊ GIÁ TRƯỚC XỬ LÝ

```
count           7,203
mean      49,237,142
std       1,603,410,679
min            0
25%     8,500,000
50%    16,500,000
75%   32,500,000
max  136,000,000,000
Name: Giá, dtype: object
```

THỐNG KÊ GIÁ SAU XỬ LÝ

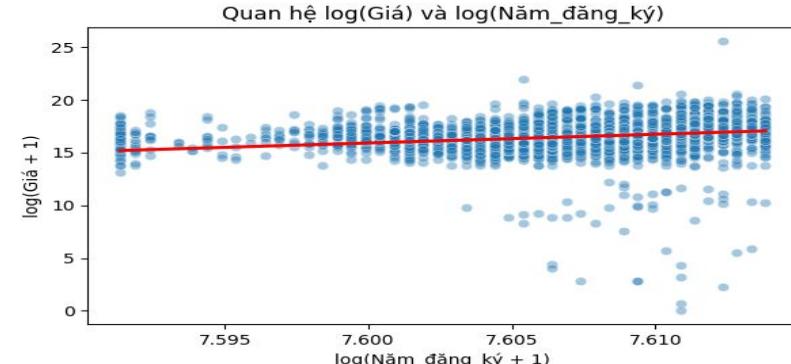
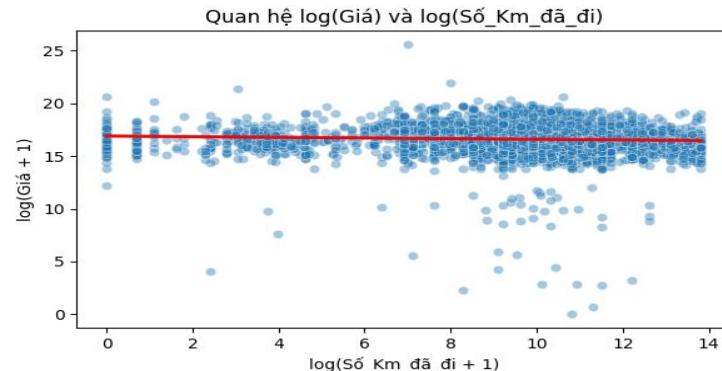
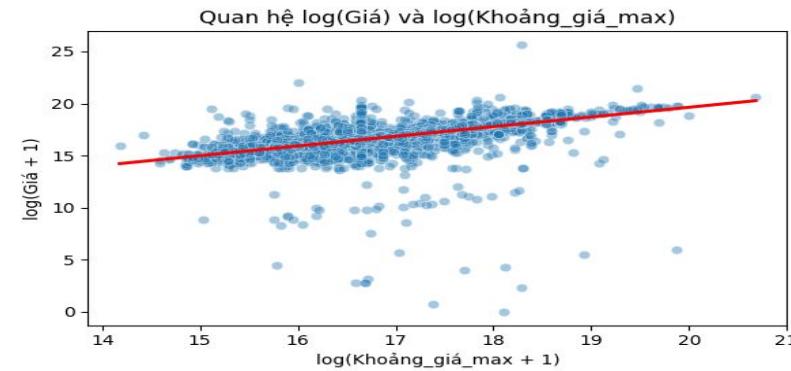
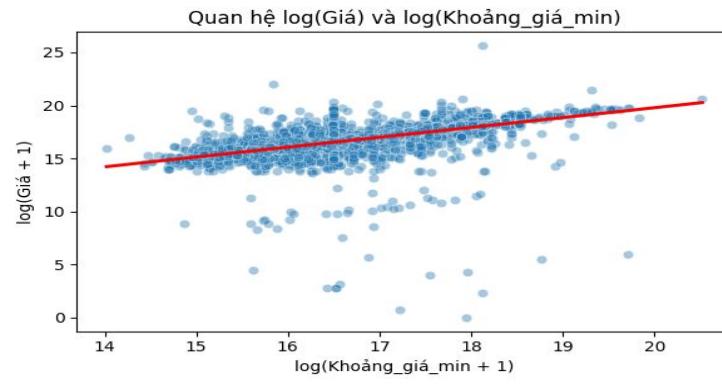
```
count           6,544
mean      19,200,902
std       15,044,528
min            0
25%     8,000,000
50%    15,000,000
75%   26,922,500
max   68,500,000
Name: Giá, dtype: object
```





EDA – PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

Các biến số (numeric variable) ảnh hưởng đến biến giá (xài log)

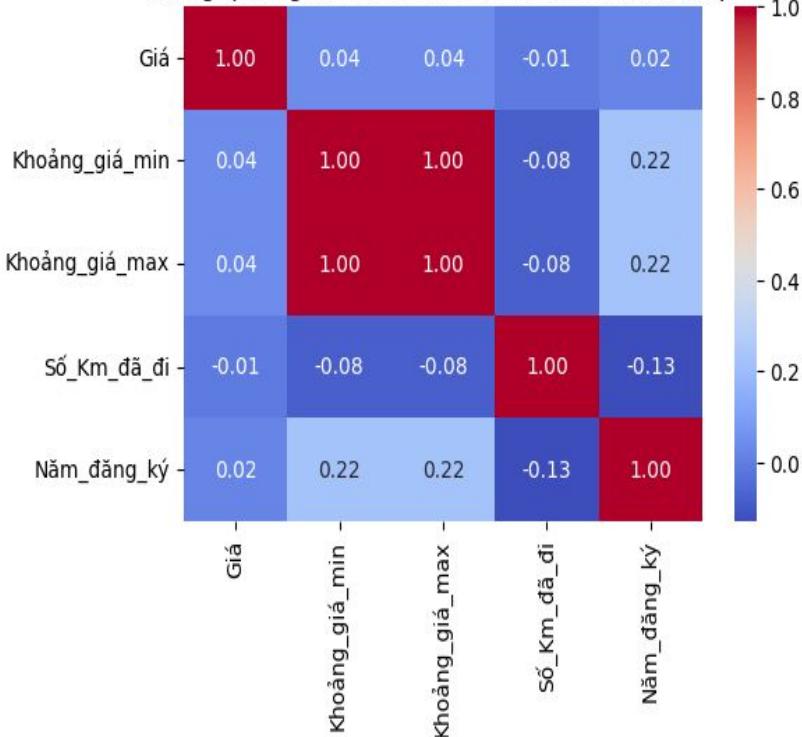




EDA – PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

Các biến số (numeric variable)

Tương quan giữa các biến số (Correlation Heatmap)

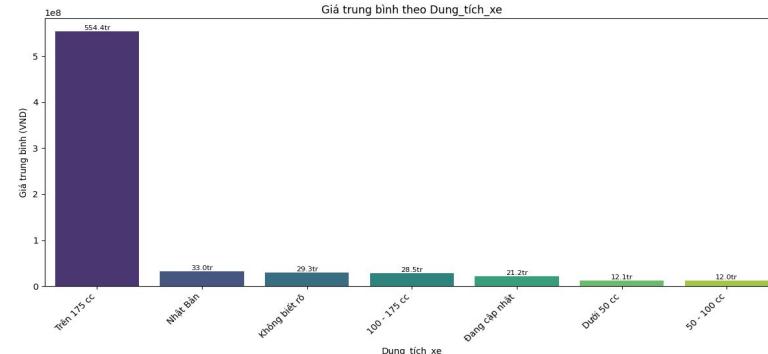
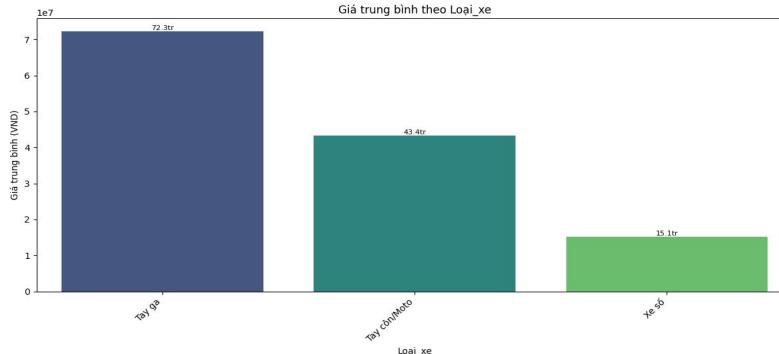
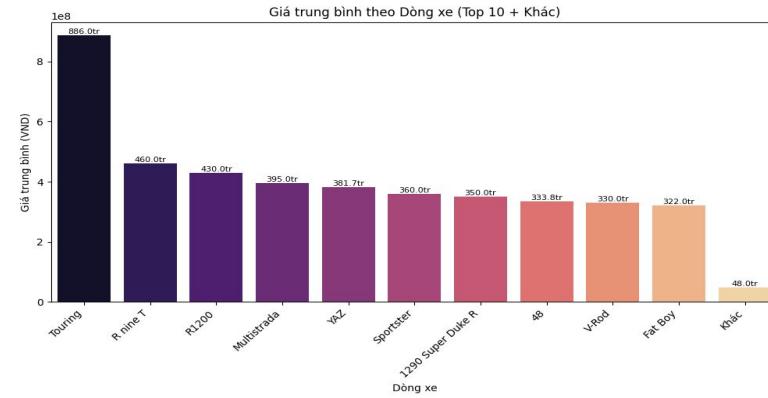
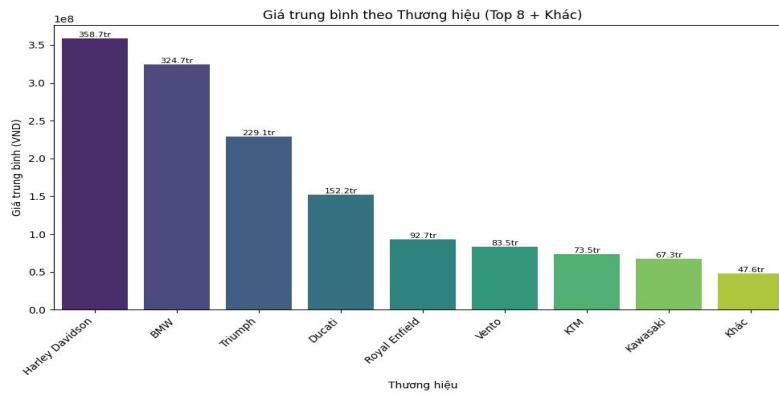


Cặp biến	Hệ số tương quan (r)	Nhận xét
Giá ↔ Khoảng_giá_min	0.04	Quan hệ rất yếu, gần như không tuyến tính
Giá ↔ Khoảng_giá_max	0.04	Tương tự, không có mối quan hệ tuyến tính rõ ràng
Giá ↔ Số_Km_dã đi	-0.01	Cực kỳ yếu, số km gần như không ảnh hưởng đến giá
Giá ↔ Năm_dăng_ký	0.02	Gần như không có mối quan hệ tuyến tính đáng kể



EDA – PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

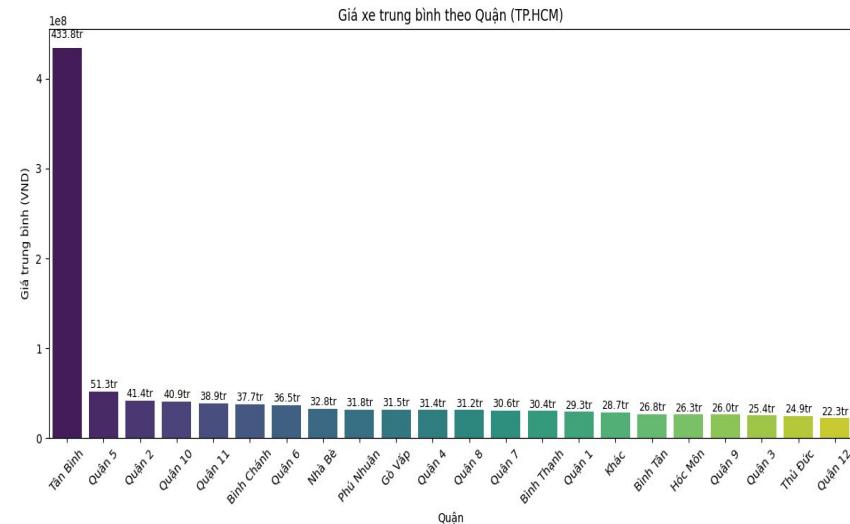
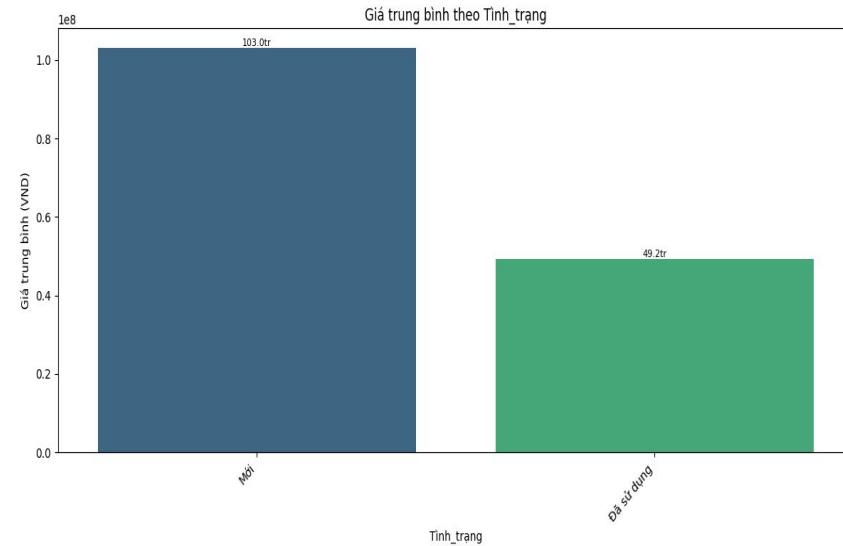
Các biến phân loại (categorical variable) ảnh hưởng đến biến giá





EDA – PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

Các biến phân loại (categorical variable) ảnh hưởng đến biến giá





EDA – PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

Kiểm định Chi-square giữa các biến phân loại

Dựa vào giá trị $p\text{-value}$, với ngưỡng $\alpha = 0.05$:

- Nếu $p\text{-value} < 0.05 \rightarrow$ có mối quan hệ
- Nếu $p\text{-value} \geq 0.05 \rightarrow$ độc lập

Cặp biến	p-value	Kết luận
Thương_hiệu ↔ Dòng_xe	0	✓ Có mối quan hệ chặt chẽ
Thương_hiệu ↔ Loại_xe	~0	✓ Có mối quan hệ mạnh
Thương_hiệu ↔ Dung_tích_xe	0	✓ Có mối quan hệ rõ rệt
Thương_hiệu ↔ Tình_trạng	0	✓ Có mối quan hệ đáng kể
Thương_hiệu ↔ Xuất_xứ	0	✓ Có mối quan hệ rất mạnh
Dòng_xe ↔ Tình_trạng	1.0	⚠ Gần như độc lập
Loại_xe ↔ Tình_trạng	0.13	⚠ Gần như độc lập



EDA – PHÂN TÍCH DỮ LIỆU KHÁM PHÁ

Nguồn phân chia giá cao/thấp: 16,500,000 VND

Xe GIÁ CAO (từ khóa nổi bật)

Xe GIÁ THẤP (từ khóa nổi bật)

Nhóm	Từ khóa đặc trưng	Hàm ý nội dung
Giá cao	<i>chính chủ, sang tên, công chứng, liên hệ, xem xe</i>	Nhấn mạnh uy tín, pháp lý và chất lượng → dễ làm tăng giá trị cảm nhận
Giá thấp	<i>bán xe, tờ đày, zin, đày đủ, cần bán</i>	Nhấn mạnh tình trạng sử dụng và nhu cầu bán gấp → thể hiện xu hướng giảm giá



DATA INSIGHT

Dữ liệu sau làm sạch phản ánh thị trường xe cũ chính xác hơn, giảm nhiễu từ giá trị cực đoan.

Yếu tố ảnh hưởng giá mạnh: thương hiệu, dung tích, tình trạng xe; các yếu tố kỹ thuật khác (năm đăng ký, số km) tác động yếu.

Phân khúc cao cấp (Harley, BMW, Triumph, Touring...) có giá vượt trội, tạo ranh giới rõ rệt giữa xe phổ thông & xe sang.

Xe tay ga & phân khối lớn: biên độ giá rộng → phù hợp mô hình dự đoán & phát hiện bất thường.

Giá theo quận chủ yếu 30–50 triệu; vượt xa vùng này có thể là dấu hiệu bất thường hoặc thị trường riêng.

Hành vi người bán từ mô tả:

- Giá cao → mô tả chi tiết, đáng tin cậy
- Giá thấp → mô tả ngắn, nhấn mạnh bán nhanh

Ứng dụng: sử dụng các yếu tố này để huấn luyện mô hình dự đoán giá hợp lý và phát hiện tin bất thường (giá ảo, nhập sai, gian lận).

DATA PREPARATION

🧹 **Làm sạch dữ liệu:** loại bỏ giá trị null, bản ghi trùng lặp, ngoại lệ và các cột không cần thiết.

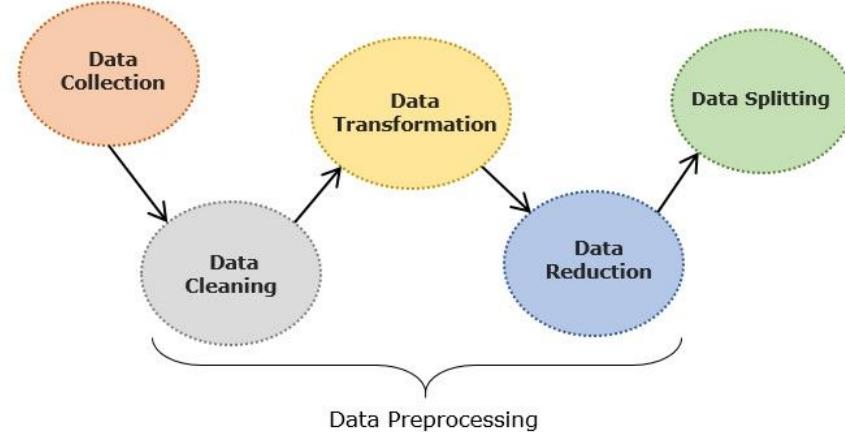
📝 **Xử lý dữ liệu văn bản:** chuẩn hóa, làm sạch text để sử dụng cho mô hình.

🔢 **Mã hóa biến phân loại:** áp dụng Label Encoding cho các cột categorical.

⚙️ **Tạo đặc trưng mới:** ví dụ tuổi xe, giá/km, quận để tăng hiệu quả mô hình.

📊 **Chia tập dữ liệu:** phân tách **train/test** để huấn luyện và đánh giá mô hình.

DATA PREPARATION





DATA PREPARATION

Quyết định giữ / bỏ cột

Quyết định	Cột	Lý do
<input checked="" type="checkbox"/> Giữ	Khoảng_giá_min	Là mốc giá thị trường, đại diện cho giá trung bình hợp lý nhất.
<input checked="" type="checkbox"/> Giữ	Tuổi_xe, log_Km, Km_trên_năm	Phản ánh hao mòn và mức độ sử dụng xe — tác động mạnh đến giá.
<input checked="" type="checkbox"/> Giữ	Dòng_xe_top, Loại_xe, Dung_tích_xe, Tình_trạng, Xuất_xứ, Quận, Phân_khúc, Thương_hiệu	Thể hiện phân khúc, đặc điểm kỹ thuật và vị trí xe — giúp mô hình học được khác biệt giữa các nhóm.
— Bỏ	Giá	Được dùng làm biến mục tiêu riêng (label_log).
<input checked="" type="checkbox"/> Bỏ	id, Href, Địa_chỉ, Tiêu_dè, Mô_tả_chi_tiết, tieu_de_clean, mo_ta_chi_tiet_clean, text_all_clean	text_all_clean không dùng trong PySpark vì TF-IDF làm phình dữ liệu, dễ gây lỗi bộ nhớ; tuy nhiên vẫn hữu ích cho mô hình ML truyền thống (RandomForest, XGBoost) nếu xử lý riêng.
<input checked="" type="checkbox"/> BỎ	Khoảng_giá_max, Thương_hiệu_top, Chính_sách_bảo_hành, Trọng_lượng	Tương quan cao hoặc ít giá trị phân biệt, gây nặng pipeline.



MODELING (Regression)

ML truyền thống

Mô hình	Ưu điểm nổi bật	Giải thích & Ứng dụng
Random Forest	Ôn định, ít cần tuning, kháng nhiễu tốt	Dễ triển khai, phù hợp cho bài toán tabular như dự đoán giá xe; cho baseline đáng tin cậy.
Gradient Boosting	Chính xác cao, xử lý tốt quan hệ phi tuyến	Xây chuỗi cây liên tiếp để tối ưu lỗi còn lại; phù hợp khi cần độ chính xác cao hơn RF.
CatBoost	Tối ưu cho dữ liệu phân loại, không cần mã hóa	Tự xử lý biến dạng category; tránh overfitting tốt, thường outperform GBT khi dữ liệu có nhiều biến rời rạc.
XGBoost	Nhanh, hiệu quả, nhiều tham số tối ưu	Rất phổ biến, dễ tinh chỉnh; phù hợp khi muốn tối đa hóa hiệu suất và độ chính xác.
LightGBM	Huấn luyện nhanh, tiết kiệm RAM, mạnh với dữ liệu lớn	Dùng khi cần tốc độ cao, đặc biệt với dữ liệu nhiều chiều hoặc tập lớn.
Ridge Regression	Đơn giản, dễ giải thích, kiểm soát overfitting	Mô hình tuyến tính cơ bản, dùng làm baseline hoặc phân tích ảnh hưởng của từng biến đầu vào.



MODELING (Regression)

Pyspark

Mô hình	Ưu điểm	Giải thích & Ứng dụng
 Random Forest	Ôn định, dễ song song hóa, chống overfitting	Thích hợp cho dữ liệu lớn và phức tạp; mô hình tạo nhiều cây quyết định độc lập nên chạy tốt trong Spark cluster, cho kết quả ổn định và ít cần tinh chỉnh.
 Gradient Boosting (GBTRegressor)	Độ chính xác cao, tối ưu hóa mạnh trong Spark MLlib	Kết hợp dãy cây liên tiếp để giảm sai số còn lại; phù hợp khi cần dự đoán chính xác giá trị liên tục như giá xe ; hiệu quả hơn RF nếu dữ liệu đã được xử lý kỹ.
 Ridge Regression	Nhanh, nhẹ, dễ giải thích hệ số	Là mô hình tuyến tính có regularization giúp giảm nhiễu; phù hợp làm baseline, kiểm tra quan hệ tuyến tính giữa biến độc lập và giá bán, hoặc chạy thử trên tập lớn để benchmark.

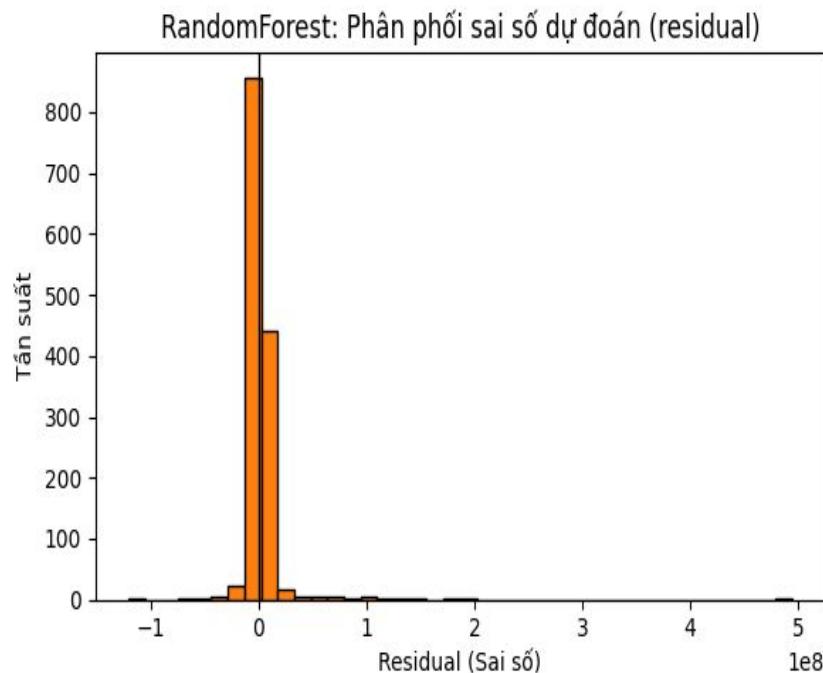
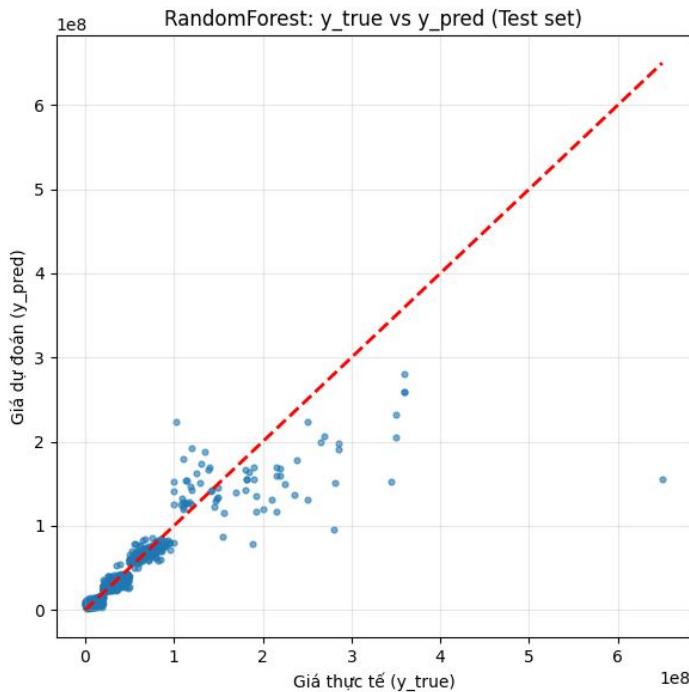
→ KẾT QUẢ ĐÁNH GIÁ: Dựa trên giá thật (không biến đổi log)

Pyspark

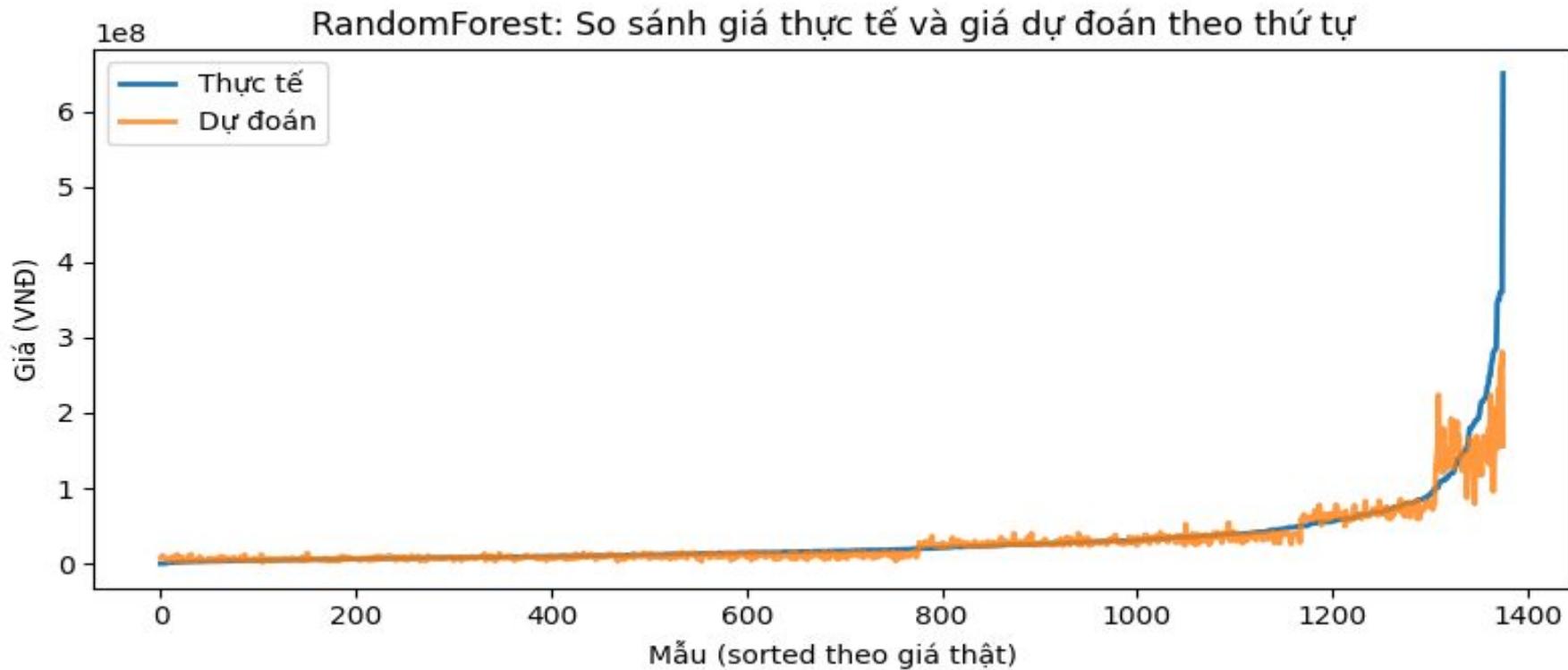
Mô hình	R ² (Test)	MAE (Test)	RMSE (Test)	Thời gian (s)
 Random Forest	0.799	6,525,991	20,885,617	107.35
 Ridge Regression	0.533	12,584,408	31,862,088	14.07
 Gradient Boosting	0.412	8,112,962	35,745,903	114.18

→ KẾT QUẢ ĐÁNH GIÁ:

Pyspark



So sánh giá thực tế và giá dự đoán theo thứ tự (trên pyspark)



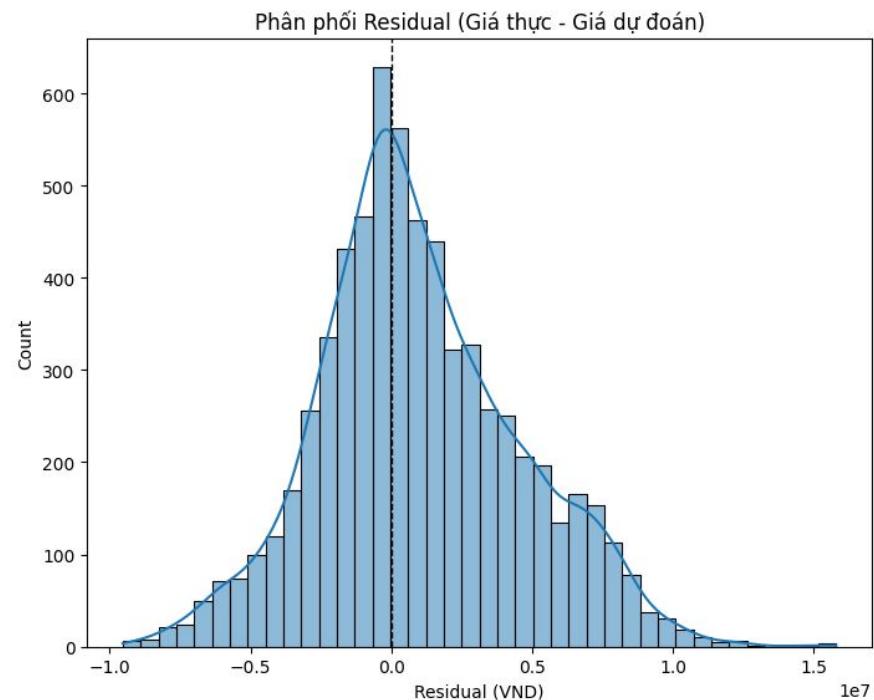
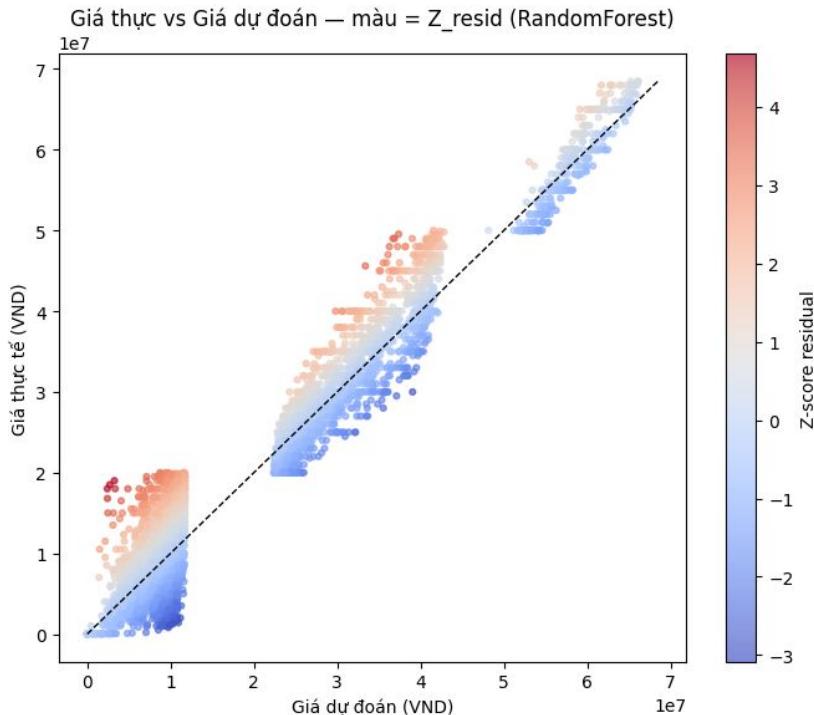
→ KẾT QUẢ ĐÁNH GIÁ: Dựa trên giá thật (không biến đổi log)

ML truyền thống

Mô hình	R ² (Test)	MAE (Test)	RMSE (Test)	Thời gian (s)
 Random Forest	0.890	3,687,591	23,934,808	293.2
 Gradient Boosting	0.875	3,682,696	27,208,377	243.2
 CatBoost	0.863	4,002,870	29,976,868	396.5
 XGBoost	0.858	3,832,443	30,924,787	95.8
 LightGBM	0.285	4,939,162	156,238,045	45.1
 Ridge Regression	0.260	7,806,546	161,607,611	2.5

→ KẾT QUẢ ĐÁNH GIÁ:

ML truyền thống





So sánh hiệu năng giữa hai môi trường

Tiêu chí	ML truyền thống	PySpark
Độ chính xác (R^2)	Cao: ~0.86–0.89 nhờ tinh chỉnh & boosting mạnh	Ôn định: ~0.8, phù hợp dữ liệu lớn
Mô hình hiệu quả nhất	Random Forest: cân bằng độ chính xác & thời gian huấn luyện	Random Forest: ổn định trên dữ liệu phân tán
Thời gian & hiệu suất	Nhanh trên tập nhỏ, nhưng tốn bộ nhớ khi mở rộng	Tối ưu cho dữ liệu lớn, hỗ trợ phân tán & song song hóa
Các mô hình khác	LightGBM, XGBoost: nhanh, chính xác, phù hợp local/batch nhỏ	Hai mô hình còn lại thì chưa phù hợp với yêu cầu

♦ Nhận xét tổng quan

- ML truyền thống cho độ chính xác cao hơn trên dữ liệu nhỏ và khả năng tinh chỉnh mạnh.
- PySpark ổn định trên dữ liệu lớn, tối ưu hiệu suất phân tán.
- Random Forest là lựa chọn cân bằng giữa chính xác và thời gian huấn luyện trong cả hai môi trường.
- Các mô hình boosting (LightGBM, XGBoost) linh hoạt, nhanh và chính xác → phù hợp cho chạy local hoặc batch nhỏ.

Phát hiện bất thường (Anomaly Detection)

→ Phương pháp:
Isolation Forest, LOF, One-Class SVM.

Mô hình	AUC (weak)	AP (weak)	Thời gian (s)
LOF (Local Outlier Factor)	0.742	0.746	0.62
Isolation Forest	0.713	0.726	1.45
One-Class SVM	0.543	0.583	0.25



Quy luật phát hiện giá bất thường

Sử dụng ML truyền thống kết hợp với phương pháp xử lý bất thường

Nguyên tắc chung

- So sánh Giá thực với Giá dự đoán (Random Forest).
- Kiểm tra tin đăng khác biệt so với thị trường (LOF – Local Outlier Factor).
- Kết hợp giá và mẫu tin để đánh giá bất thường.

Cách tính điểm

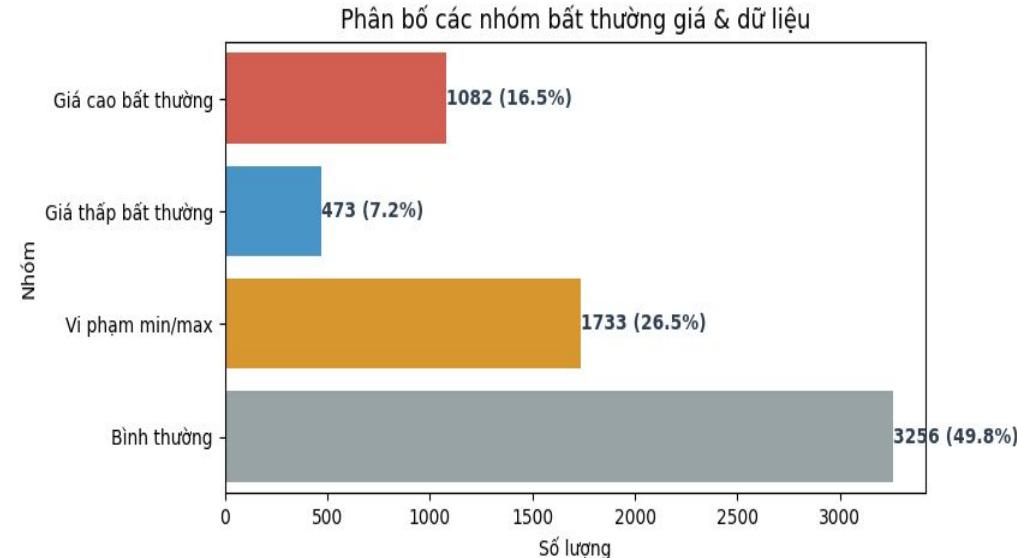
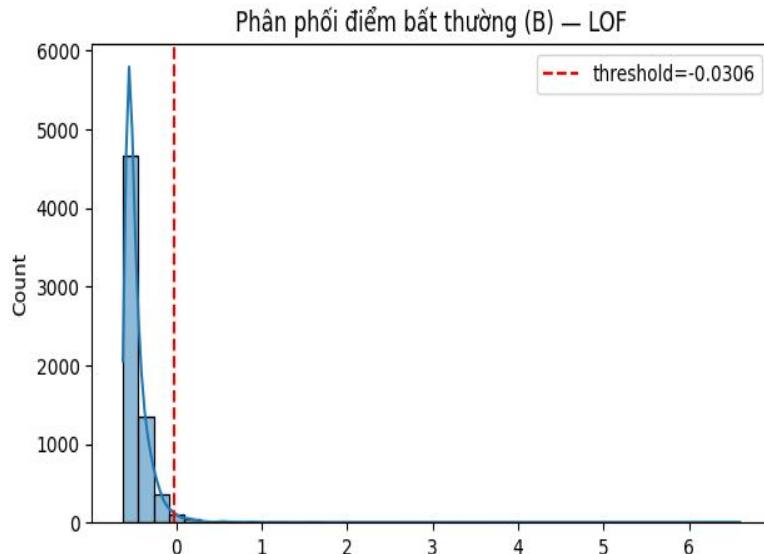
Điểm	Mô tả
A-score (Z-score)	Mức độ lệch giữa Giá thực ↔ Giá dự đoán
B-score (LOF)	Mức độ khác biệt của tin đăng so với dữ liệu chung
Abnormal_score	$0.3 \times A + 0.7 \times B \rightarrow 0-100$ (càng cao → càng bất thường)

Khi nào xem là bất thường?

- **Giá cao bất thường:** $Z \geq +3$ hoặc Giá thực $>$ Giá dự đoán $+40\%$
- **Giá thấp bất thường:** $Z \leq -3$ hoặc Giá thực $<$ Giá dự đoán -40%
- **Vị phạm min/max:** Giá nằm ngoài khung min–max $\pm 15\%$
- **LOF B_flag = 1** → tin có đặc điểm lạ, ưu tiên kiểm tra

Evaluation & Comparison

- Regression: Random Forest
- Anomaly Detection: Isolation Forest ổn định nhất.
- Các chỉ số đánh giá tổng hợp (RMSE, MAE, Recall anomaly).



DEPLOYMENT

<https://dudoangiaxevaphathienbatthuong-vha-nmxb.streamlit.app/>

Đồ án tốt nghiệp
Data Science

Thời gian: 29/11/2025

Người thực hiện

- Võ Thị Hoàng Anh
anvo.bio@gmail.com
- Nguyễn Mai Xuân Bách
bachxdn@gmail.com

Giảng viên hướng dẫn

- Cô Khuất Thùy Phương

Điều hướng

Chọn mục điều hướng:

Dataset Input

Đăng nhập Admin

Username (admin)

Password

Đăng nhập

Dự đoán giá & Phát hiện bất thường giá xe máy



chợ TỐT

CHỢ TỐT XE MÁY - VUI LÊN LÀ CÓ XE NGON!

Săn xe độc - Giá hời - Chốt đơn cái vèo!

HÀNG NGÀN XE CHỜ BẠN

Dự đoán giá + Phát hiện bất thường

Giới thiệu & Quy trình

DEPLOYMENT

Trang Giới thiệu, Quy trình

Đồ án tốt nghiệp
Data Science

Thời gian: 29/11/2025

Người thực hiện

- Dataset Input
- Business + Data Un...
- EDA Numeric
- EDA Categorical
- WordCloud
- Modeling
- Anomaly Detection
- Evaluation
- Dataset Input

Đọc dữ liệu ban đầu (Dataset Input)

5 dòng đầu tiên của dữ liệu:

	ID	Tiêu đề	Giá	Khoảng giá min	Khoảng giá max	Địa chỉ	Mô tả chi tiết
0	1	Bán Vespa Sprint 125cc 2024 xanh dương, xe đẹp 95%	66.000.000 đ	72.53 tr	85.14 tr	Phường Bến Thành, Quận 1, Tp Hồ Chí Minh	Bán xe #Vespa S
1	2	SH 150i Thắng ABS 2019 BSTP Chính Chủ	79.500.000 đ	62.76 tr	73.68 tr	Phường Tân Định, Quận 1, Tp Hồ Chí Minh	Bán SH 150i Th
2	3	CC Vision Thể Thao 2023 Đen+bộ đèn Demi audi A7	37.000.000 đ	28 tr	32.86 tr	Phường Cầu Kho, Quận 1, Tp Hồ Chí Minh	Chính chủ bán VI
3	4	Vespa Sprint 2019 -125- Đen Đỏ Sport -CHÍNH CHỦ	45.000.000 đ	43.1 tr	50.6 tr	Phường Bến Nghé, Quận 1, Tp Hồ Chí Minh	XE CÁ NHÂN BÁN
4	5	Xe tay ga Yamaha Latte 125 - Đăng ký 2021	23.000.000 đ	17.02 tr	19.98 tr	Phường Tân Định, Quận 1, Tp Hồ Chí Minh	Thông tin xe:

Đánh giá mô hình

Biểu đồ đánh giá

Phân bố các nhóm bất thường giá & dữ liệu

Nhóm	Số lượng	Tỷ lệ (%)
Giá cao bất thường	1082	16.5%
Giá thấp bất thường	473	7.2%
Vì phạm min/max	1733	26.5%
Bình thường	3256	49.8%

Nhận xét

Kết quả:

- Bình thường: chiếm đa số** → Phần lớn dữ liệu có mức giá hợp lý, cho thấy hệ thống đánh giá hoạt động ổn định.
- Vì phạm min/max: nhóm lớn thứ hai** → Giá rao nằm ngoài khoảng giá tham chiếu (cao hơn hoặc thấp hơn khung hợp lý). Nhóm này không hẳn sai, nhưng là vùng rủi ro cần được xem xét kỹ khi kiểm duyệt (xe độ, xe hiếm, xe bán gấp...).
- Giá bất thường: chiếm tỷ lệ nhỏ** → Những tin đăng có mức giá cao hoặc thấp khác thường, thường liên quan tới nâng giá, nhập sai, hoặc mô tả bất thường.

Ứng dụng:

- Gợi ý mức giá hợp lý cho người bán.

DEPLOYMENT

Về phía người đăng tin

Dự đoán giá + Phát hiện bất thường Giới thiệu & Quy trình

Thực hiện dự đoán giá & kiểm tra bất thường

Chọn cách nhập dữ liệu:

Nhập tay từng xe Tải file CSV/XLSX

Thương hiệu	Năm đăng ký	Khoảng giá_min (VND) – có thể bỏ trống
Honda	2020	0
Dòng xe	Số km đã đi	Khoảng giá_max (VND) – có thể bỏ trống
SH	20000	0
Loại xe	Giá thực (VND) – dùng để đánh giá bất thường	
Tay ga	5000000	
Dung tích xe		
100 - 175 cc		

Kết quả dự đoán

Giá dự đoán: **61,427,873 VND**

Đánh giá bất thường về giá

⚠ Kết luận: Giá thấp bất thường

> Xem lý do

Bảng chi tiết

	Thương_hiệu	Dòng_xe	Loại_xe	Dung_tích_xe	Quận	Khoảng_giá_min	Khoảng_giá_max	Năm_đăng_ký	Tuổi_xe	Số_Km_đã_di	Giá	Giá_dự đoán	Lý do
0	honda	sh	tay ga	100 - 175 cc	Quận 1	56,355,000	89,504,500	2020	5	20000	50,000,000	61,427,873	(*)

Thực hiện dự đoán giá & kiểm tra bất thường

Chọn cách nhập dữ liệu:

Nhập tay từng xe Tải file CSV/XLSX

Chọn file dữ liệu:

Drag and drop file here
Limit 200MB per file • CSV, XLSX

demo.xlsx 10.8KB X

Dự đoán giá cho file Phát hiện bất thường cho file

Có 4 dòng bất thường.

	id	Thương_hiệu	Dòng_xe	Loại_xe	Dung_tích_xe	Quận	Giá	Giá_dự đoán	Kết_luận_cuối	Lý do_ngắn_gọn
0	0	piaggio	vespa	tay ga	100 - 175 cc	Quận 1	66,000,000	37,729,638	Giá cao bất thường	cao hơn dự đoán, đặc điểm khác biệt
4	4	yamaha	latte	tay ga	100 - 175 cc	Quận 1	23,000,000	10,986,527	Giá cao bất thường	cao hơn dự đoán, ngoài min/max, đặc điểm khác biệt
6	6	honda	air blade	tay ga	100 - 175 cc	Quận 1	16,000,000	10,789,071	Giá cao bất thường	cao hơn dự đoán, đặc điểm khác biệt
7	7	honda	wave	xe số	50 - 100 cc	Quận 1	7,300,000	9,460,457	Vi phạm min/max	thấp hơn dự đoán, ngoài min/max, đặc điểm khác biệt

Các tin còn lại là BÌNH THƯỜNG

	id	Thương_hiệu	Dòng_xe	Loại_xe	Dung_tích_xe	Quận	Giá	Giá_dự đoán	Kết_luận_cuối
1	1	honda	sh	tay ga	100 - 175 cc	Quận 1	79,500,000	62,037,959	Bình thường
2	2	honda	vision	tay ga	100 - 175 cc	Quận 1	37,000,000	34,028,742	Bình thường
3	3	piaggio	vespa	tay ga	100 - 175 cc	Quận 1	45,000,000	41,794,834	Bình thường
5	5	sym	elegant	xe số	dưới 50 cc	Quận 1	7,500,000	8,491,316	Bình thường

DEPLOYMENT

Về phía admin (người quản lý)

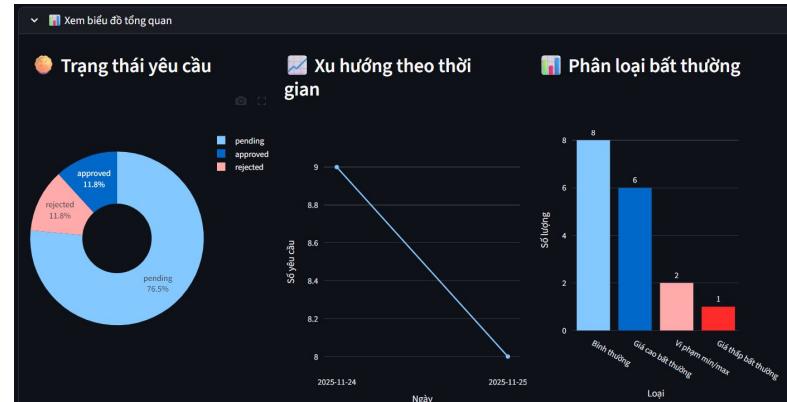
Danh sách yêu cầu			
Tin BẤT THƯỜNG (Pending) ↗			
	id_yêu_cầu	thời_gian	kết_quả_mô_hình
0	9c8811dc-7672-4cb1-8a69-89359873c7b7	2025-11-24 14:47:30	Giá_dự_đoán: 37,729,638 Kết_luận: Giá cao bất thường Lý_do: - Giá thực cao hơn giá dự đoán k
4	cd455344-8c4c-4571-8a59-969a8b11d8bf	2025-11-24 14:47:33	Giá_dự_đoán: 10,986,527 Kết_luận: Giá cao bất thường Lý_do: - Giá thực cao hơn giá dự đoán k
6	82be3e1e-2545-46a6-8ae9-de42a24dd9ea	2025-11-24 14:47:34	Giá_dự_đoán: 10,789,071 Kết_luận: Giá cao bất thường Lý_do: - Giá thực cao hơn giá dự đoán k
7	74076d05-9156-487a-bb07-25c69d8b9f87	2025-11-24 14:47:35	Giá_dự_đoán: 9,460,457 Kết_luận: Vì phạm min/max Lý_do: - Giá thực thấp hơn giá dự đoán k
9	d09e6dc3-6de8-451b-8601-049e1e44ba7e	2025-11-25 13:11:16	Giá_dự_đoán: 37,729,638 Kết_luận: Giá cao bất thường Lý_do: - Giá thực cao hơn giá dự đoán k
13	f83f3251-f72f-48a4-9fd4-3dbd87d71813	2025-11-25 13:11:20	Giá_dự_đoán: 10,986,527 Kết_luận: Giá cao bất thường Lý_do: - Giá thực cao hơn giá dự đoán k
15	d1d18815-0a08-4c83-8385-75e87bb1109c	2025-11-25 13:11:22	Giá_dự_đoán: 10,789,071 Kết_luận: Giá cao bất thường Lý_do: - Giá thực cao hơn giá dự đoán k

Trung tâm quản trị

Có 13 yêu cầu mới đang chờ xử lý!

Thống kê ↗

Tổng yêu cầu	Đang chờ	Đã xử lý
17	13	2 ✓ / 2 X





KẾT QUẢ & INSIGHT

🏁 Kết quả chính

- Xây dựng thành công hệ thống dự đoán giá xe máy cũ và phát hiện bất thường từ dữ liệu thực.
- RandomForest cho kết quả dự đoán ổn định và phù hợp với bài toán định giá.
- LOF phát hiện hiệu quả các tin có giá lệch chuẩn hoặc đặc điểm khác thường trong mô tả.
- Hoàn thiện web-app trực quan, giúp người dùng tự định giá và hỗ trợ đội kiểm duyệt nhanh hơn.

🚀 Hướng cải thiện

- Nâng cấp mô hình xử lý văn bản: TF-IDF → BERT để hiểu mô tả chính xác hơn.
- Xây dựng mô hình min/max động, thay vì sử dụng giới hạn cố định $\pm 15\%$.
- Phát triển API để tích hợp vào hệ thống thực tế ở quy mô lớn.
- Tối ưu UI/UX và cải thiện tốc độ xử lý khi tải file dung lượng lớn.

NHÓM THỰC HIỆN



VÕ THỊ HOÀNG ANH

🎓 ThS. Công nghệ Sinh học
✉ Email: anhvo.bio@gmail.com
📱 SĐT: 0935 719 426
🐙 GitHub: <https://github.com/anhvobio>
💼 Nội dung truyền thông

🔧 **Nhiệm vụ đã thực hiện:** Xây dựng mô hình dự đoán giá; soạn bài thuyết trình và viết README.md hướng dẫn sử dụng ứng dụng cho người dùng; cùng phát triển ứng dụng Streamlit.

🎯 **Kết quả & Điều học được:** Qua quá trình làm project, Hoàng Anh hiểu rõ hơn từng bước, phương pháp và ý nghĩa trong toàn bộ quy trình Data Science. Đồng thời biết cách thực hiện trọn vẹn các giai đoạn từ xử lý tiền dữ liệu, xây dựng mô hình đến triển khai giao diện người dùng bằng GUI Streamlit.



NGUYỄN MAI XUÂN BÁCH

🎓 CN. Kỹ Thuật Điều Khiển Và Tự Động Hóa
✉ Email: bachxdn@gmail.com
📱 SĐT: 0982 308 974
🐙 GitHub: <https://github.com/BachNguyen2kx>
💼 Hỗ trợ kỹ thuật

🔧 **Nhiệm vụ đã thực hiện:** Khám phá & xử lý dữ liệu, phát hiện bất thường; kiểm tra code và viết README.md cho final project; cùng tham gia xây dựng ứng dụng Streamlit.

🎯 **Kết quả & Điều học được:** Bách nhận ra tầm quan trọng của việc hiểu sâu dữ liệu và lựa chọn mô hình phù hợp — chỉ cần sai vài cột hoặc chọn sai mô hình là toàn bộ kết quả có thể lệch đi.

Việc triển khai ứng dụng lên Streamlit Cloud và Google Sheets mang lại kinh nghiệm thực tế về xây dựng giải pháp end-to-end, không chỉ dừng lại ở việc xây dựng mô hình.

Bên cạnh đó, Bách học được cách quản lý trạng thái, xử lý lỗi, thiết kế UI/UX và viết code sao cho người dùng không rành kỹ thuật vẫn có thể thao tác dễ dàng.