

Báo cáo kỹ thuật

Dự án: Xây dựng mô hình hồi quy tuyến tính dự đoán giá xe ô tô cũ

Giới thiệu bài toán:

Dự án này sử dụng bộ dữ liệu Car Price Prediction (Kaggle) để xây dựng các mô hình học máy dự đoán giá xe ô tô cũ.

Dữ liệu bao gồm hãng xe, năm xe được bán ra và một số thông tin của chiếc xe và người bán.

Mục tiêu của dự án là xây dựng một số mô hình hồi quy hỗ trợ người mua xe cũ tốt nhất dựa trên những đặc trưng của chiếc xe được bán, giúp người mua và người bán có thể định giá chiếc xe

Phân tích và khám phá dữ liệu:

Bộ dữ liệu Car Price Prediction bao gồm 8128 dòng và 13 cột mô tả các thông tin về chiếc xe như là hãng xe, năm sản xuất, loại nhiên liệu sử dụng, người bán, loại hộp số (tự động hoặc thủ công),... và giá bán

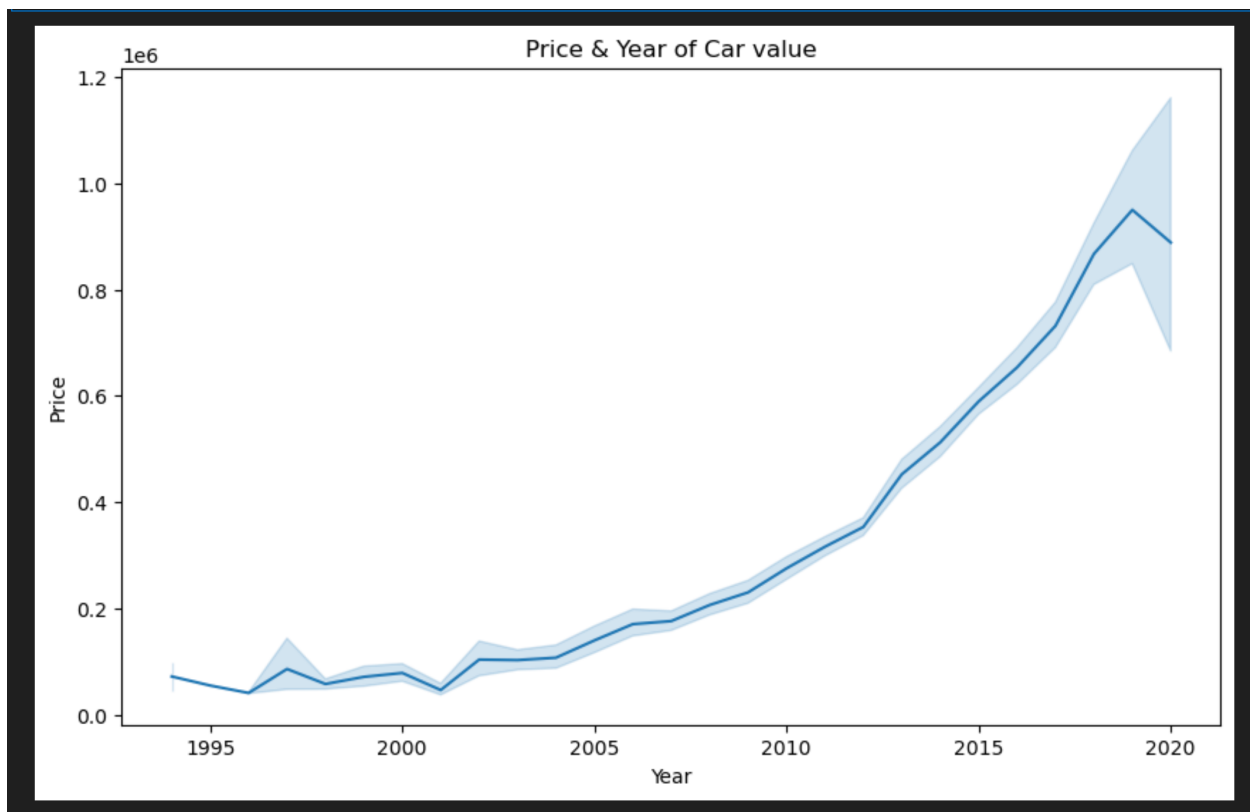
Trong giai đoạn phân tích và khám phá dữ liệu, dữ liệu đã được làm sạch bằng việc kiểm tra các giá trị bị trùng hoặc thiếu để đảm bảo chất lượng dữ liệu trước khi tiếp tục các bước phân tích khác.

Sau khi tiến hành xử lý, dữ liệu còn lại 6926 dòng và 12 cột

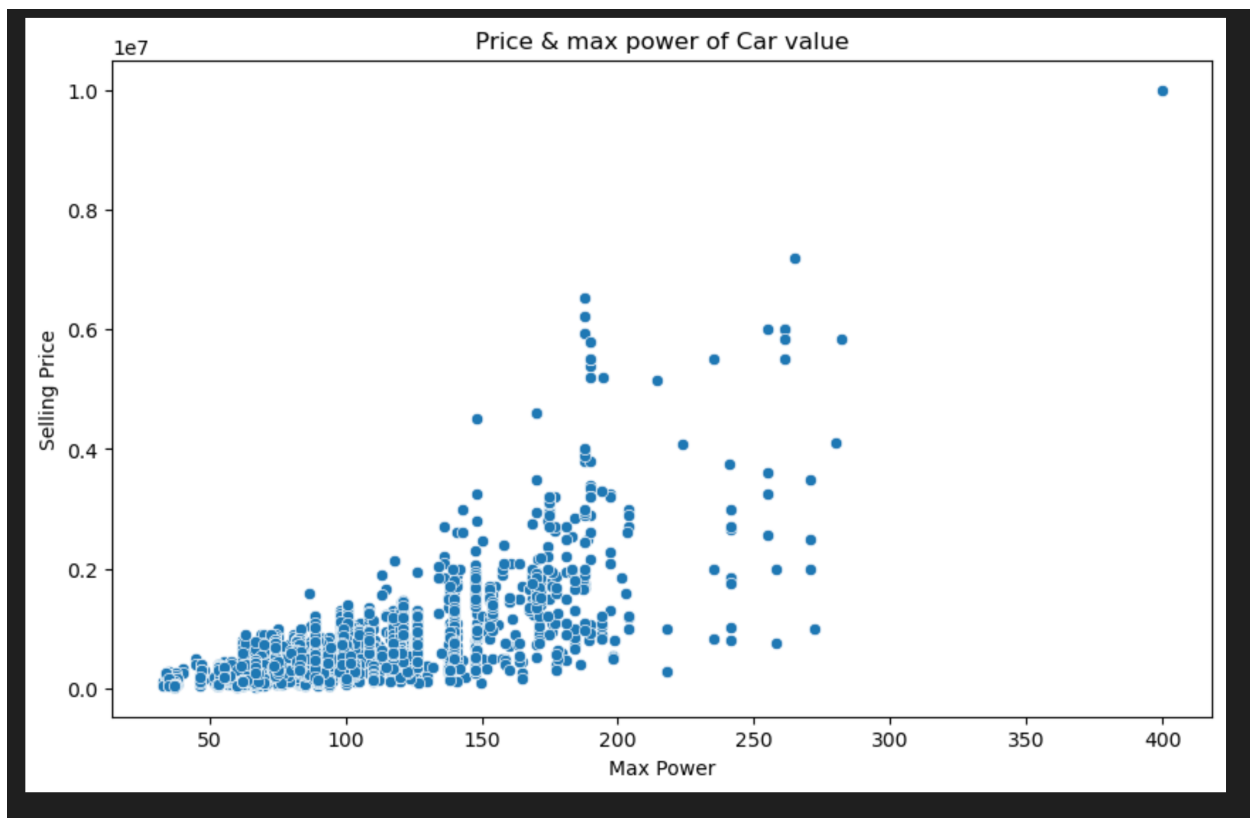
Bên cạnh đó, việc kiểm tra các cột như mileage, engine, max power,... giúp xác định các thang đo khác nhau trong bộ dữ liệu có thống nhất hay không để biến đổi các đơn vị đo này thành biểu diễn dạng số để có thể dễ dàng để mô phỏng.

Sau khi phân tích và biến đổi dữ liệu, việc mô phỏng bằng một số biểu đồ để kiểm tra sự liên kết giữa giá bán và một số thông tin khác của xe đã tối ưu hơn

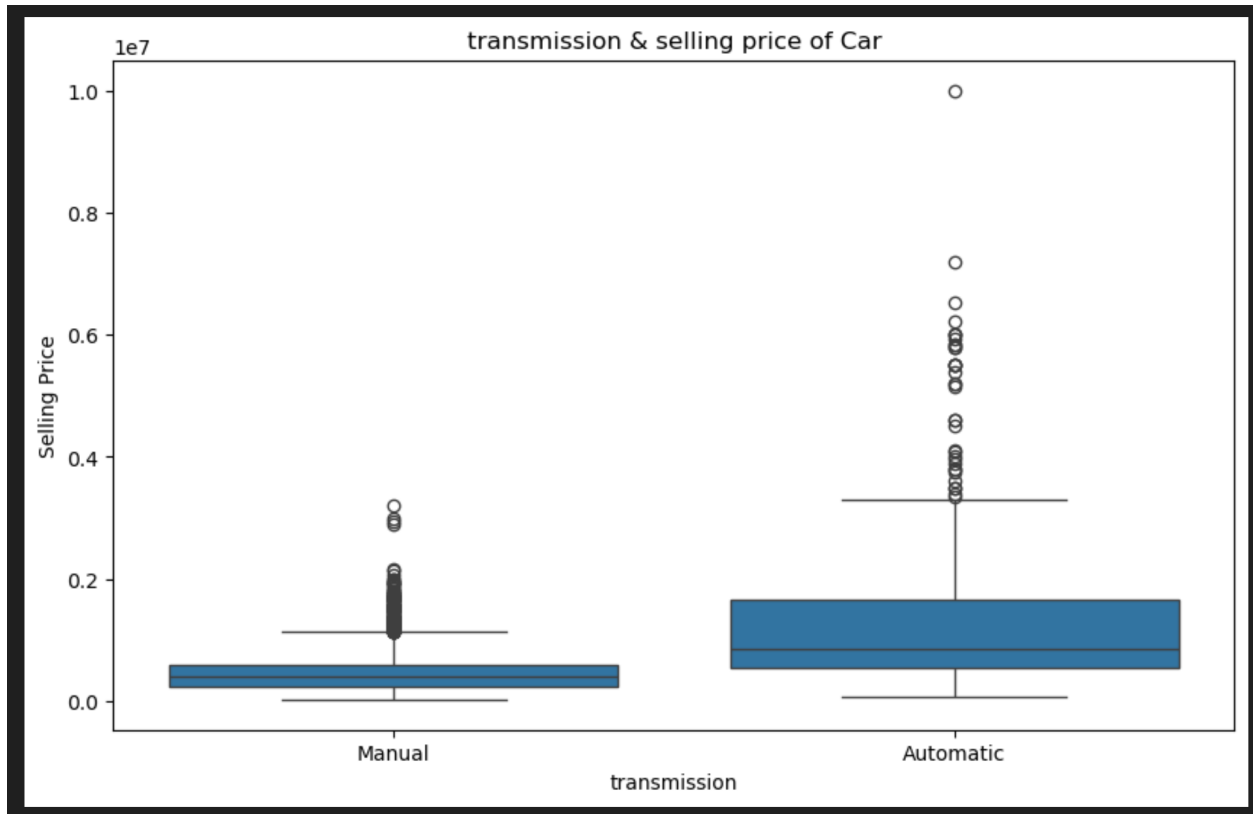
Như: năm sản xuất và giá bán xe



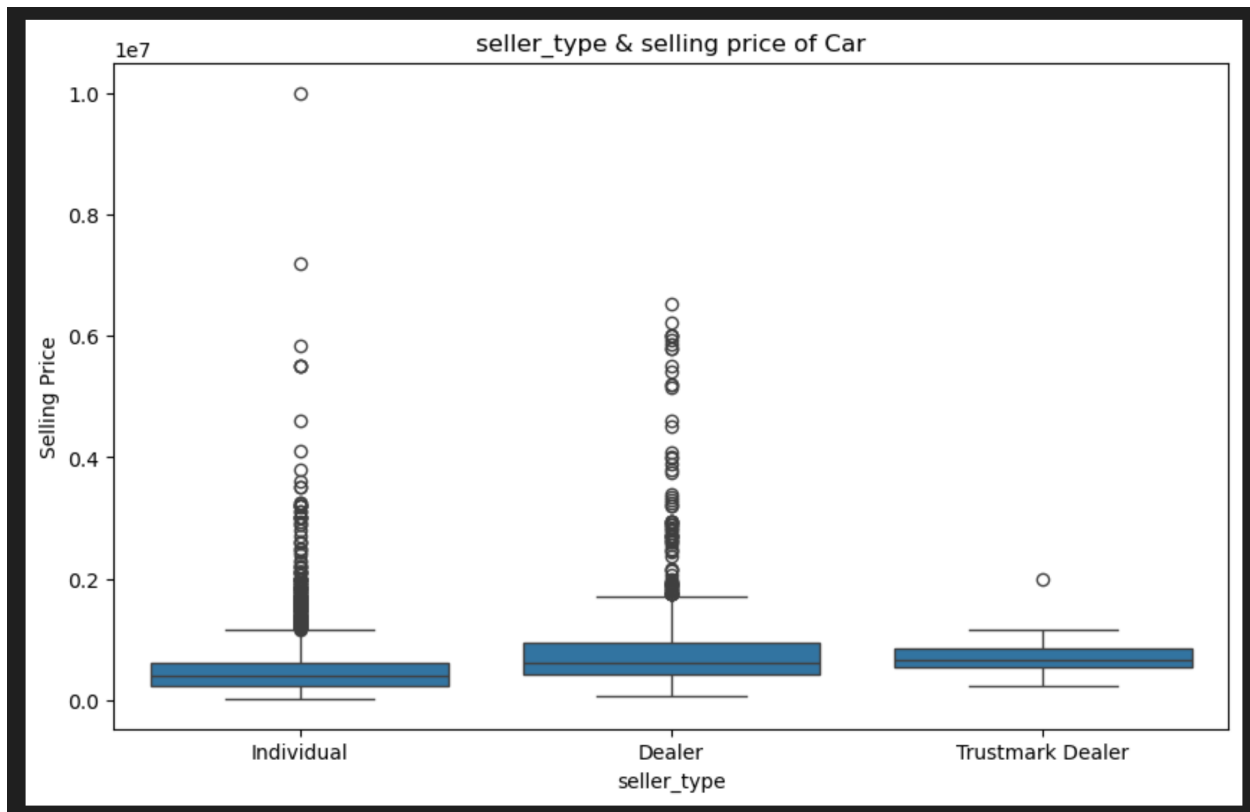
Hiệu năng và giá bán xe:



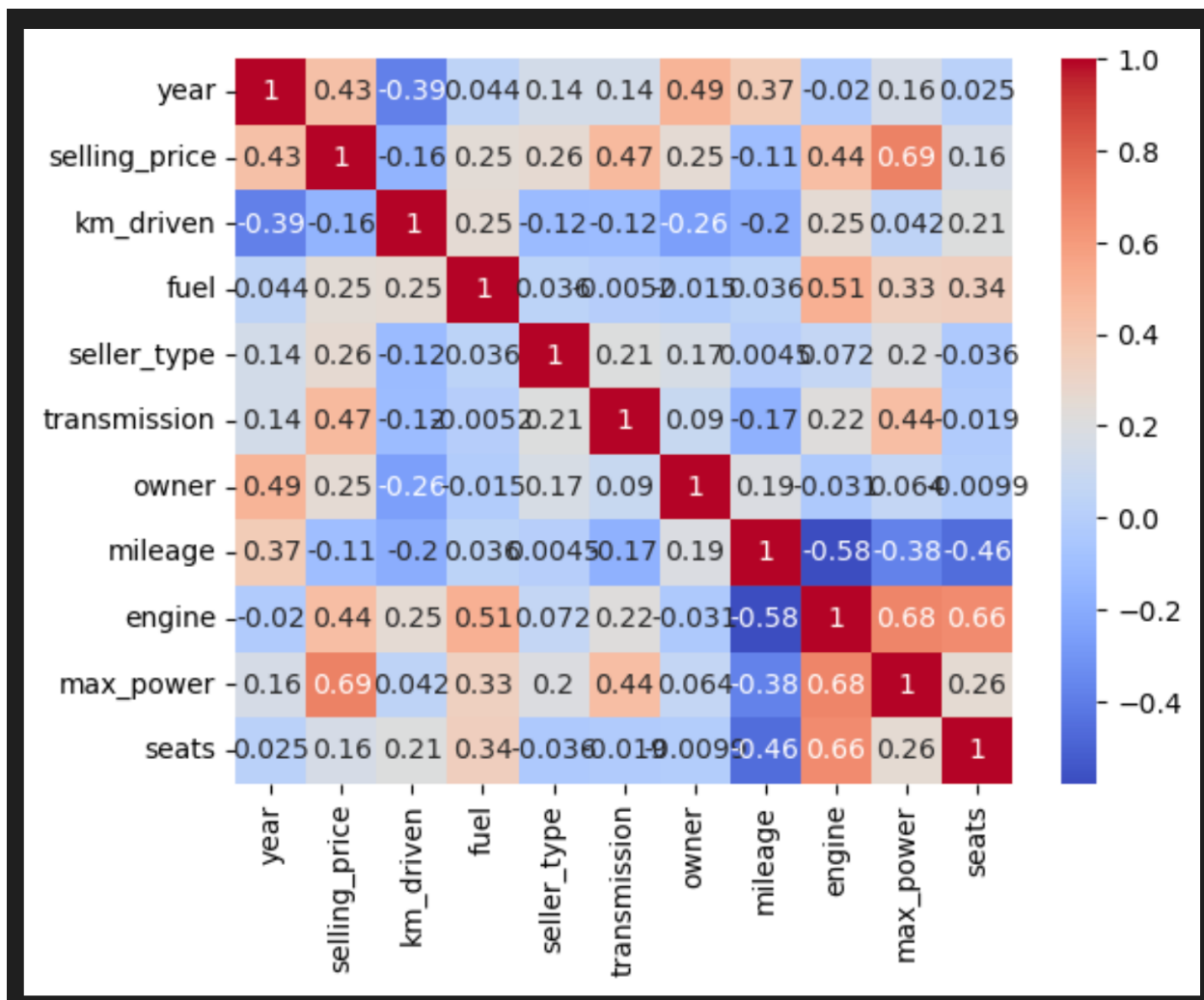
loại hộp số và giá xe:



Loại người bán và giá tiền của họ đề ra:



Việc xây dựng biểu đồ nhiệt dựa trên những quan sát này giúp chúng ta thấy được quan hệ giữa các đặc trưng có thể ảnh hưởng đến giá bán xe như thế nào



Vậy ta rút ra được hiệu năng của xe (0.69), loại hộp số (0.47), loại động cơ (0.44), năm sản xuất (0.43) liên quan mật thiết đến giá bán của xe và một số lượng các đặc trưng khác cũng giúp cho việc dự đoán giá xe chính xác hơn

Pipeline:

Tiền xử lý và chuyển đổi đặc trưng:

Tiền xử lý:

Bỏ đi cột torque do đặc trưng không nhất quán và có thể gây nhiễu

Kiểm tra số lượng dòng bị trùng và loại bỏ chúng

Kiểm tra số lượng dòng mà một số đặc trưng quyết định bị thiếu, vì chỉ chiếm 3.1% tổng số dữ liệu, một số lượng khá nhỏ, việc loại bỏ những dòng có giá trị thiếu sẽ giúp việc xây dựng mô hình dễ dàng hơn

Sau đó, kiểm tra thang đo của các đơn vị trong dữ liệu gốc và chuyển đổi các đơn vị đặc trưng thành dạng số thực để dễ tính toán

Sau đó, chuyển đổi tên xe thành những tên hãng xe để dễ dàng phân loại

Biến đổi đặc trưng:

Áp dụng One-hot encoding cho đặc trưng hãng xe

Áp dụng Ordinal Encoding để biến đổi các đặc trưng người bán, loại hộp số, loại người sở hữu, và loại xăng thành các số 0,1,...

Áp dụng chuẩn hoá RobustScaler để chuẩn hoá cột max_power và engine dựa trên training data và áp dụng vào testing data, biến đổi các năm thành số năm tuổi của xe với 2020 là năm lớn nhất

Áp dụng log transformation lên giá bán của xe để co lại khoảng giá tiền outlier trong dữ liệu gốc để mô hình ổn định hơn

Lựa chọn đặc trưng:

Dựa trên heatmap, các đặc trưng như hiệu năng của xe, loại hộp số, người bán, năm sản xuất, loại nhiên liệu, người sở hữu, và động cơ ảnh hưởng lớn đến giá tiền của xe

Vậy nên các đặc trưng được chọn là: max_power, transmission, seller_type, owner, year, fuel, engine sẽ được đưa vào mô hình để học và dự đoán

Huấn luyện mô hình:

Phân chia tập dữ liệu với tỉ lệ 80/20 (với 80% là training data, 20% là test data)

Sử dụng 3 mô hình là Linear Regression, Random Forest, và Neural Network để tiến hành dự đoán và lựa chọn mô hình tốt nhất

Linear Regression: Mô hình được xây dựng áp dụng thư viện scikit-learn để làm baseline

Random Forest:

Cài đặt: số lượng cây quyết định: 100, trạng thái ngẫu nhiên: 42 (đảm bảo nhất quán trong quá trình huấn luyện)

Neural Network:

Cài đặt: (cài đặt bằng PyTorch)

Biến đổi các đặc trưng như X_{train} , X_{test} , y_{train} , y_{test} thành dạng Tensor để PyTorch có thể xử lý

Mạng nơ-ron bao gồm 4 lớp: 1 lớp nhận dữ liệu, 2 lớp ẩn, và 1 lớp output. Tất cả được áp dụng Batch Normalization để giúp quá trình huấn luyện ổn định hơn và tránh vấn đề explosive gradient

Với các lớp là hoàn toàn được kết nối với nhau được cài đặt như sau:

Input: Nhận đầu vào với kích thước là cột của dữ liệu, đầu ra có kích thước 256

Lớp ẩn thứ 1: Nhận đầu vào với kích thước 256, đầu ra có kích thước 128

Lớp ẩn thứ 2: Nhận đầu vào với kích thước 128, đầu ra có kích thước 64

Output: Nhận đầu vào với kích thước 64, đầu ra có kích thước 1 (là kết quả cuối cùng)

Mạng nơ-ron được huấn luyện với 500 epoch, sử dụng loại mất mát MSE, loại kích hoạt ReLU (để giúp mô hình học được các dạng dữ liệu phức tạp hơn), và sử dụng tối ưu hoá Adam (adaptive learning rate) với lr bắt đầu ở 0.01

Các mô hình sẽ được so sánh dựa trên kết quả huấn luyện để được chọn serving

Kết quả của quá trình huấn luyện: (dựa trên tập kiểm tra)

(Trên dữ liệu thực tế, đơn vị tiền tệ được tính theo đồng rupee - Ấn Độ)

Chỉ số	R^2	RMSE	MAE
Linear Regression	0.822	161219	93408
Random Forest	0.92	132284	77122
Neural Network	0.91	144320	84070

Việc kiểm tra sử dụng cross-validation (gồm 5 lần chia) để đảm bảo đầu ra đã chuẩn: (với giá trị trung bình trong 5 lần)

Chỉ số	R^2	RMSE	MAE
Linear Regression	0.83	213013	99488
Random Forest	0.9	163648	80157
Neural Network	0.88	179855	84324

(RMSE được tính dựa trên dữ liệu khi đã được chuẩn hoá)

Vậy mô hình Random Forest cho ra kết quả tốt nhất, vậy nên random forest được dùng để triển khai.

Cách chạy ứng dụng:

Bước 1: Cài đặt những thư viện cần dùng để chạy:

pip install -r requirements.txt

```
conda activate base
(base) lebinhminh@Shions-Macbook-Air submission % pip install -r requirements.txt
Requirement already satisfied: fastapi in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 1)) (0.128.0)
Requirement already satisfied: uvicorn in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 2)) (0.40.0)
Requirement already satisfied: pydantic in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 3)) (2.10.3)
Requirement already satisfied: pyngrok in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 4)) (7.5.0)
Requirement already satisfied: streamlit in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 5)) (1.45.1)
Requirement already satisfied: requests in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 6)) (2.32.3)
Requirement already satisfied: pandas in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 7)) (2.2.3)
Requirement already satisfied: numpy in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 8)) (2.1.3)
Requirement already satisfied: torch in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 9)) (2.9.1)
Requirement already satisfied: scikit-learn in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 10)) (1.6.1)
Requirement already satisfied: joblib in /opt/anaconda3/lib/python3.13/site-packages (from -r requirements.txt (line 11)) (1.4.2)
```

Bước 2: Chạy file huấn luyện (để lưu mô hình mới): (không bắt buộc)

python -m src.train

```
10:50:57 lebinhminh@Shions-Macbook-Air submission % python -m src.train
metrics for linear regression, with RMSE, MSE, and MAE for train_and_test
0.8573286843195098 0.8815493270822563
202578.40779114678 163219.48683130462
97948.39404749655 93408.79276977726
[[0.8572197192194877, 0.8580799119674413], [0.8673681916553295, 0.7761498622041125], [0.8552565407616584, 0.8755460961219707], [0.866429106450189, 0.832437684164
9241], [0.8604388396518055, 0.8191014039004321]]
[[98432.24332171642, 95364.18272301958], [95994.63497604289, 102316.35230028478], [96921.72437226822, 98871.70066343644], [96875.32986395528, 99220.5014365399],
[96658.78507475171, 101671.52754812177]]
[[199968.18349177408, 188392.34768358883], [192595.65496090706, 237274.3036962648], [198427.1236183697, 187414.22237599595], [190015.0805513745, 220076.295673958
6], [193486.3524879159, 231911.23157262686]]

metrics for random forest, with RMSE, MSE, and MAE for train_and_test
0.9745170018860927 0.9202516114558377
85615.04342817182 132284.7333876817
47685.325744750655 77122.98158238827
[[0.9735789263951878, 0.8958907296932568], [0.9744087365551072, 0.9210897775115411], [0.9718677029347136, 0.9159784893435712], [0.9723629414571852, 0.90654574707
77247], [0.9771928232699933, 0.8685862704425774]]
[[47753.4580944456, 78664.53277286238], [47676.55412456391, 77147.65941346141], [47256.129919879015, 80924.22996209176], [47736.45330294609, 83026.06077946187],
[47228.771904304194, 81022.5198916232]]
[[86020.51060451134, 161356.33598338102], [84599.63773533347, 140876.53413187474], [87479.07062740167, 153990.35707237537], [86432.69680308408, 164355.6451388110
3], [78217.47926144357, 197662.62311766951]]

metrics for Neural Network, with RMSE, MSE, and MAE for train_and_test
0.9539594054222107 0.9050804376602173
115078.82515910562 144320.0425720558
67507.140625 84070.1796875
[[0.961050808429718, 0.9328145384788513], [0.9619346857070923, 0.8491305112838745], [0.9617573618888855, 0.9086839556694031], [0.9655402898788452, 0.797531306743
6218], [0.9658370018005371, 0.903423011302948]]
[[64337.2578125, 83551.9296875], [63546.90625, 87995.9609375], [63913.015625, 83460.296875], [61681.7265625, 87498.4609375], [63944.21875, 79115.09375]]
[[101854.47187040931, 143073.72846193673], [102723.14512319023, 198623.58498426114], [103902.85185691489, 148337.89061463697], [94042.56125818778, 263788.7154523
4835], [99154.2862411900, 145449.8953729428]]
Random Forest performs best
(Random Forest) Model exported
(base) lebinhminh@Shions-Macbook-Air submission %
```

Bước 3: Khởi tạo API:

uvicorn app.API:app --reload


```
○ (base) lebinhminh@Shions-Macbook-Air submission % uvicorn app.API:app --reload
INFO: Will watch for changes in these directories: ['/Users/lebinhminh/Downloads/FoML-main/submission']
INFO: Uvicorn running on http://127.0.0.1:8000 (Press CTRL+C to quit)
INFO: Started reloader process [74623] using StatReload
INFO: Started server process [74626]
INFO: Waiting for application startup.
INFO: Application startup complete.
```

Bước 4: Khởi động ứng dụng:

streamlit run app/interface.py

```
conda activate base
● (base) lebinhminh@Shions-Macbook-Air submission % conda activate base
○ (base) lebinhminh@Shions-Macbook-Air submission % streamlit run app/i
nterface.py

You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://192.168.200.26:8501

For better performance, install the Watchdog module:

$ xcode-select --install
$ pip install watchdog
```

Giao diện ứng dụng sẽ được khởi động

Bước 5: Nhập các thông tin và bấm nút Predict để tiến hành dự đoán giá xe (với các hãng xe trong 33 hãng đề cập trong file README.md trên github)

Enter the car you want to predict the price

Enter a car brand

Ford

Car Year

2018

Transmission

Automatic

Seller

Individual

Fuel

Petrol

Owner

Second

Input engine (CC) from 624 - 3604

1000.00

Input engine horsepower (bph) from 33 - 400

60.00

Predict

Estimated Price (converted to USD): \$4,221.27

(Trong ứng dụng, giá tiền được biến đổi thành USD để dễ ước lượng)

Kết luận và hướng phát triển:

Kết luận:

Dự án này đã xây dựng hoàn chỉnh quy trình xây dựng mô hình học máy để dự đoán giá xe ô tô cũ. Quy trình bao gồm các bước làm sạch dữ liệu, tiền xử lý, biến đổi đặc trưng và lựa chọn đặc trưng cũng như việc lựa chọn các kiến trúc tối ưu để giải quyết bài toán dự đoán.

Hướng phát triển:

Dự án này còn có thể tối ưu việc dự đoán giá xe cũ dựa trên việc tiến hành tạo thêm dữ liệu bằng kỹ thuật data augmentation. Lý do của việc sử dụng data augmentation là để sinh ra thêm nhiều dữ liệu huấn luyện phong phú hơn, giúp mô hình hiểu được các mẫu dữ liệu phức tạp và đa dạng hơn, giúp cho việc dự đoán tốt hơn cũng như áp dụng các mô hình tiên tiến hơn, ví dụ như TabNet, giảm thiểu khả năng overfitting trên mô hình phức tạp hơn.

Việc áp dụng các kiến trúc DNN (mô hình mạng nơ-ron sâu) để tối ưu việc lựa chọn các đặc trưng và khái quát như việc transfer learning các kiến trúc Neural Network tối ưu cho dữ liệu bảng như mô hình TabNet [1]. Lý do việc lựa chọn TabNet là vì TabNet sử dụng cơ chế attention, giúp mô hình có thể tùy biến lựa chọn đặc trưng phù hợp qua nhiều bước suy luận, giúp mô hình có thể tối ưu và khái quát hoá tốt hơn trong thực tế.

Tham khảo:

1. Arik, S. Ö., & Pfister, T. (2021). TabNet: Attentive Interpretable Tabular Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(8), 6679-6687.