



Trường Đại học Công nghệ Thông tin

Lớp DS108.P21

Giảng viên hướng dẫn:

TS.Nguyễn Gia Tuấn Anh - CN.Trần Quốc Khanh

## HỆ THỐNG HỎI ĐÁP HỖ TRỢ TƯ VẤN LAPTOP CHO NGƯỜI DÙNG

- Nguyễn Trần Ngọc Ty  
23521758
- Đoàn Thị Ngọc Anh  
23520042
- Nguyễn Thị Minh Tiến  
23521579

NHÓM 22



# CONTENT

**01** Data Collection



**03** Modelling



**05** Conclusion



**02** Data transformation  
Data flow



**04** Deploy





# DATA COLLECTION



## Data sources

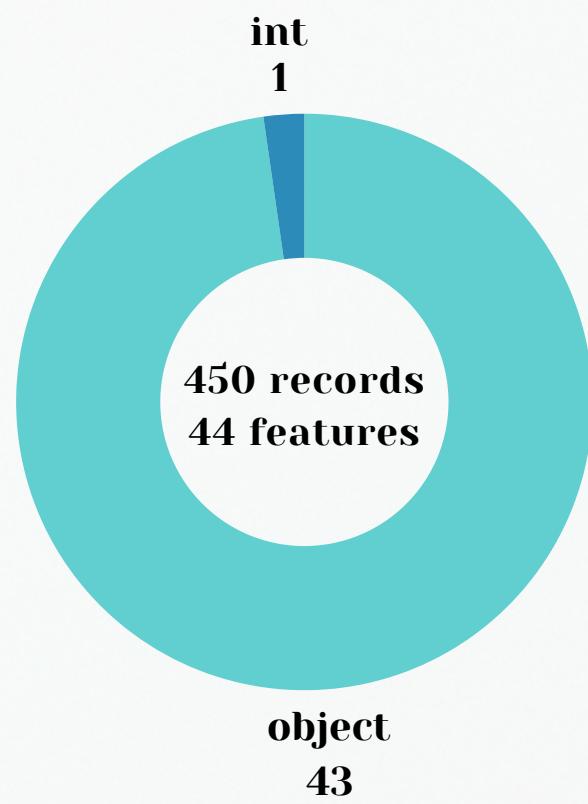
Thegioididong.com (Scrapy)  
cellphones.com.vn (API + Scrapy)



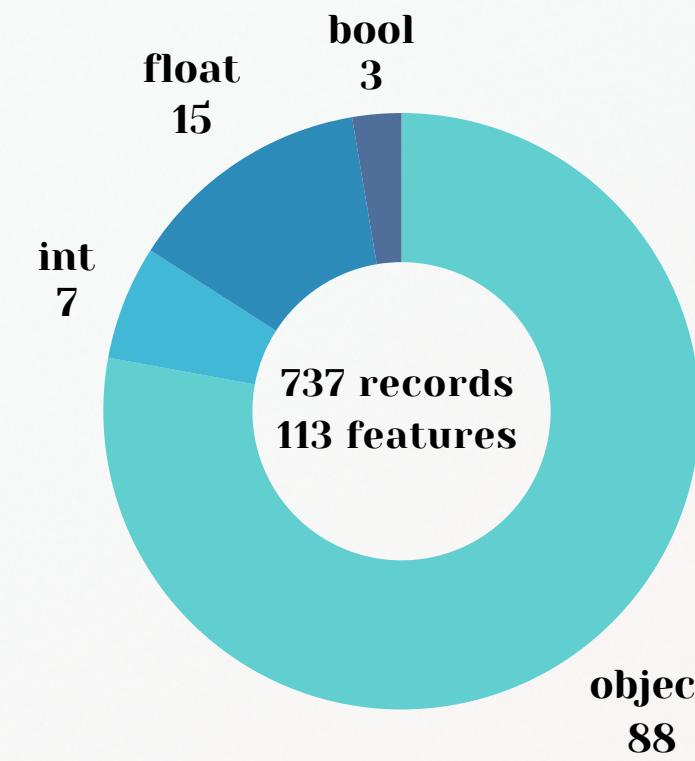
## Data extraction

Page source → CSV files

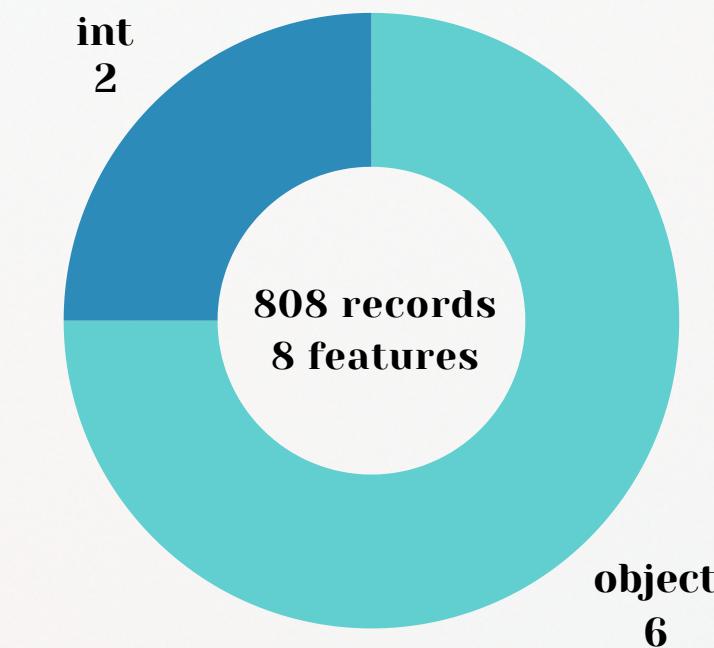
**Thegioididong Laptop**



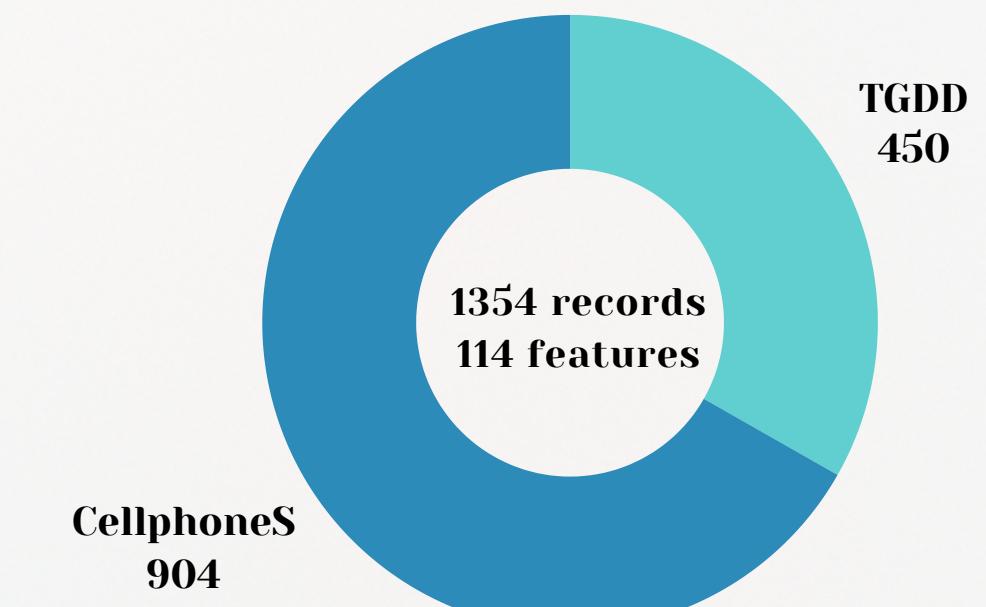
**CellphoneS Laptop**



**CellphoneS Laptop Variant**

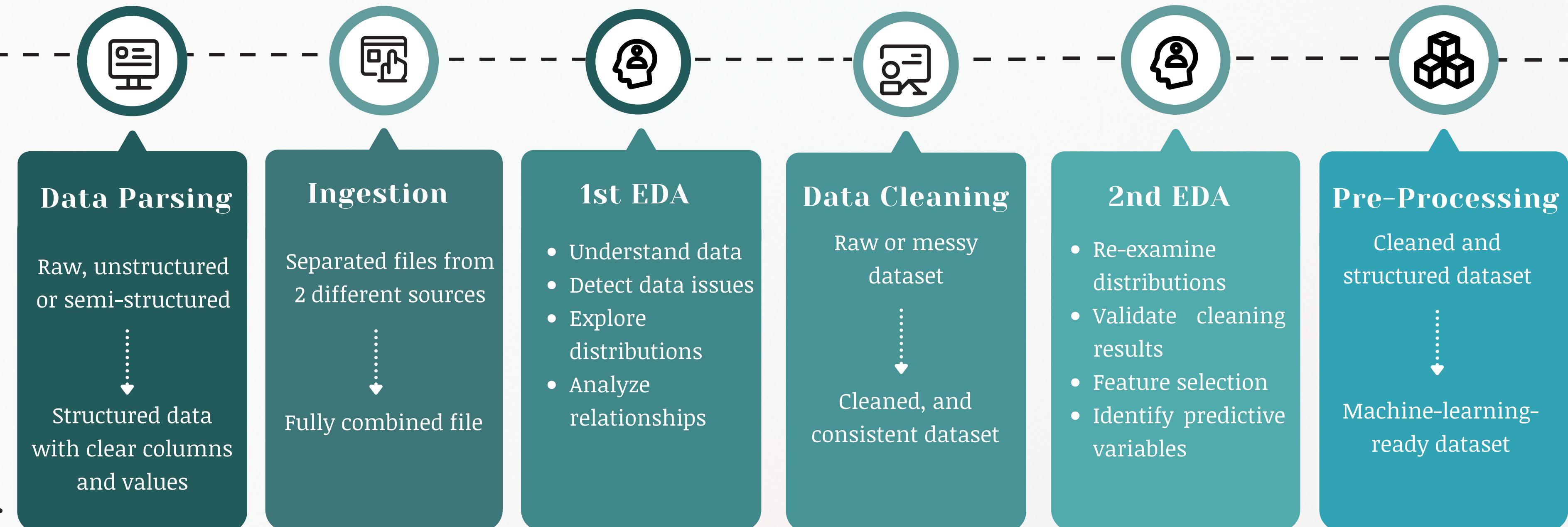


**Total**





# DATA TRANSFORMATION



## DATA PARSING

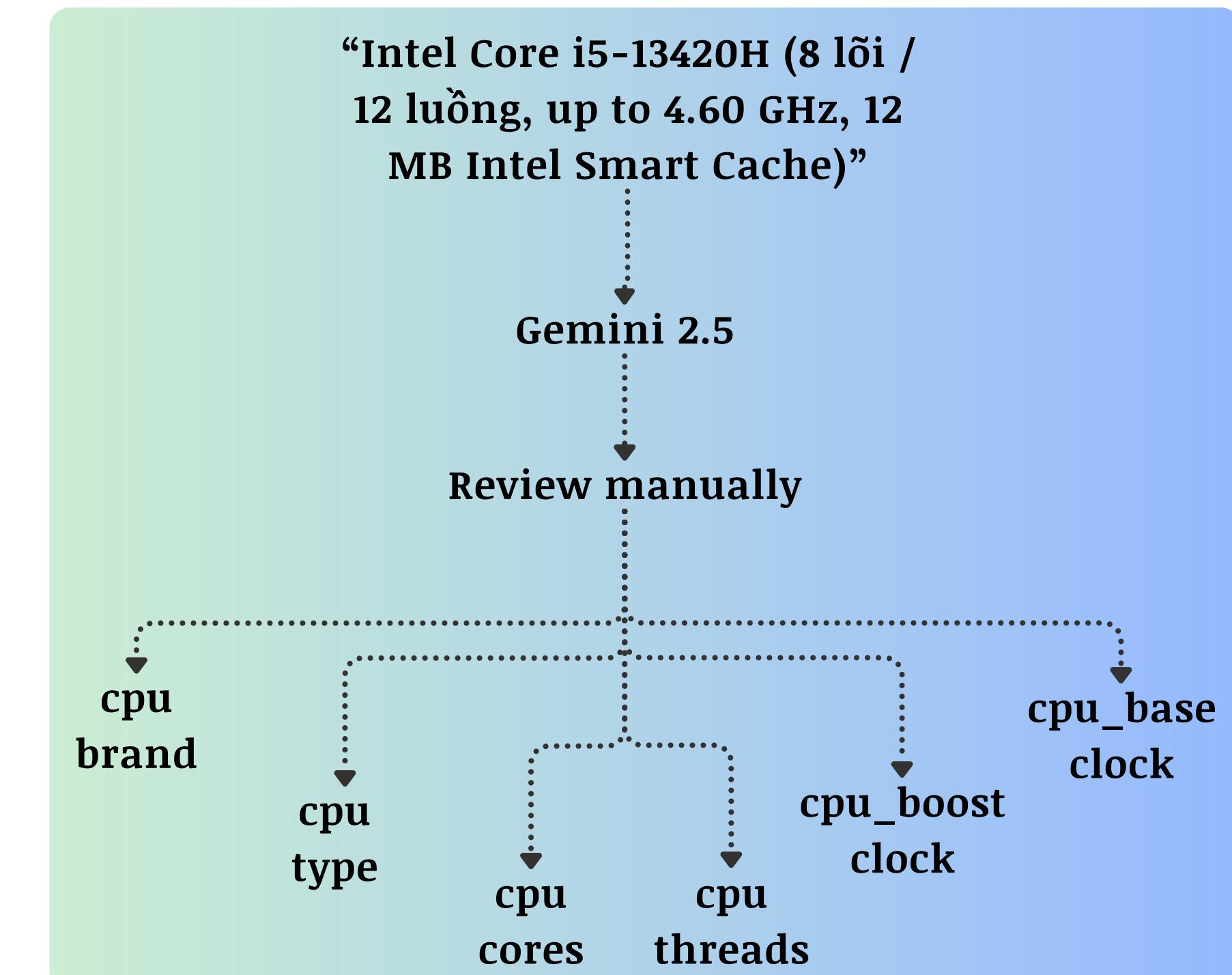
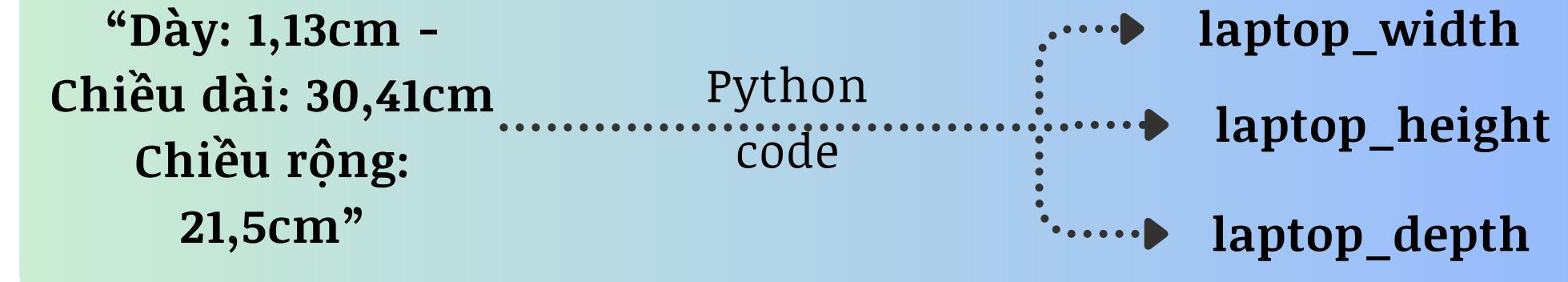
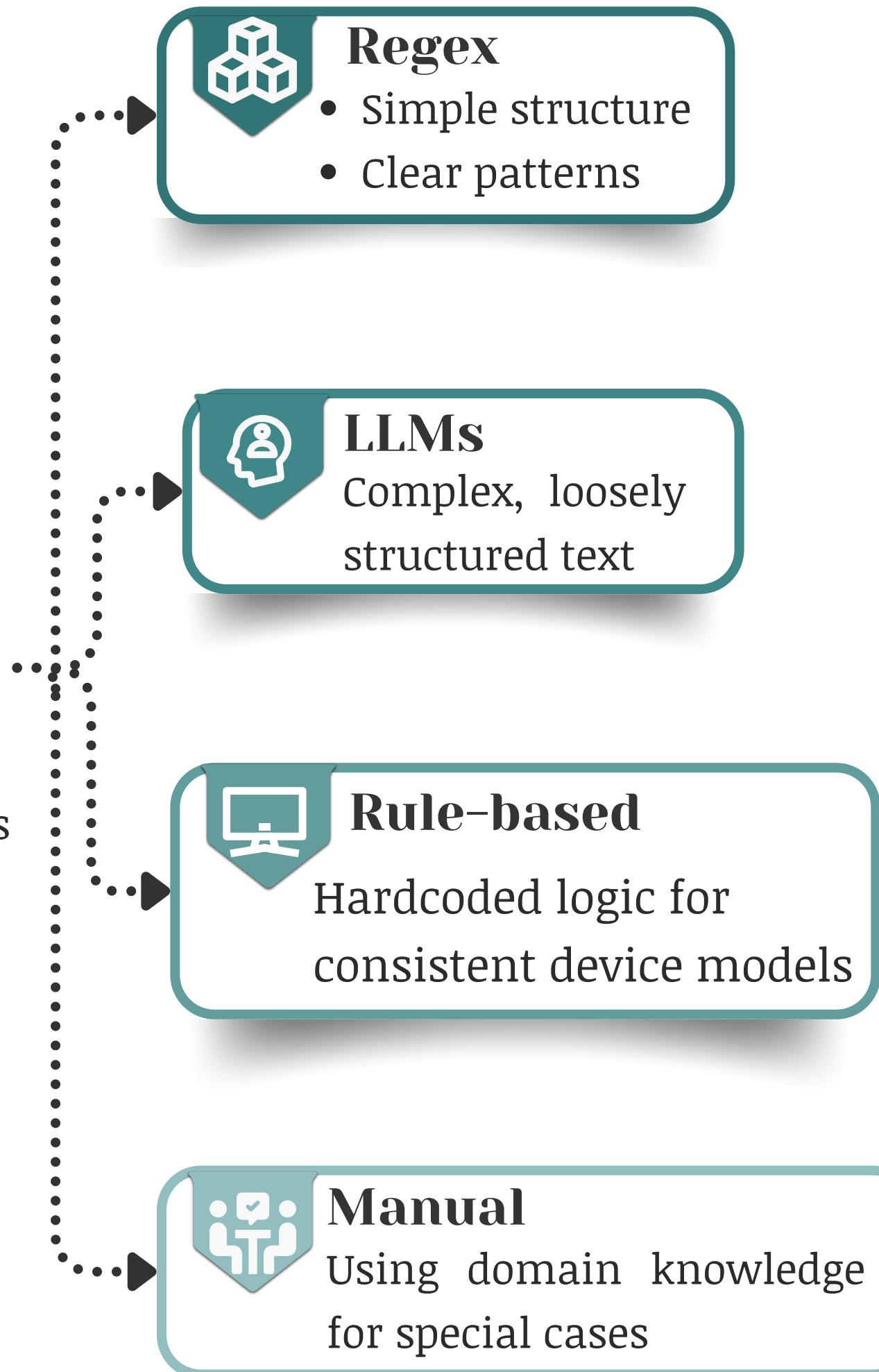
114 → 52 features

Keep

Delete

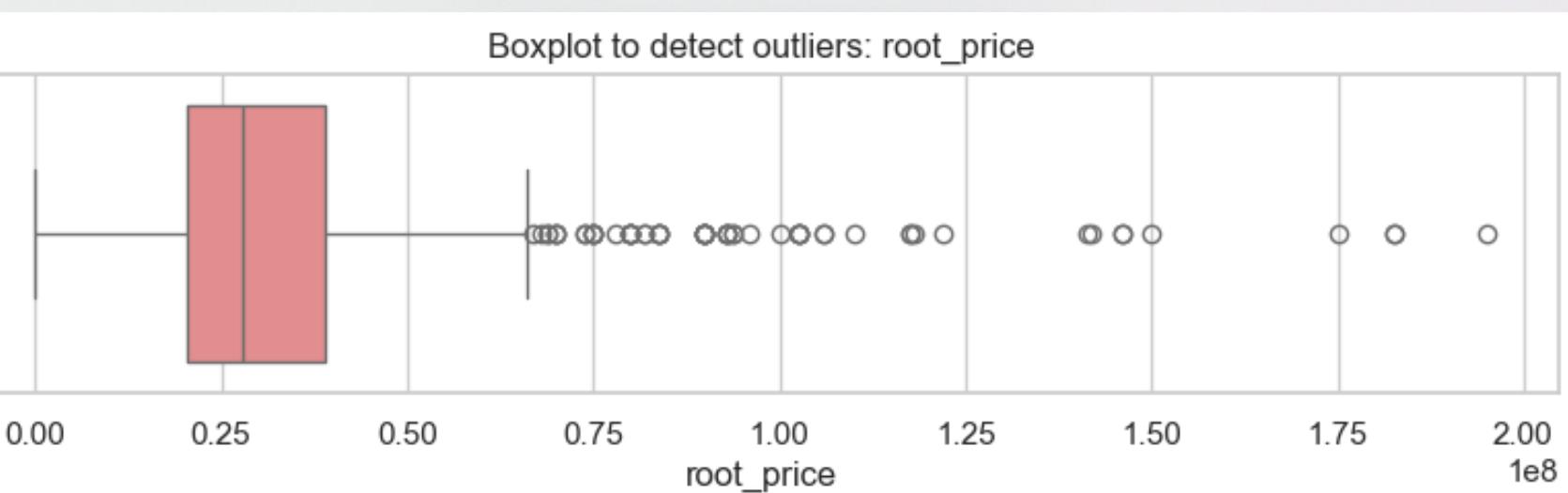
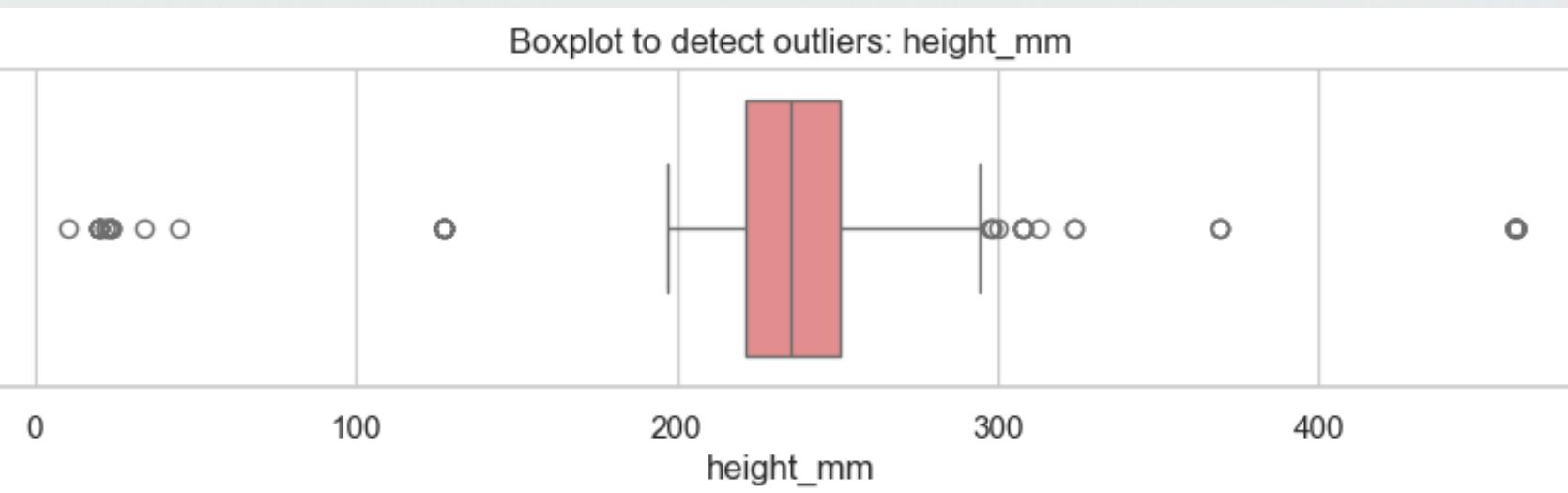
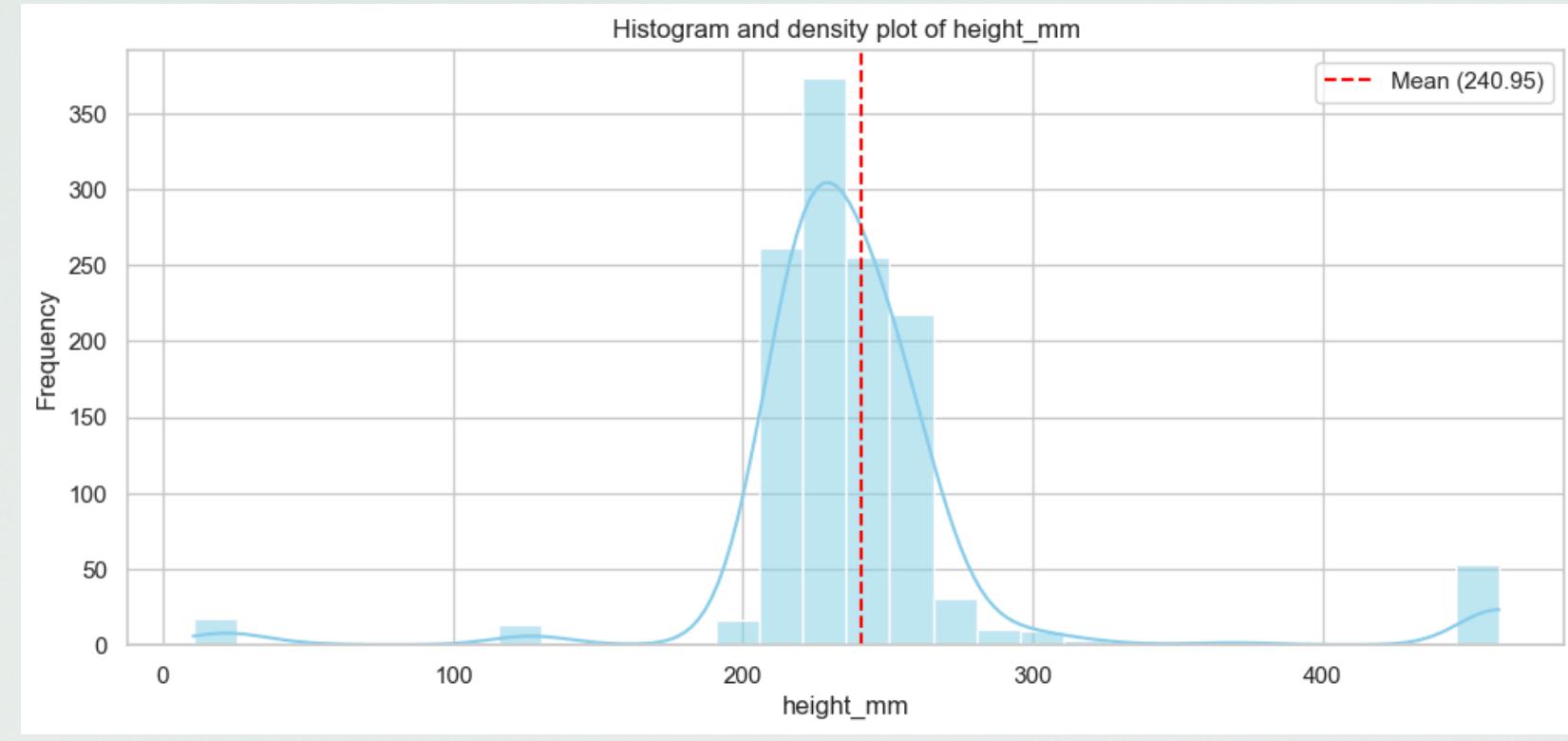
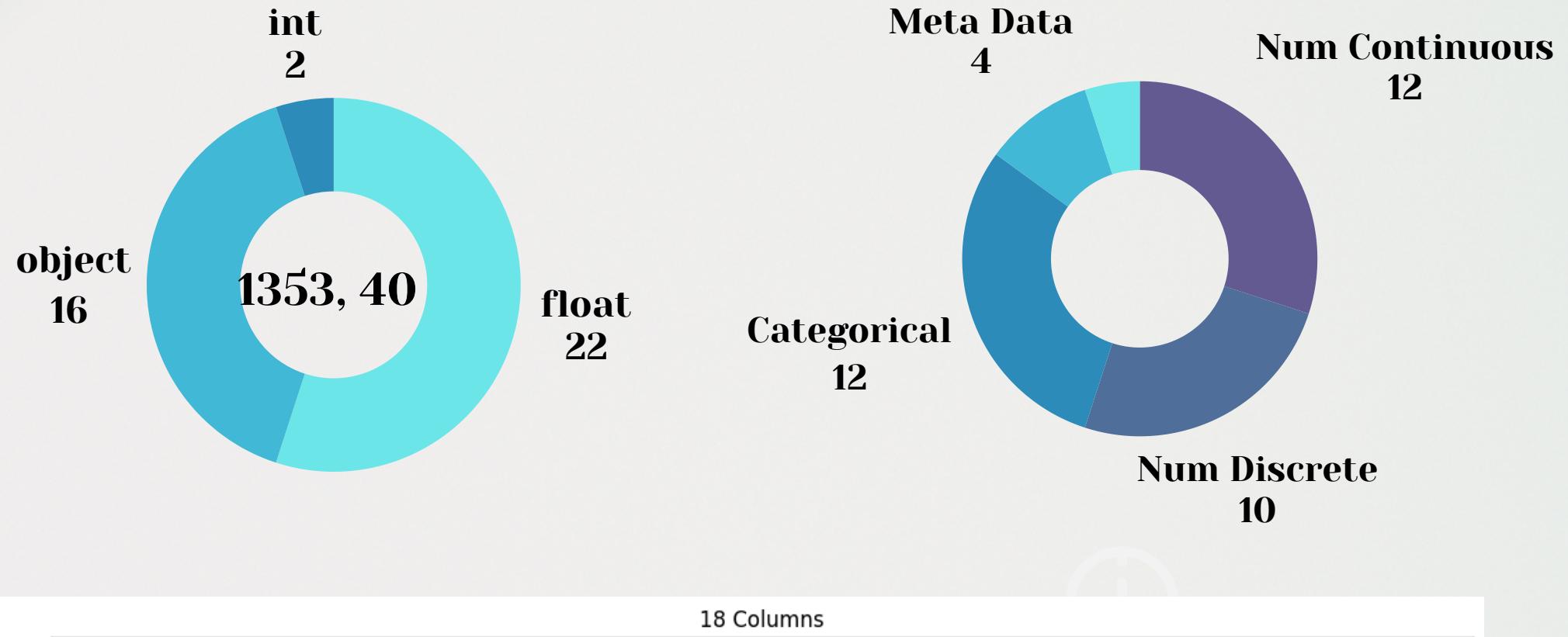
Not sure

2.1



## 2.2

# DATA INGESTION → 1ST EDA



## 2.3. DATA CLEANING



### Outlier Handling

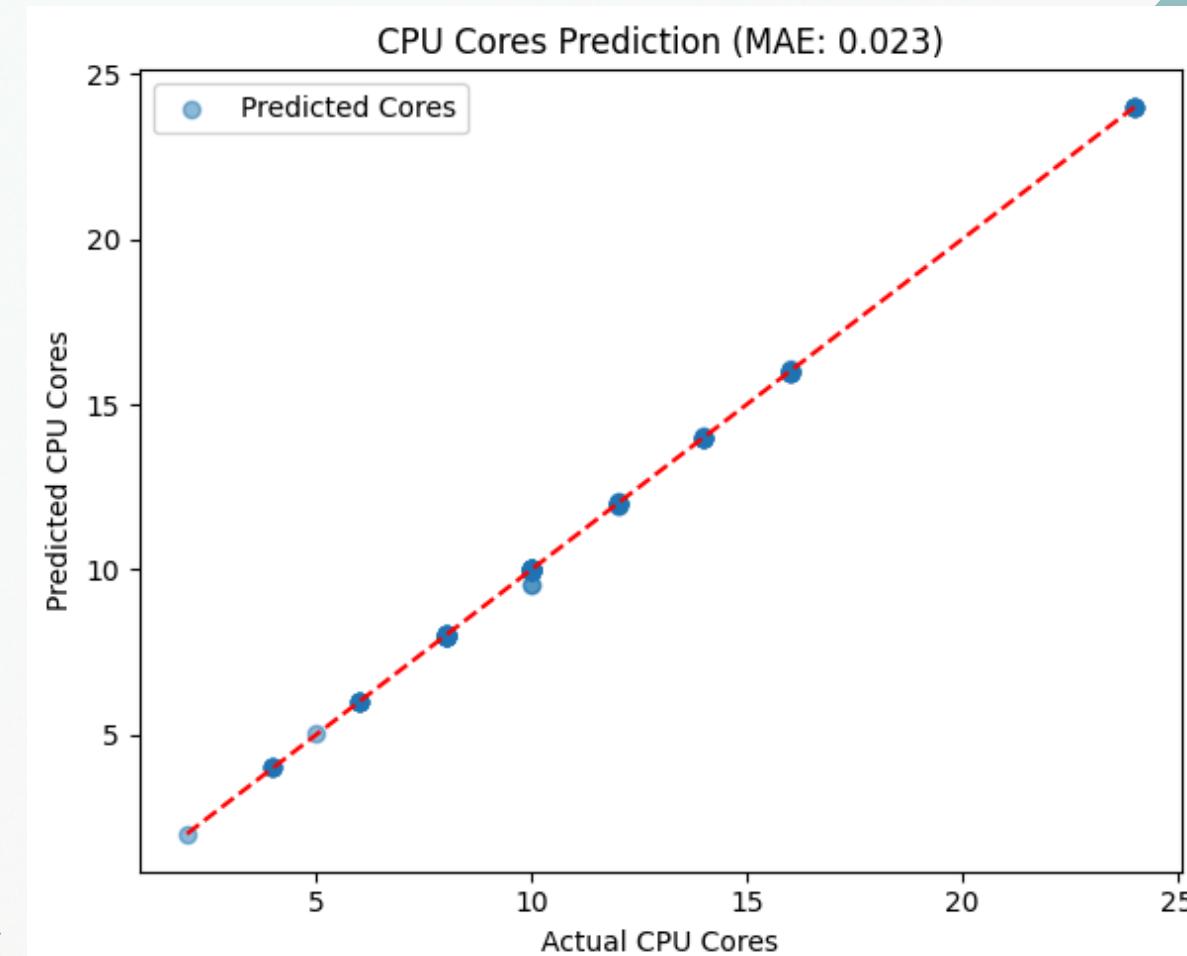
- Manual detection
- Domain-specific correction



### Rebuilding missing data

- MAR
  - Statistical imputation  
(mean, mode, model-based  
- XGBoost và embeddings)

LLMs	Accuracy	Continuous	Discrete
Gemini 2.5 flash	98,42%	97,76%	99,08%
Gemini 2.0 flash	97,60%	98,68%	96,52%
Llama 3.3	62,90%	60,74%	65,01%



### Standardize data format

Convert values to consistent and correct types (e.g. number, text)

8M

Min of root\_price

195M

Max of root\_price

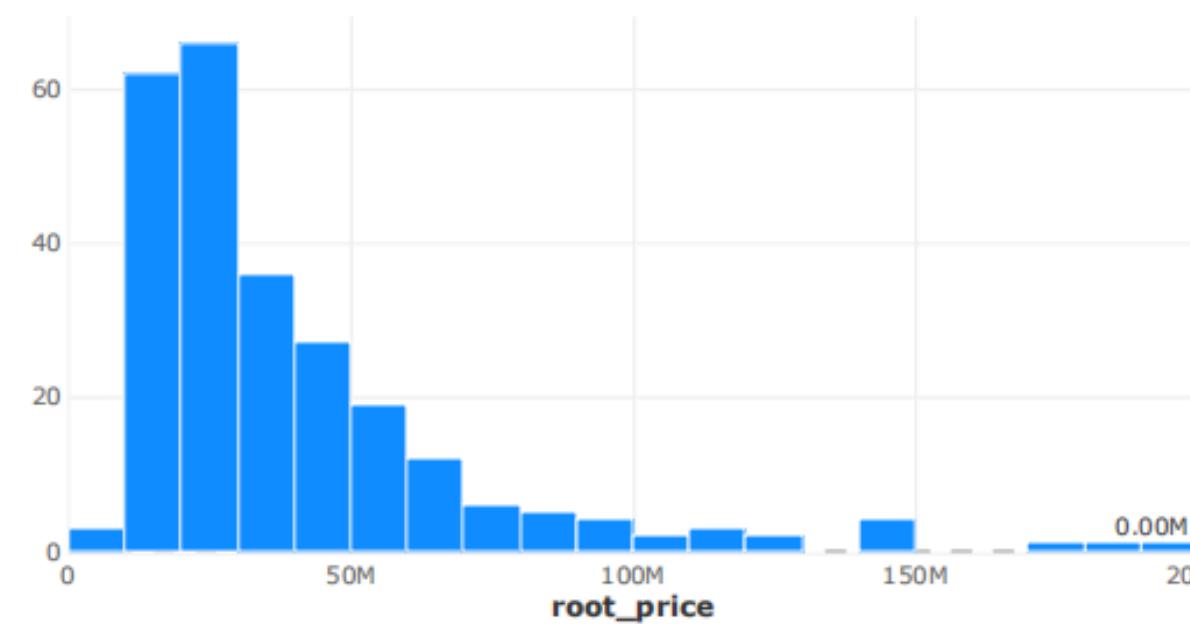
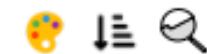
33.90M

Average of root\_price

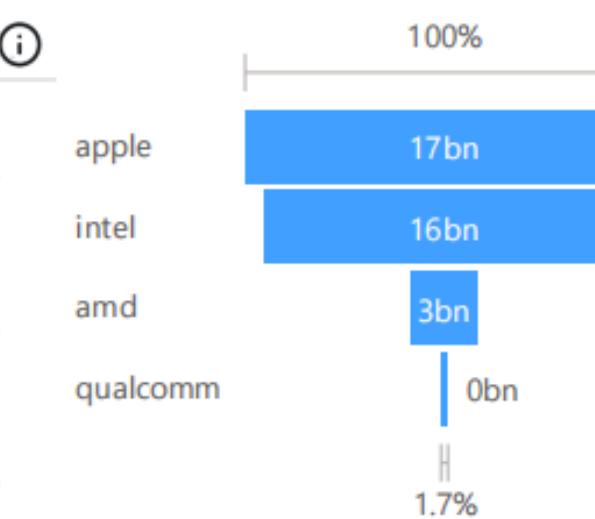
29M

Median of root\_price

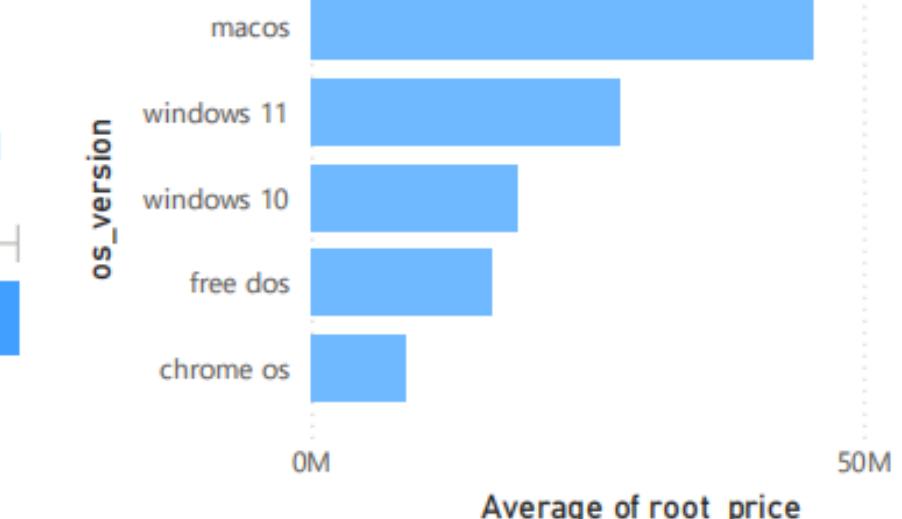
Price of laptop



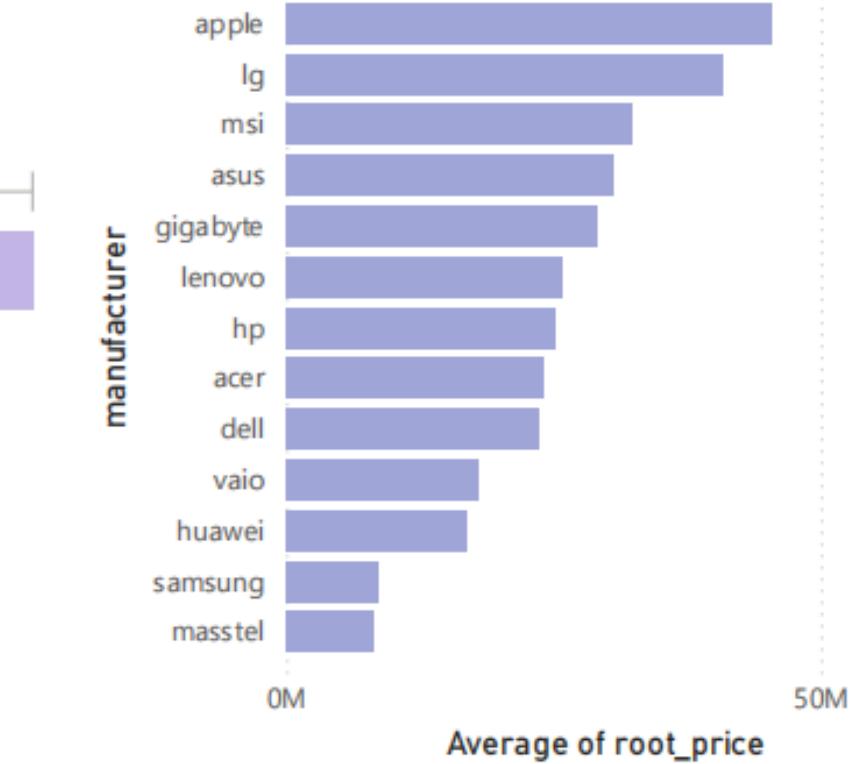
Sum of root\_price by cpu\_brand



Average of root\_price by os\_version



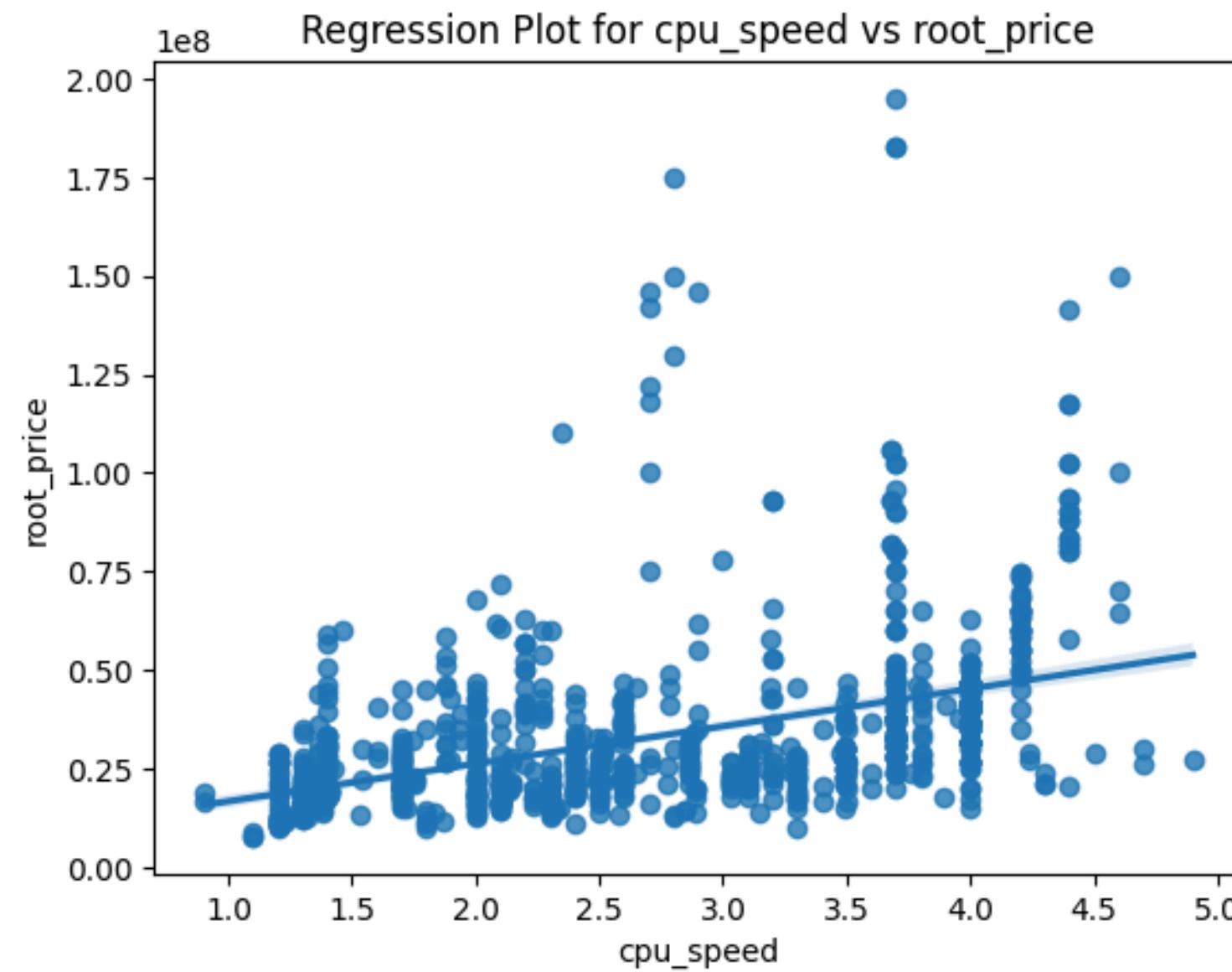
Average of root\_price by manufacturer



Phân tích đa biến giữa price và các feature số khác:

- Hầu hết các cột đều có tương quan dương với price
- Trừ các biến liên quan đến kích thước thì không có mối quan hệ rõ ràng

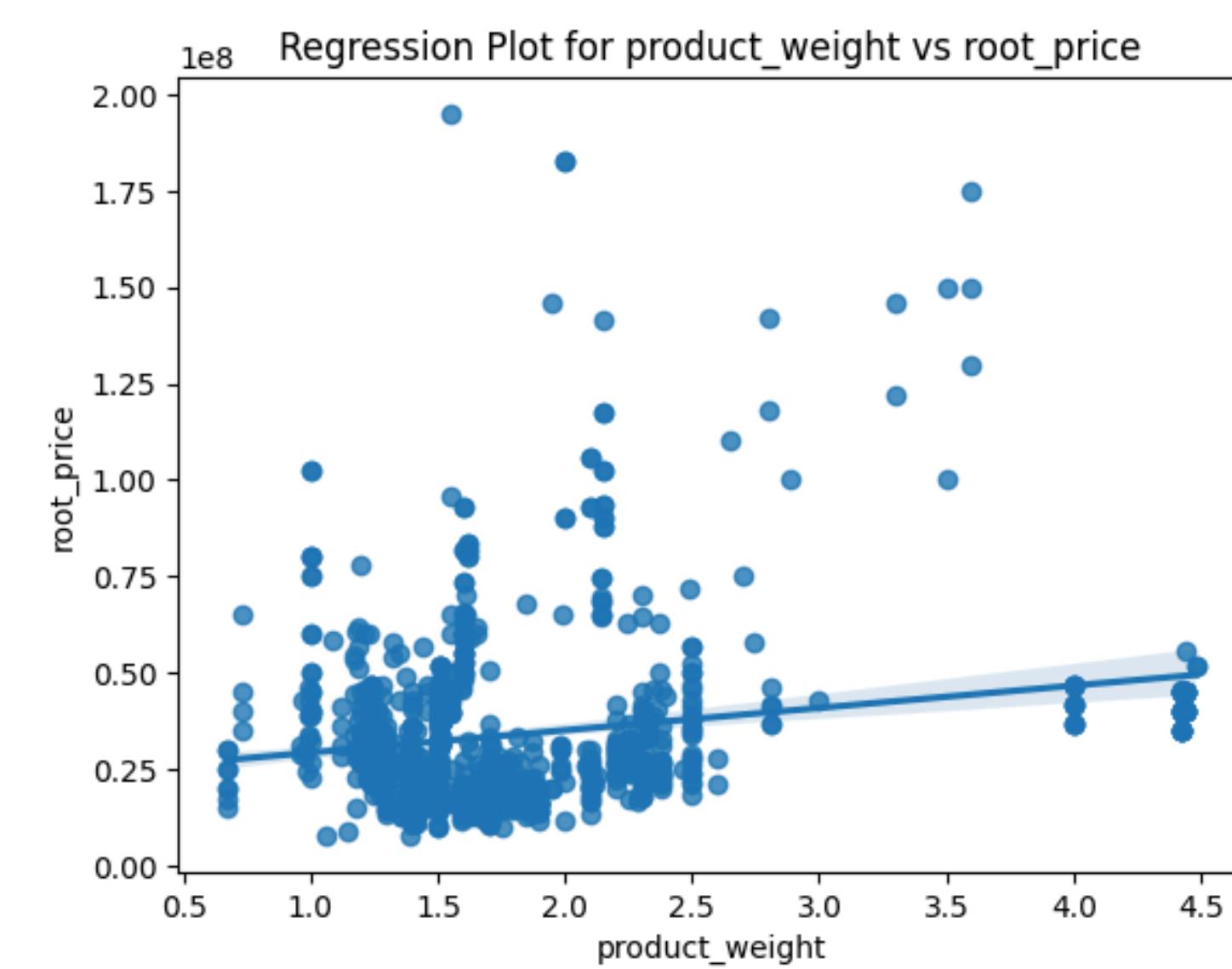
# 2ND EDA



Có ảnh hưởng

Correlation Coefficient: 0.4508

P-value: 5.4644e-55

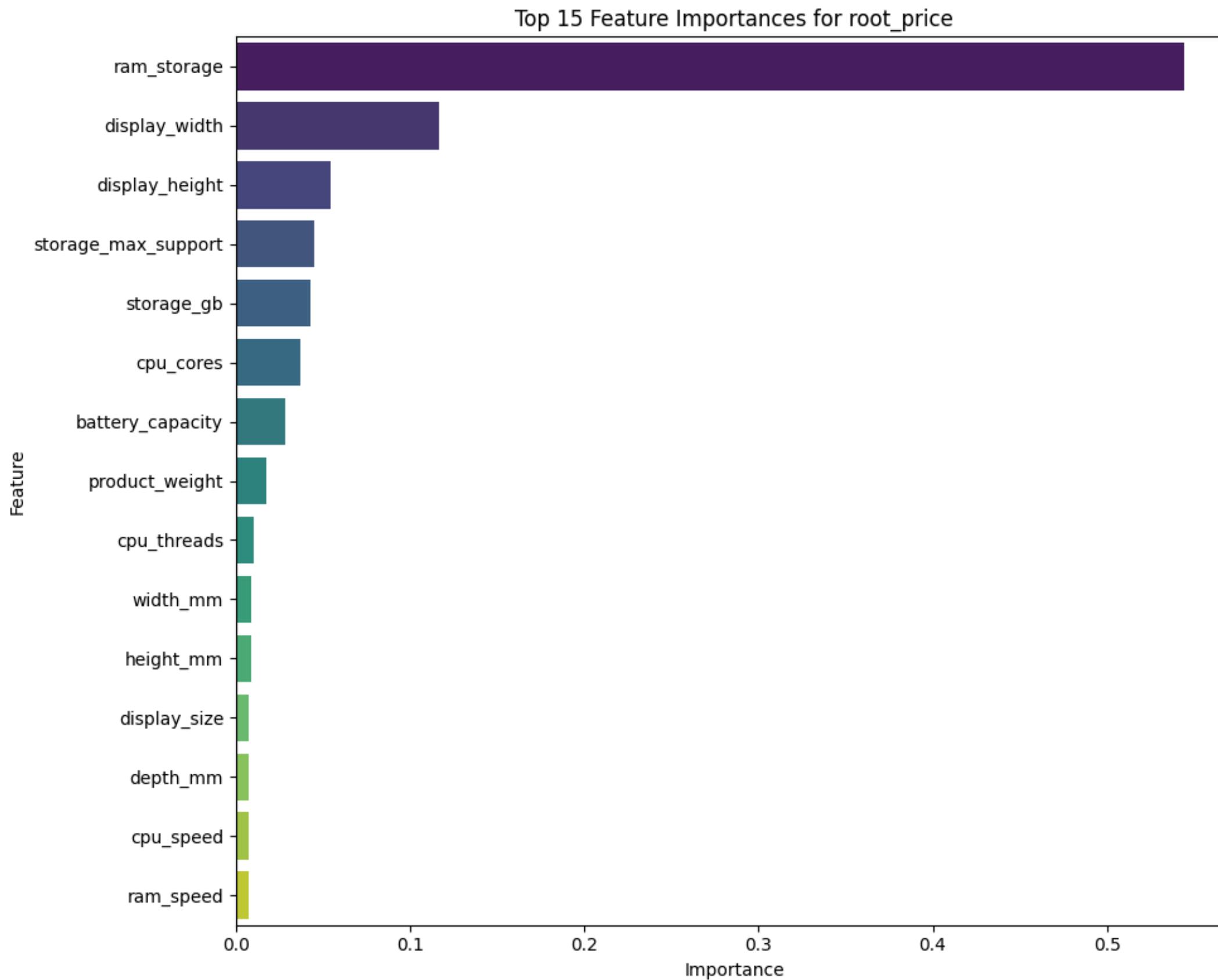


Không ảnh hưởng

Correlation Coefficient: 0.1905

P-value: 2.9614e-10

# 2ND EDA



## Feature importances for root\_price

Features	Feature importances
ram_storage	0,544631
display_width	0,116766
display_height	0,054055
storage_max_support	0,045027
storage_gb	0,042426
cpu_cores	0,036955
battery_capacity	0,028198
product_weight	0,017607
cpu_threads	0,01035

**Thực hiện kiểm định ANOVA với các categorical features với giả thiết với biến target “root\_price”:**

- H0: hai biến độc lập
- H1: hai biến phụ thuộc nhau

**45 features → 22 features**

Index	Features	F-statistic	P-value	Đánh giá
0	material	92,1582	0	Ảnh hưởng
1	manufacturer	19,3462	0	Ảnh hưởng
2	ram_type	46,2786	0	Ảnh hưởng
3	os_version	52,236	0	Ảnh hưởng
4	laptop_color	2,1626	0,018	Ảnh hưởng
5	vga_type	0,017	0,8962	Không ảnh hưởng
6	vga_brand	66,6873	0	Ảnh hưởng
7	laptop_camera	247,5311	0	Ảnh hưởng
8	cpu_brand	67,3149	0	Ảnh hưởng

## One-hot Encoding

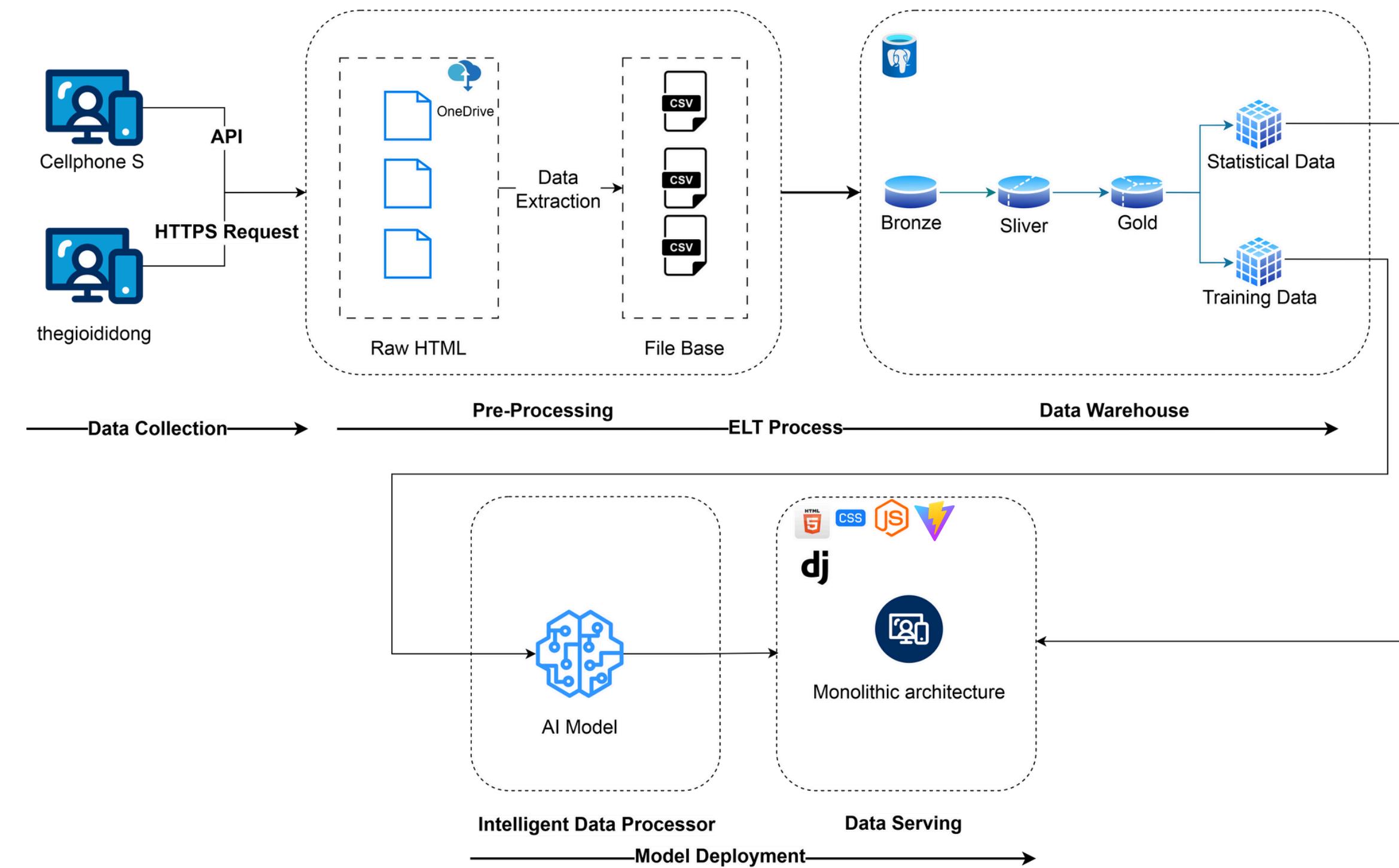
```
manufacturer', 'cpu_brand',  
'material', 'os_version',  
'laptop_color', 'vga_brand',  
'laptop_camera', 'ram_type',  
'laptop_sang_tao_noi_dung',  
'do_hoa_ky_thuat',  
'cao_cap_sang_trong'
```

## Standard Scaler

```
'storage_max_support',  
'storage_gb', 'display_width',  
'cpu_threads', 'cpu_cores',  
'ram_speed', 'cpu_speed',  
'ram_storage', 'ram_slots',  
'battery_capacity', 'display_height'
```

2.6

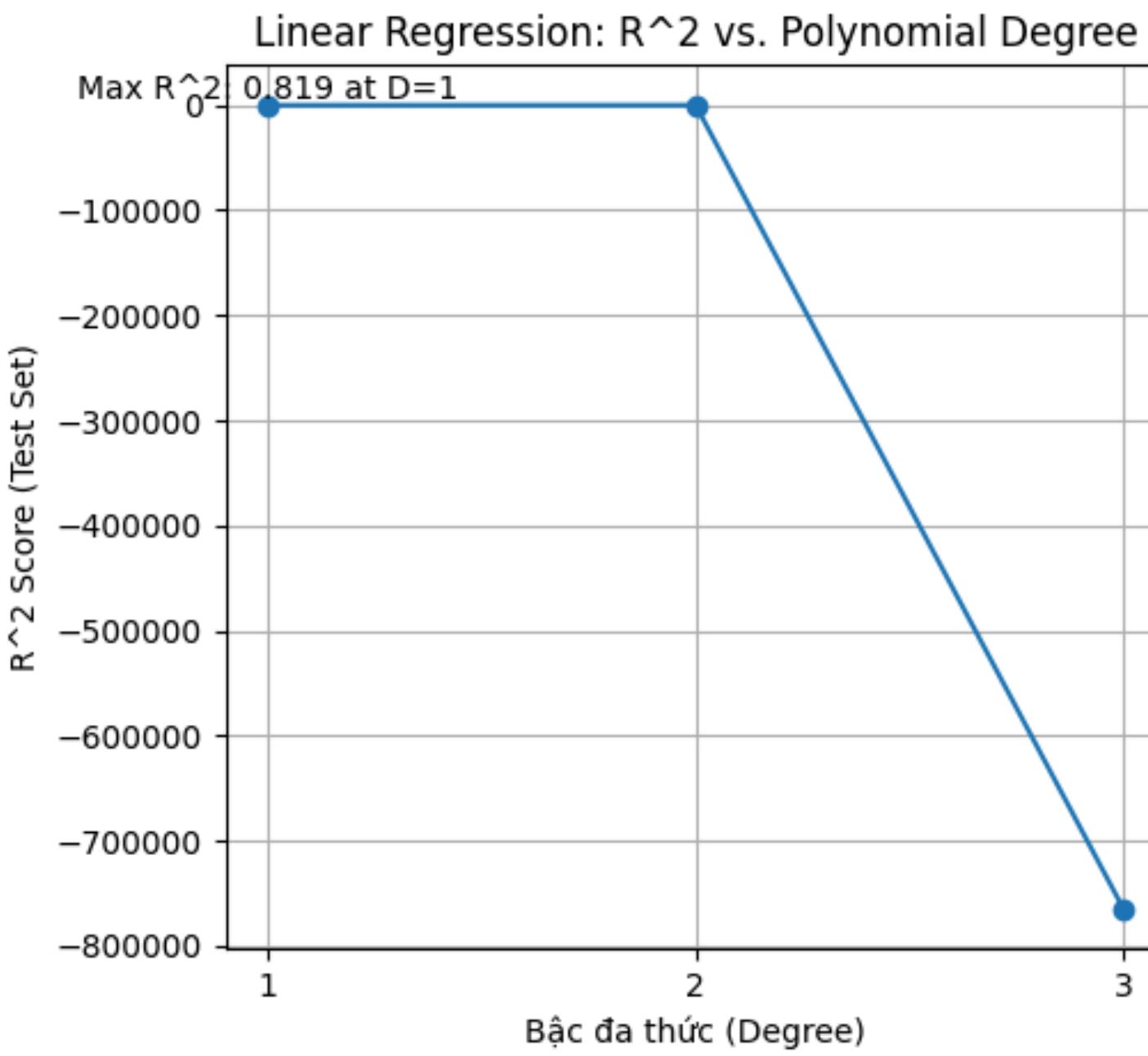
# DATA FLOW



## 03. MODELLING

Linear Regression với  
PolynomialFeatures  
xác định bậc để nắm  
mối quan hệ phi tuyến

Kết quả: bậc n = 1 cho kết quả  
dự đoán mô hình ổn định nhất

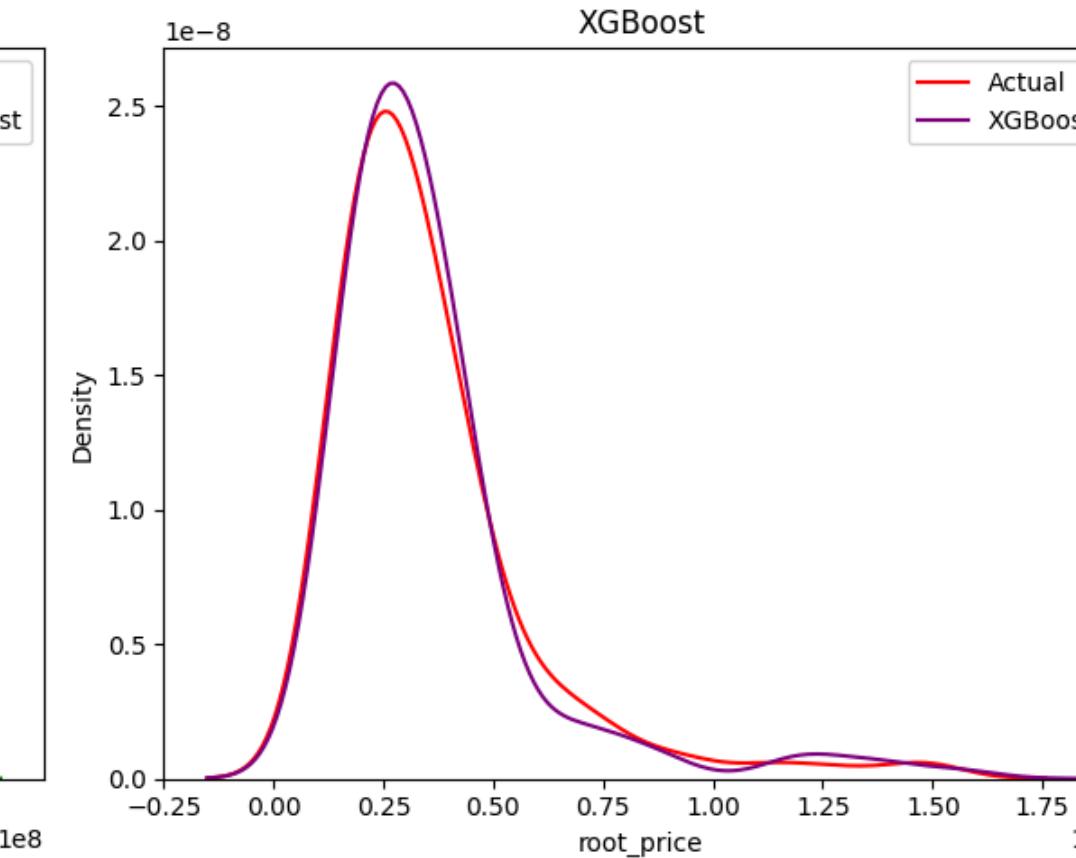
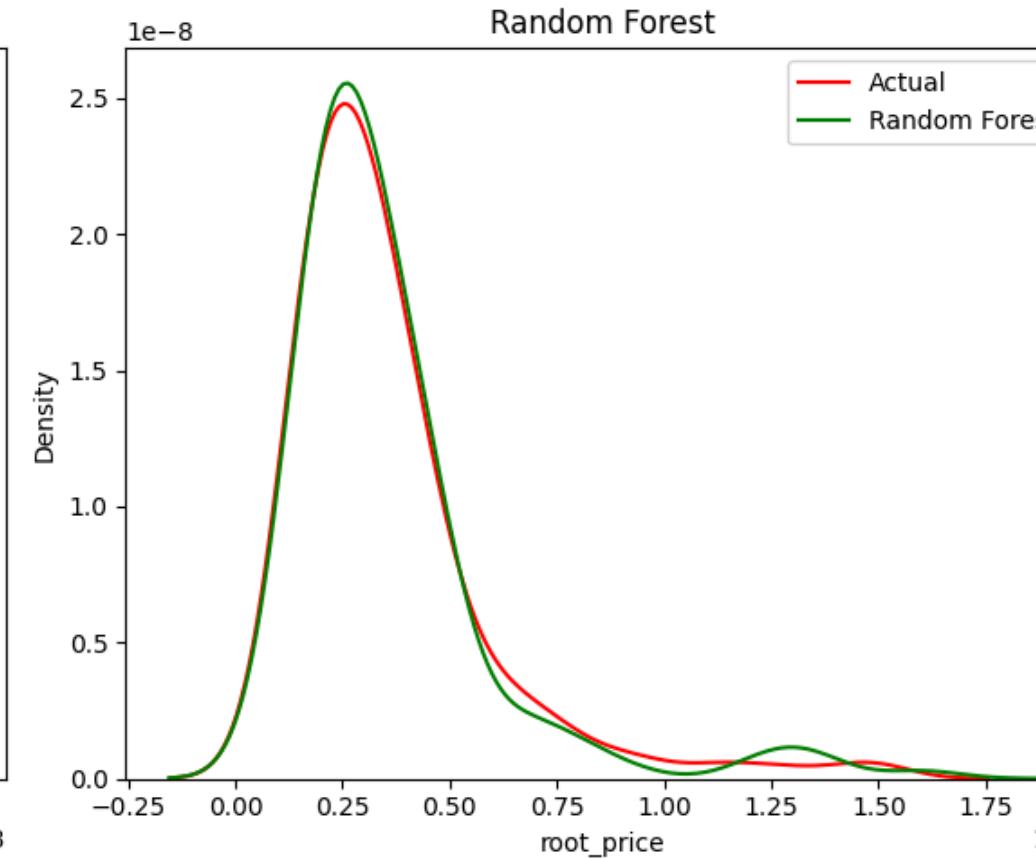
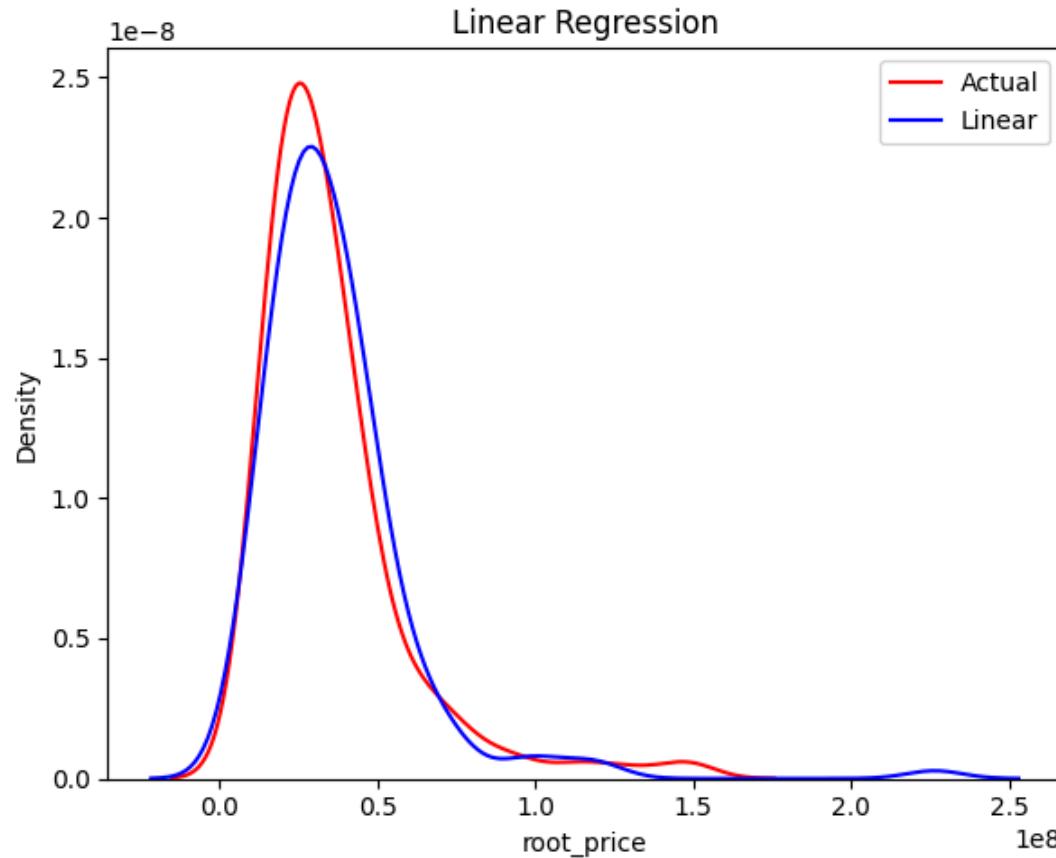


# 03. MODELLING

## Training và Testing với 3 mô hình với Cross Validation là 5

Trực quan hóa

Mô hình	R <sup>2</sup> (Train)	R <sup>2</sup> (Test)	RMSE (Test)	Cross-validation R <sup>2</sup>
Linear Regression	0,891	0,819	10,060,950.43	0,753
Random Forest	0,989	0,916	6,853,476.04	0,840
XGBoost	<b>0,983</b>	<b>0,937</b>	<b>5,938,349.67</b>	<b>0,853</b>

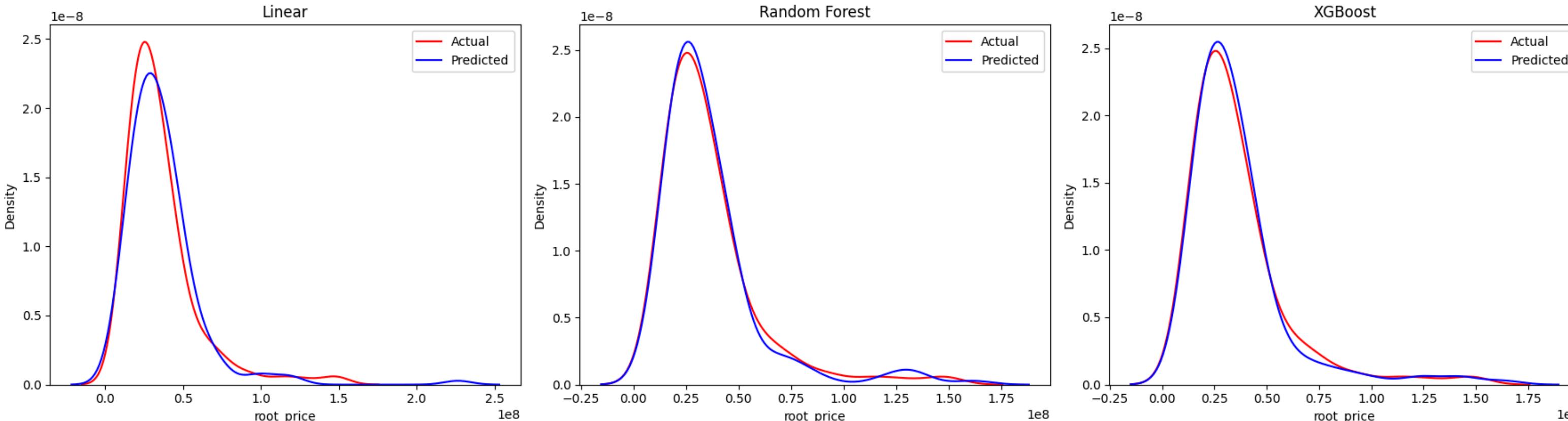


# 03. MODELLING

## Tunning XGBoost và RandomForestRegressor sử dụng GridSearchCV

Mô hình	Train R <sup>2</sup>	Test R <sup>2</sup>	RMSE (Test)	CV R <sup>2</sup> (mean)
XGBoost	<b>0,99080</b>	<b>0,94923</b>	5,320,958	<b>0,87542</b>
Random Forest	0,98882	0,91869	6,734,008	0,84166
Linear Regression	0,89092	0,81850	10,060,950	0,75257

## Trực quan hóa



4

# DEMO APPLICATION

LAPBOT  
Laptop Computer Market

Home Products Chat

Giới thiệu cho tôi các laptop Nvidia, có giá trên 30 triệu

Dưới đây là một số laptop chúng tôi gợi ý dựa trên thông tin bạn đã cung cấp :

Với nhóm Cỗ máy Gaming Cao cấp, giá cần chuẩn bị khoảng : **39.536.348đ**

Laptop Gigabyte G6 KF-9RC56 KFOHJJANIVN000-Đen  
**26.990.000đ** 30.990.000đ  
Đồ họa - Kỹ thuật | Mạng nhẹ | Gaming

Laptop MSI Cyborg 15 A1VEK-05 3VN -Đen  
**27.590.000đ** 30.990.000đ  
Cao cấp - Sáng trọng | Học tập - Văn phòng | Gaming

Laptop Asus TUF Gaming A15 F A507NV-LP061W-Xám  
**27.590.000đ** 30.990.000đ  
Gaming

Với nhóm Workstation Sáng tạo chuyên nghiệp, giá cần chuẩn bị khoảng : **104.928.376đ**

Laptop MSI Prestige 14 AI Studio C1VEG-056VN-Xám  
Nhập tin nhắn...  
Delete

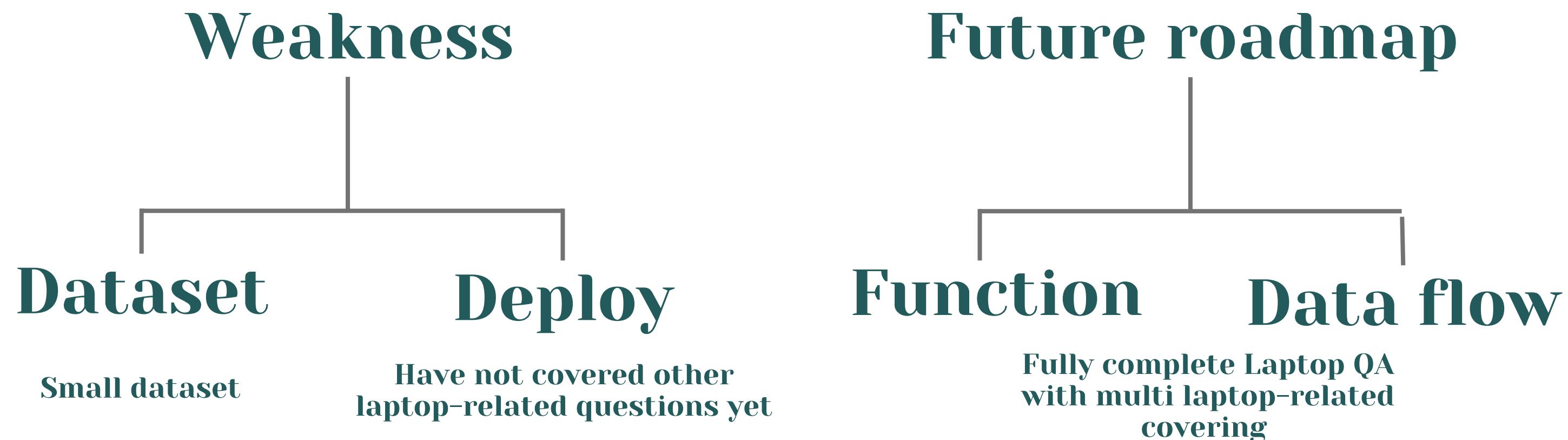
Laptop MSI Prestige 14 AI Studio C1VEG Ultra 7 155H/32GB/1TB/6GB RTX4050/Win11 (056VN)

Laptop ASUS ROG Strix G16 G61 4JU-N4450W-Xám

Lapbot - AI thông minh tìm kiếm laptop phù hợp nhất với bạn

# 5. CONCLUSION

## LAPQA



Thank you for listening!