

13-09-2025

Agenda:

FE - T

- backward elimination → p-value
 - remove highest p-value column
 - Train the model again (ML)
- forward selection / elimination → R² value
 - R²-absolute
 - add the relevant ones
 - filter the irrelevant ones
 - p-value

p-value: measure the prob of observing the data that you have

< 0.0

H₀: states that coefficient for that feature is 0

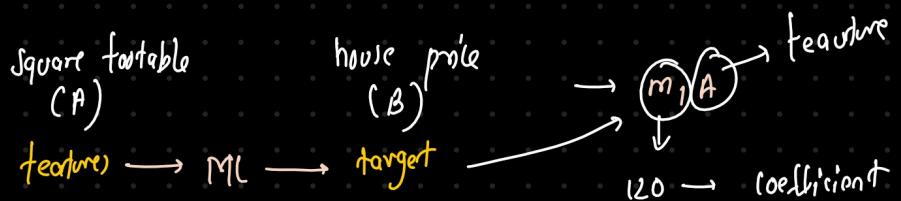
$$y = m_0 + m_1 x_1 + \dots + m_n x_n$$

< 0.05 → reject H₀

> 0.05 → fail to reject H₀

Example:

regression model → Linear Regression
 ↓ ↓
 Type Name of
 algo



$$SE = 40$$

$$T\text{-statistic} = \frac{\text{coefficient}}{SE} = \frac{120}{40} = 3 \rightarrow \underline{T\text{-table}} \rightarrow \underline{3.182}$$

$p(\text{value})$
 ↓
 not convert
 value

$p\text{-value} < 0.05$
 → fail to reject H_0

→ forward selection → $\chi^2 (0-1)$

	Actual_Target_Column (A)	Model predicted (B)	Mean (C)
	2	2.5	3.5
	5	4.0	3.5
	3	4.0	3.5
	4	3.5	3.5

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{total}}} \rightarrow$$

Sum of squares of Residual

$$\underline{SS_{\text{total}}} : \sum (A - C)^2$$

$$: (2 - 3.5)^2 + (5 - 3.5)^2 + (3 - 3.5)^2 + (4 - 3.5)^2$$

$$: 5.0$$

$$SS_{\text{res}} : \sum (A - B)^2$$

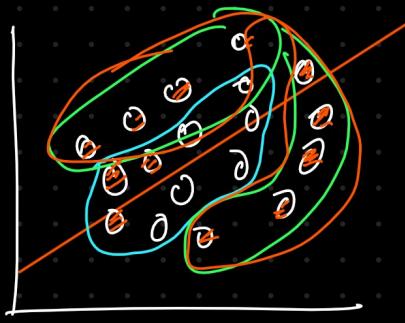
$$: (2 - 2.5)^2 + (5 - 4.0)^2 + (3 - 4.0)^2 + (4 - 3.5)^2$$

$$: 2.5$$

$$R^2 = 1 - \frac{2.5}{5.0}$$

$$= 1 - \frac{1}{2} = 1 - 0.5 = \underline{\underline{0.5}}$$

means our model explain 50% of the variance of the data.



50%

$$\underline{\underline{R^2 = 0.5}}$$

$$\underline{\underline{t^2 = 0.1}}$$

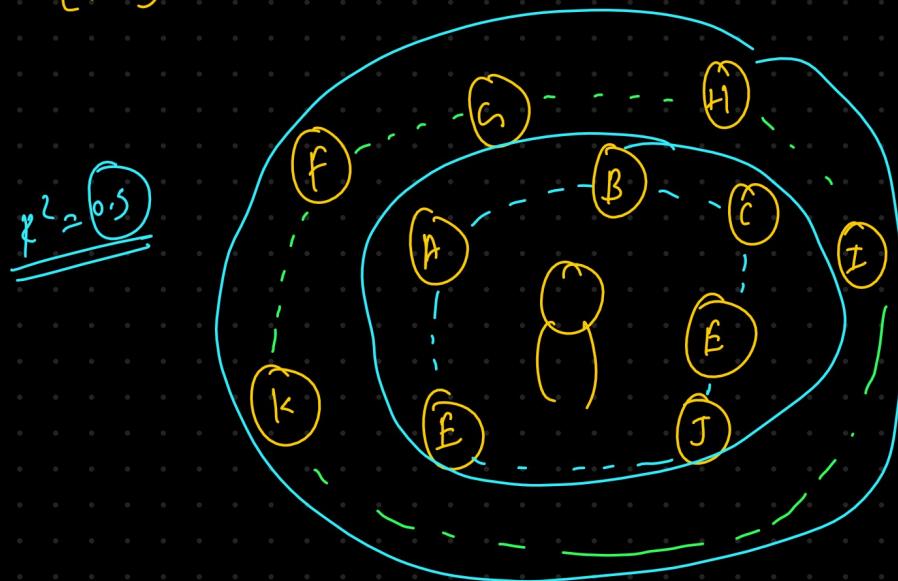
A B C D E F G H Target-variable

$A \rightarrow T \rightarrow$

$B \rightarrow T$

$C \rightarrow D$

$\rightarrow p^2 \uparrow$



10 day

Sol. ✅

Sol. ✗

A, B, C, D, E

```

def forward_selection(X, y, threshold=0.05):
    selected_features = []
    remaining_features = list(X.columns)
    while remaining_features:
        pvals = {}
        # Keep adding all the features one by one and get the p_value of each selected features
        for feature in remaining_features:
            model = sm.OLS(y, sm.add_constant(X[selected_features] + [feature])).fit()
            pvals[feature] = model.pvalues[feature]
        min_pval_feature = min(pvals, key=pvals.get)
        min_pval = pvals[min_pval_feature]
        if min_pval < threshold:
            selected_features.append(min_pval_feature) # add the feature we selected to select_features
            remaining_features.remove(min_pval_feature) # remove the feature we selected from X
            print(f"Selected {min_pval_feature} with p-value {min_pval:.4f}")
        else:
            break
    return selected_features
  
```

$A, C, B \rightarrow 0.1$
 $A, C, D \rightarrow 0.2$
 $A, C, E \rightarrow 0.3$
 $0.1 < 0.05$

(A)

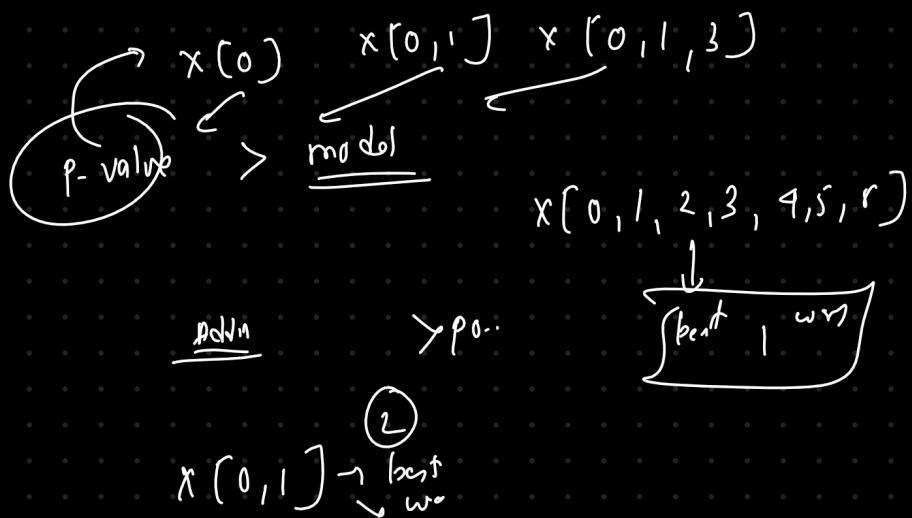
$0.01 < 0.05$

$model \rightarrow A \rightarrow$

[A, C]

$\downarrow \rightarrow model$
 $const + x[0]$

$$\underline{\underline{P\text{-value} > 0.05}}$$



RFE (Recursive Feature Elimination)

↓
method to select features

↓
algo → linear - regression
→ Random forest

$$y = m_0 + m_1 x_1 + m_2 x_2 + m_3 x_3$$

Type_of_car	Engine_cc	Car_color	price	Type of car >
sedan	250	red	y_1	Engine_cc >
Hatchback	200	white	y_2	car_color
sedan	300	blue	y_3	
Sport car	600	white	y_4	

$$y = m_0 + m_1 x_1 + m_2 x_2 + m_3 x_3$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

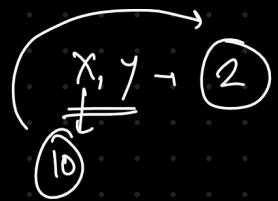
$m_1 \rightarrow x_1$
 $m_2 \rightarrow x_2$
 $m_3 \rightarrow x_3$

check the coefficient of features: $m_1 \rightarrow x_1$

$$m_2 \rightarrow x_2$$

$$m_3 \rightarrow x_3$$

RFE (algo, imp - features - to - select = 2)



RFE : -

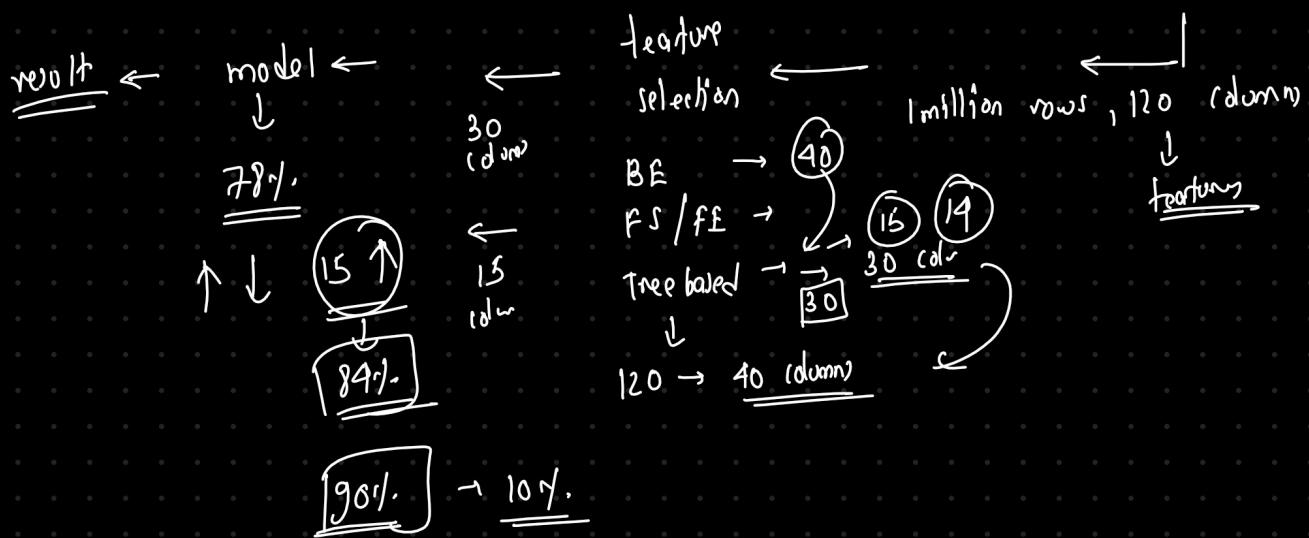
Tree based method!

linear - regression

- linear
- non-linear

Data → Data Analysis → Data cleaning → Data handling ↴

Dartq
encoding



0.7 to 1.0 or -0.7 to $-1.0 \rightarrow$ highly correlated

Feature \rightarrow Target column

LotArea \rightarrow SalePrice \rightarrow keep the feature

Feature \rightarrow Feature

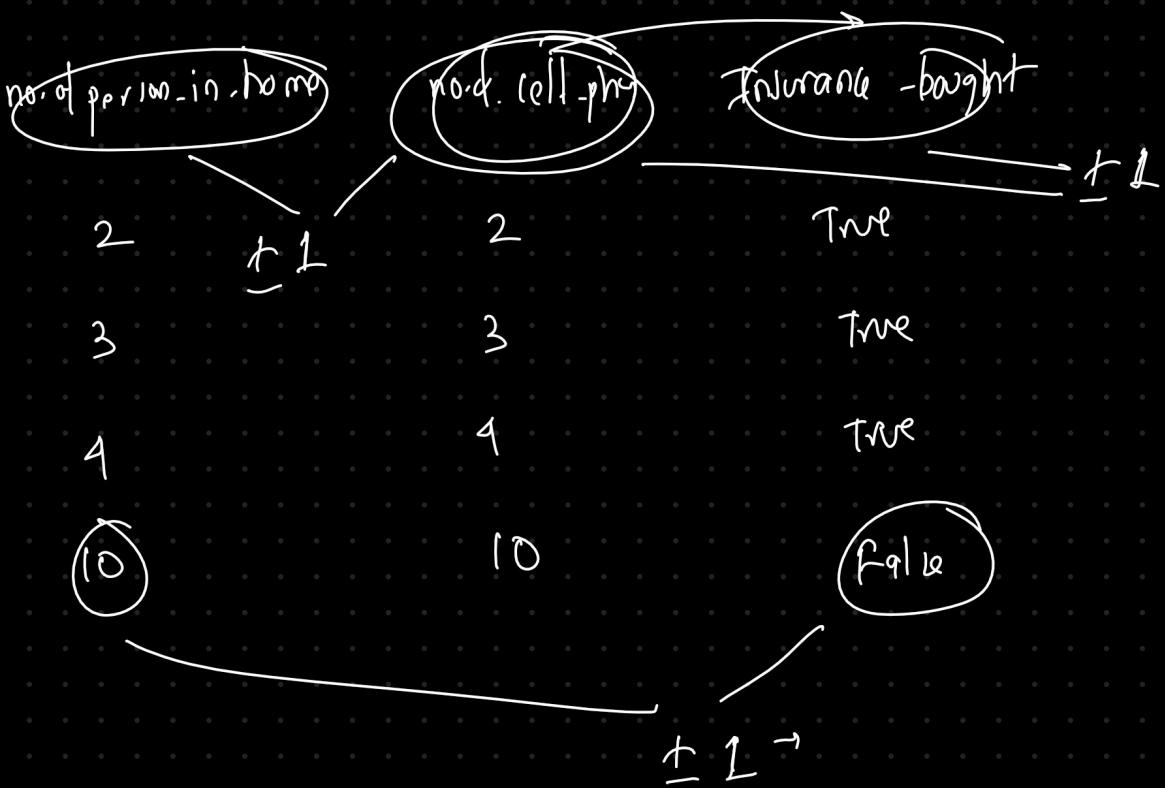
LotArea \rightarrow GrLivArea \rightarrow keep only one

0.5 to 0.7 or -0.5 to -0.7

usually fine to keep

0.0 to 0.3 or -0.3 to 0.0

weak correlation (little linear relation)



VIF: Variance Inflation Factor

→ detect multicollinearity in regression feature

→ feature → correlated ← feature)

$$VIF = \frac{1}{1 - R^2}$$

VIF: 1 : No multicollinearity

: 1-5 : Moderate correlation

: >10 : Strong multicollinearity

			Y
x_1 : A <u>bedroom</u>	x_2 : B <u>no of rooms</u>	→ C	
2	2		1200 \$
3	3		1600 \$
4	4		1700 \$
5	5		1900 \$

A, B, C, D, E

(A, B) → model
P > 0.05

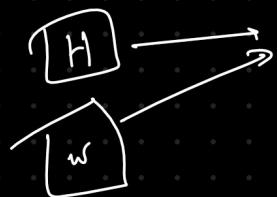
A → model

B → model

X → C → D → E → model
P > 0.05

Heg Int: $P > 0.05$ BM |
 weight $\rightarrow P > 0.05$ BB |

(H, w) $\rightarrow \beta^m |$
 $\underline{0.03}$



$\text{AT} \rightarrow \text{BE} \rightarrow \text{g}^{at}$
 $\rightarrow \text{VF}$

$\rightarrow \text{BE} \cup \text{FE} \rightarrow 92\%$

$\rightarrow \text{VF}$

$\rightarrow \text{BE}_{\text{current}} \& \text{VF} \rightarrow 90\%$