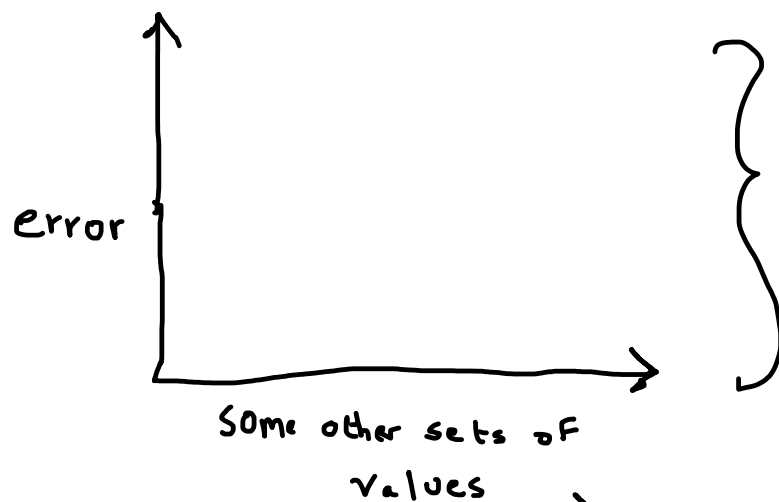


	y_actual	y_pred	error
0	50.000000	50.252809	0.252809
1	51.010101	53.509858	2.499757
2	52.020202	47.040657	-4.979545
3	53.030303	56.498296	3.467993
4	54.040404	51.948896	-2.091508

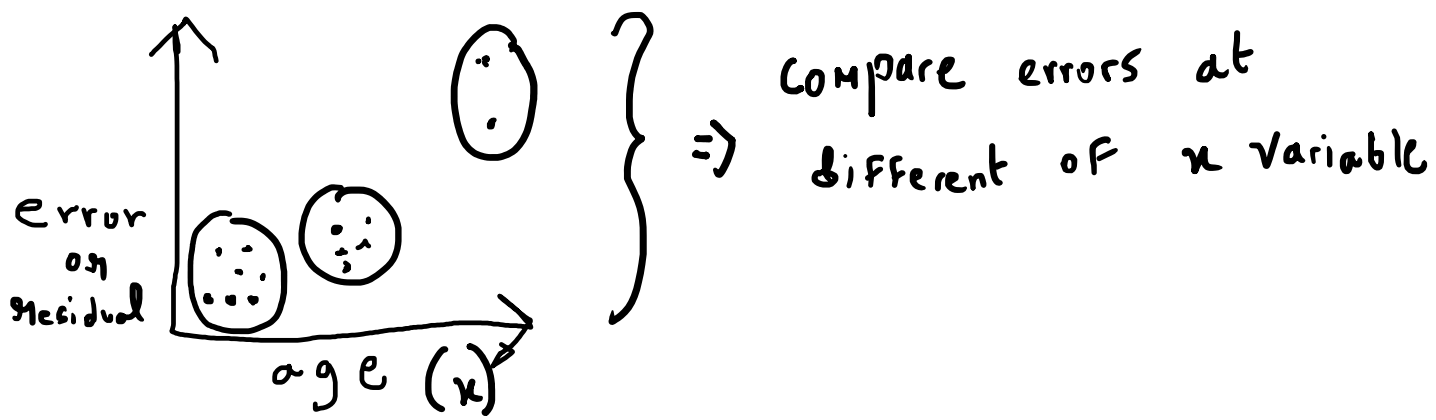
$\rightarrow y_{pred} - y_{actual}$
 \Downarrow
 Residual

The importance of error based plots!

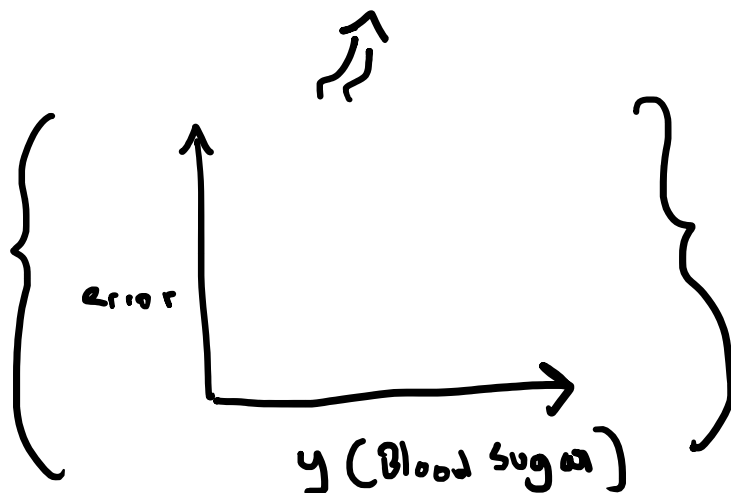


Can be input variables or target variable

Assume age is a feature

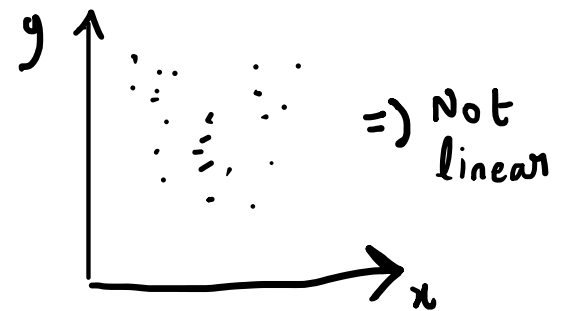
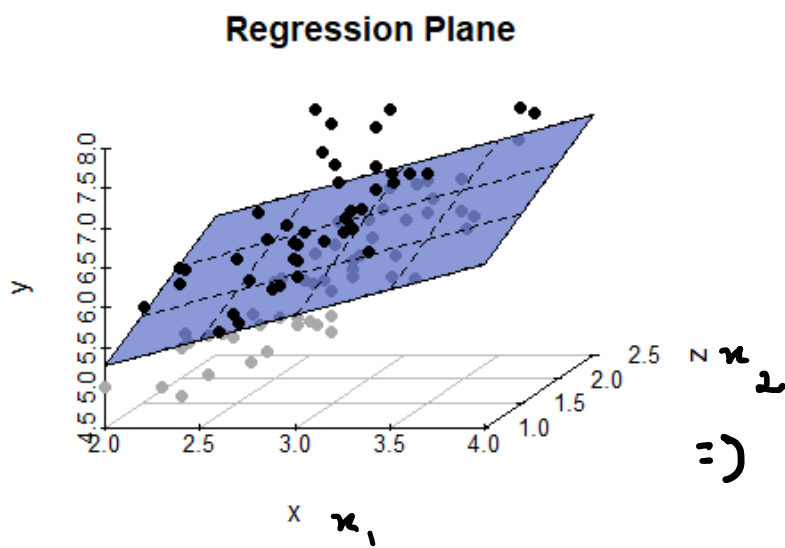
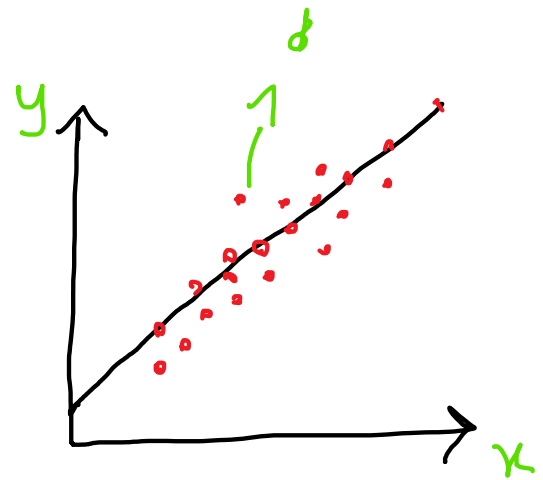
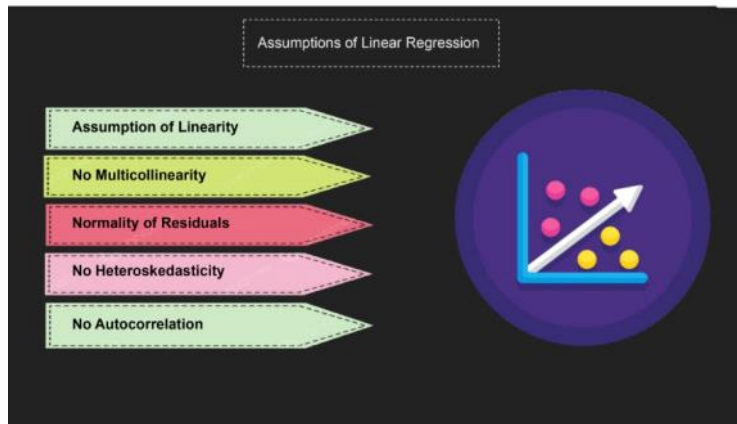


Throughout the session, we will look into importance of error vs y variable plots

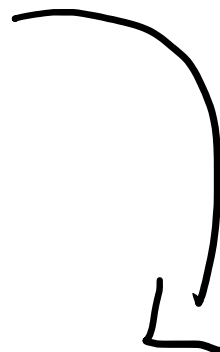
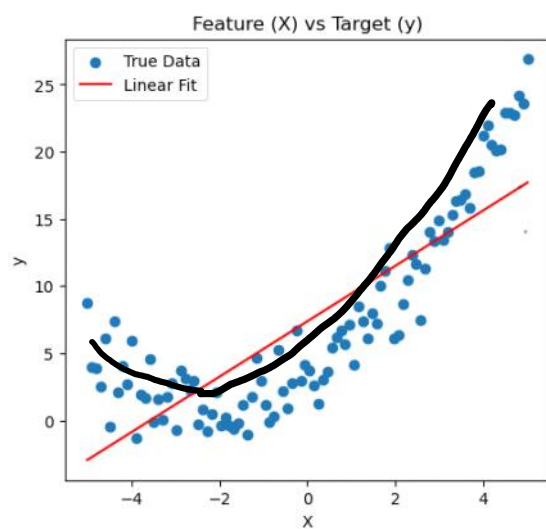


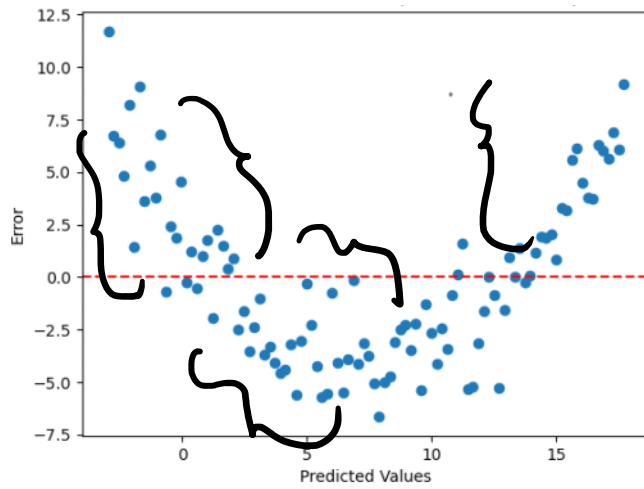
Assumption of Linearity

30 July 2025 16:02



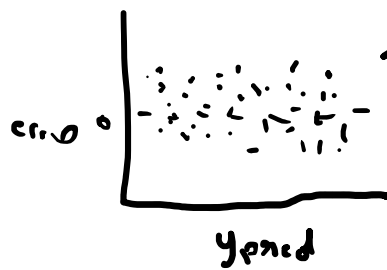
\Rightarrow hyperplane also linear





When there is a noticeable relationship b/w y_{pred} and error, linear regression is not suitable option

ideal Scenario



→ random plot } → linear regression working well

Statistical Correctness
vs
Predictive Correctness

} Sometimes a trade-off
with linear models

What is Multi-collinearity? \nearrow $\text{Travel_expense} = w_1 \times \text{Km} + w_2 \times \text{Fuel} + w_0$

Features

Km_Driven (X1)	Fuel_Consumed (X2)	target Travel_Expense (Y)
1000	80	6400
800	64	5120
1200	96	7680
600	48	3840
900	72	5760

α

\rightarrow Are X_1 and X_2 related?

$$\{x_1 = wx_2 + b\} \rightarrow \text{linearly related}$$

Why is it an issue?

\rightarrow Feature interpretability heavily impacted,
Which is more important, x_1 or x_2 ? Both?

\Uparrow

\rightarrow Sign flipping in co-efficients - Causes Confusion.

But it is not always a problem!

How can we avoid Multi-collinearity?


VIF! \rightarrow Variance Inflation Factor
 \Downarrow

How to deal with multicollinearity?

We will use **Variance Inflation Factor (VIF)**

Say, we have 'd' features $\langle f_1, f_2, f_3, \dots, f_d \rangle$

In, (VIF) we treat $\begin{cases} \text{one feature as 'y'} \\ \text{remaining features as 'x'} \end{cases}$



$f_1, f_2, f_3, \dots, f_d$	f_d
x_i	y

Handwritten labels: f_1, f_2, f_3, f_4 and f_3, f_2, f_1, f_4 are written next to the table.

$\left\{ \begin{array}{l} \text{identify} \\ \text{Multi-collinearity} \end{array} \right\}$

$x_1, x_2, x_3, x_4, \dots, y$

Now,

Train **linear regression** model with (x_i, y)

Find R^2 of the model

To Calculate VIF :

$$VIF = \frac{1}{1 - R_j^2}, R_j^2 : R^2 \text{ for } j^{\text{th}} \text{ feature}$$

target

Features

$$x_1 = w_1 x_2 + w_2 x_3 + w_3 x_4 + w_0 \Rightarrow$$

$$\hookrightarrow \frac{1}{1 - R_{sq, x_1}} \text{ For this model} = VIF \text{ for } x_1$$

$$x_2 = w_1 x_1 + w_2 x_3 + w_3 x_4 + w_0$$

$$\vdots$$

$$x_d = w_1 x_1 + w_2 x_2 + \dots + w_0$$

$1 \rightarrow R^2 = 0.9$

$\left\{ VIF = \frac{1}{1 - 0.9} = \frac{1}{0.1} = 10 \right\}$

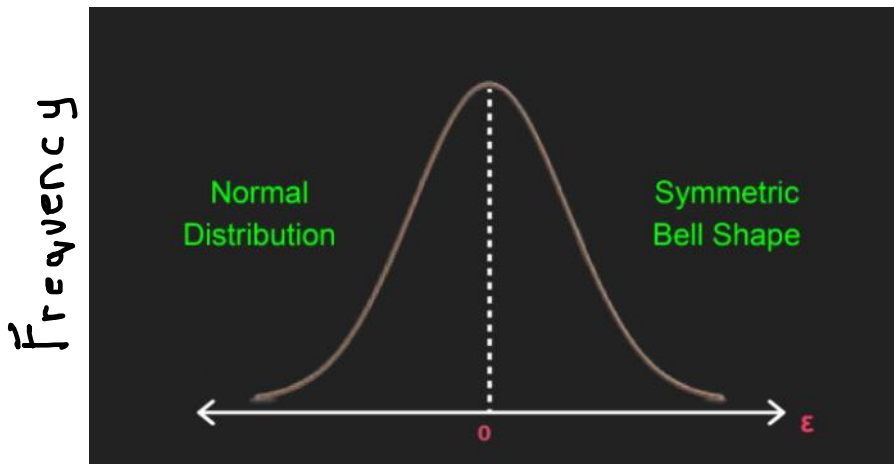
Higher the VIF For a Variable, higher the Multi-collinearity it shares with other variables.

$VIF > 5 \Rightarrow$ think of dropping

$VIF > 10 \Rightarrow \{ \text{definitely drop} \}$

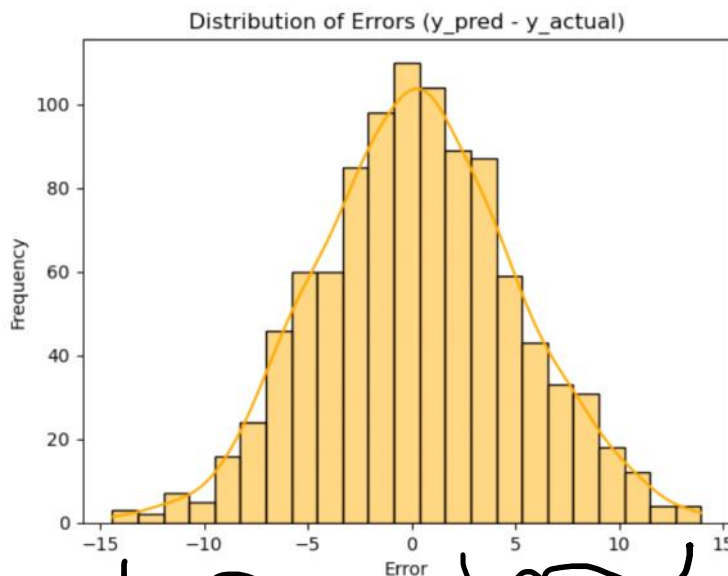
$\{ 0 \text{ to } \infty \}$

$$\text{error} = \underset{\text{Actual}}{y_i} - \underset{\text{Predicted}}{\hat{y}_i}$$



Why?

i)



Symmetric

↙

No. of overpredicted

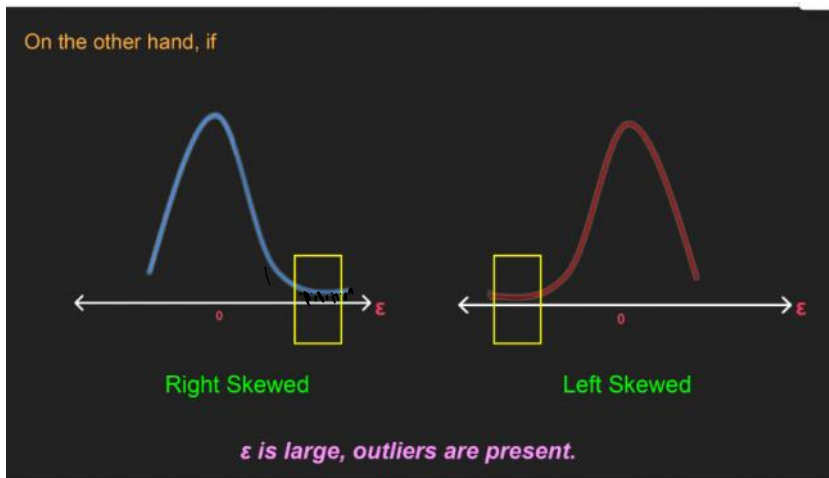
=

No. of underpredicted

Underprediction overprediction
 $\{y_{pred} - y_{actual}\}$

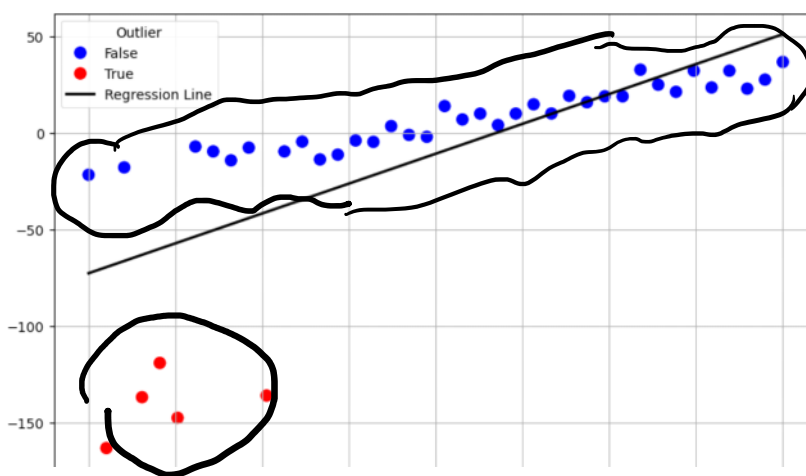
If my model is either, o.p or u.p, heavily

On the other hand, if

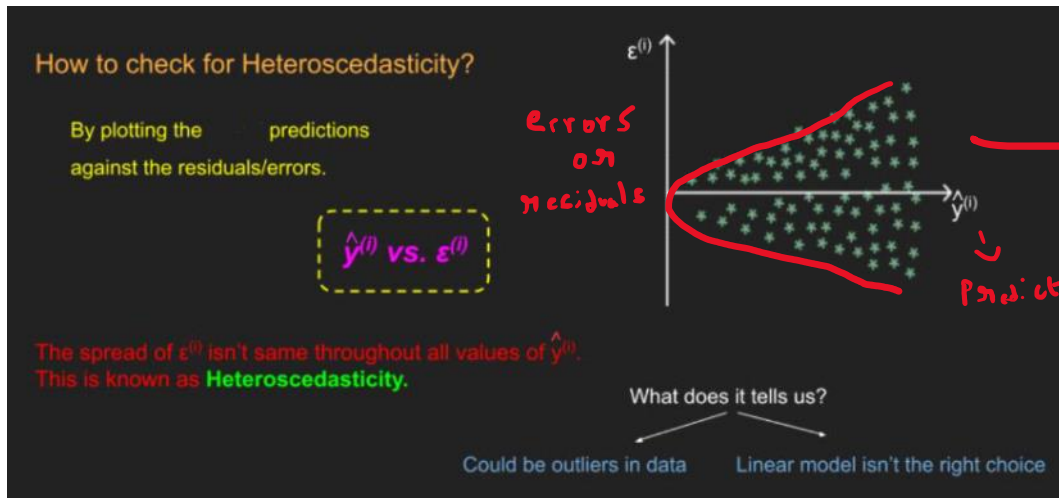


i) Model is unable to capture sufficient information

ii) There are outliers in the data



Variance of errors should
 ↗ be constant across \hat{y} -predicted



here
 it is
 not
 constant

What does Heteroscedasticity indicate?

- Linear Model may not be suitable
- Use a non-linear Model
- Introduce other variables into your model



What is self auto-correlation ?
 in errors there should be
 no auto-correlation

Note : Assumption applicable only For time-series data

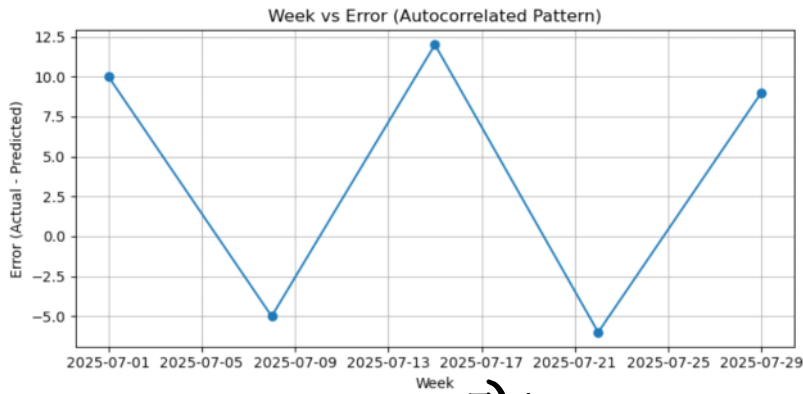
What is time-series data ?

Week	Sales (target)	Promotion	Temperature (°C)	Holiday
2025-07-01	200	1	30	0
2025-07-08	180	0	31	0
2025-07-15	220	1	29	0
2025-07-22	250	1	28	1
2025-07-29	190	0	32	0

} Predicting
 based on
 a time component
 is called
 time-series

Final Table

Week	Sales (Actual y)	Predicted Sales (y_{pred})	Error ($y - y_{pred}$)
2025-07-01	200	190	10
2025-07-08	180	185	-5
2025-07-15	220	208	12
2025-07-22	250	256	-6
2025-07-29	190	181	9



→ time component

Error

vs

time component

odd weeks errors are correlated

even week errors are correlated

How is this similar to heteroscedasticity?

trend in errors!!