

Multiple Linear Regression

$$Y = a_1x_1 + a_2x_2 + a_3x_3 + a_4x_4 + b + \varepsilon$$

(i) x_1, x_2, x_3, x_4 are independent of each other

(ii) Plot , correlation

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

if VIF calculated cross from each
20, the R^2 are
not coming from normal.

regressions

$$x_1 \sim x_2 + x_3 + x_4$$

$$x_2 \sim x_1 + x_3 + x_4$$

$$x_3 \sim x_1 + x_2 + x_4$$

$$x_4 \sim x_1 + x_2 + x_3$$

R^2

VIF = Variance inflation factor

→ If VIF calculated > 20 , that variable is dependent on
Others

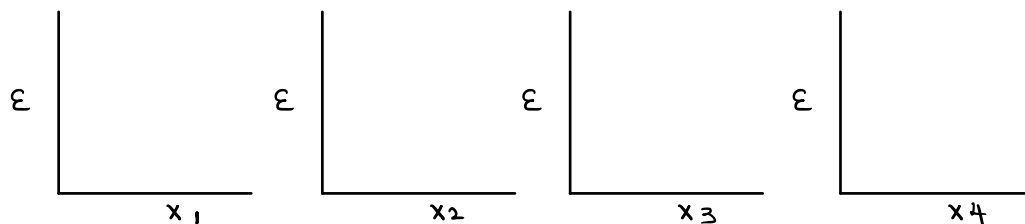
→ this dependency is called collinearity.

correlation :
$$\frac{\sum (x_i - \bar{x})(z_i - \bar{z})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (z_i - \bar{z})^2}}$$

→ If all VIF < 20

Second assumption : x_1, x_2, x_3 and $x_4 \Rightarrow$ Error is independent

calculate :



\Rightarrow There shouldn't be any pattern in these plots. That is there shouldn't be any dependency.

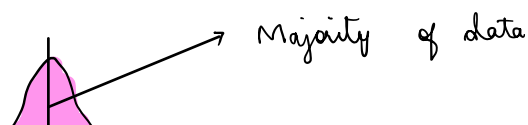
If there is dependency found then we need to perform transformation on that data to remove the dependency structure.

Assumptions

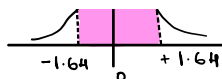
- i) We take variables are independent to remove multi-dependent-representation
- ii) the variables are independent of error

iii) Errors are coming from normal with mean zero.

\downarrow
majority of data is close to zero.



How to prove error is coming from Normal distribution



$$\begin{matrix} x_1 & y_1 & \text{pred } \hat{y}_1 \\ x_2 & y_2 & \hat{y}_2 \\ \vdots & \vdots & \vdots \\ x_n & y_n & \hat{y}_n \end{matrix} \rightarrow \epsilon_1, \dots, \epsilon_n \rightarrow N(0, \sigma^2)$$

↓
normalize

$$\left\{ \begin{array}{l} \frac{\epsilon_1 - 0}{\sigma} = \tilde{\epsilon}_1, \tilde{\epsilon}_2, \dots, \tilde{\epsilon}_n \rightarrow N(0, 1) \rightarrow \text{predicted} \\ z_1, z_2, \dots, z_n \rightarrow N(0, 1) \rightarrow \text{unpredicted} \end{array} \right.$$

Calculate the quantiles and they should be same.

Transformation of Data \Rightarrow Box-Cox transformation.

③

constant variance (Homoscedasticity)

It's a property that all calculated errors come from normal distribution and have same/constant variance.

x_1	y_1	\hat{y}_1	ϵ_1	$\rightarrow N(0, 1)$
\vdots	\vdots	\vdots	ϵ_2	$\rightarrow N(0, 1)$
\vdots	\vdots	\vdots	\vdots	\vdots
\vdots	\vdots	\vdots	\vdots	\vdots
x_n	y_n	\hat{y}_n	ϵ_n	$\rightarrow N(0, 1)$

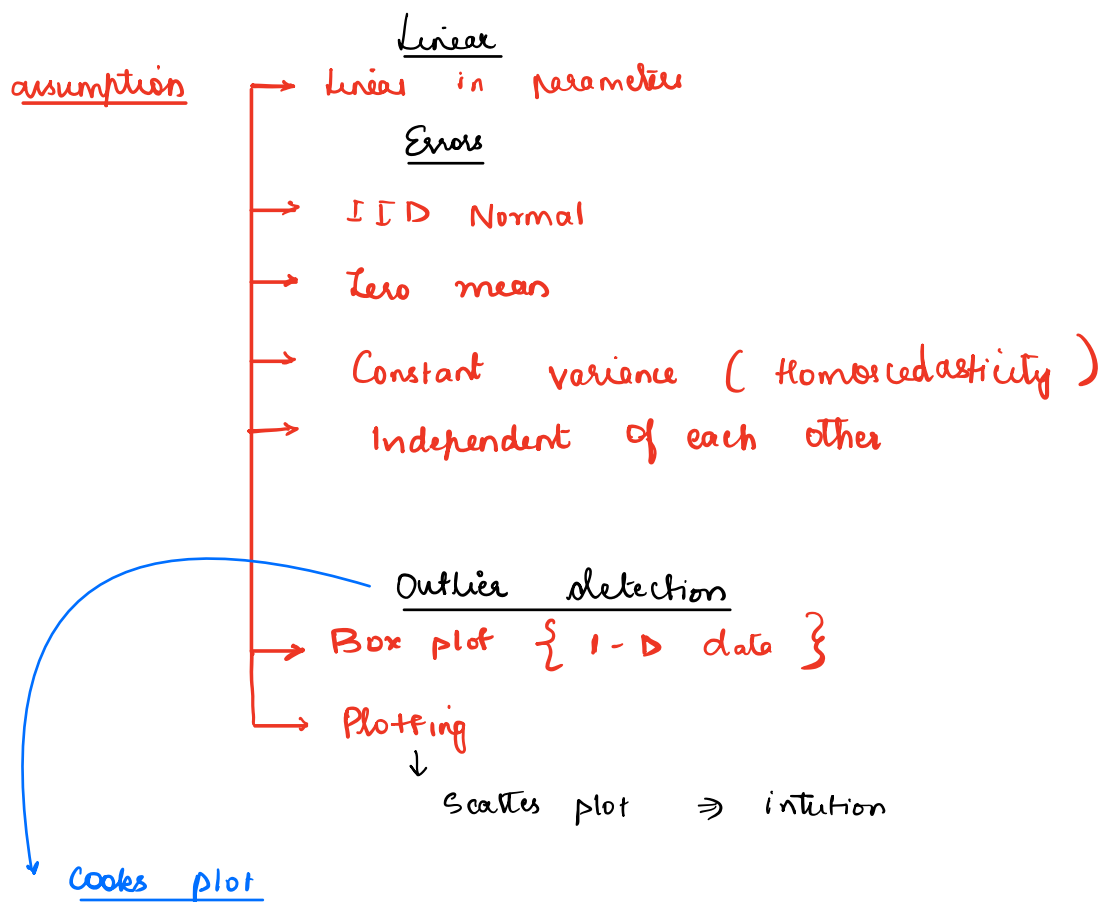
i.e. when ϵ vs \hat{y}_i is plotted, The graph should be random.



4) Partial Regression \Rightarrow Independent of each other

LR- Model ϵ Assumptions

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon$$



\rightarrow remove a row at a time a plot

\rightarrow plot all these, clearly an outlier will be visible.