# Multicollinearity

correlation  $x - y$   $\sim 0.95$

$$\boxed{X_1 \quad X_2 \quad X_3} \longrightarrow y$$

$$\Downarrow$$

$$\underline{X_1 \quad X_2 \quad X_3 \quad X_4 \quad --- \quad X_{100}}$$
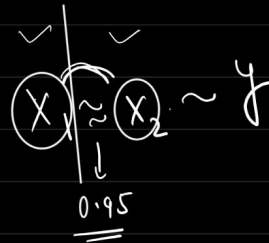
if $X_1 = X_2$

$$\left\{ \begin{array}{l} 8X_1 + 2X_2 \\ 10X_1 \\ 10X_2 \\ 2X_1 + 8X_2 \end{array} \right\}$$

$$\overset{\smile}{\cancel{X_1}} \underset{0.95}{\overset{\approx}{\phantom{x}}} \overset{\smile}{X_2} \sim y$$

$X_1 \approx X_2$

$X_1 \approx (X_2 \; X_3 \; X_4)$

$$\underline{\text{multi}} - \underline{\text{col}} - \text{linearity}$$

$$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$
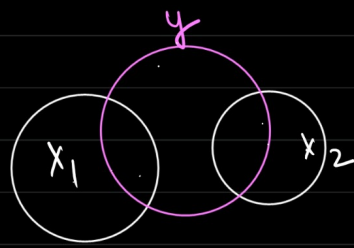
many        together      linear relationship.

$X_1 \sim X_2 \rightarrow$ Correlation

$$\left. \begin{array}{l} X_1 \approx (X_2 \; X_3) \\ X_1 \approx (X_2 \; X_3 \; X_4) \end{array} \right\} \text{multicollinearity.}$$

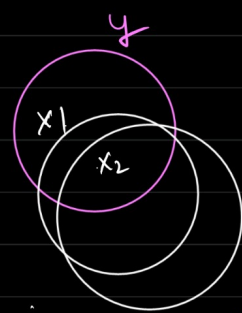$$\bigcirc\!\!\!\!X_1 \qquad X_2 \quad X_3 \; X_4 \; X_5 \qquad y$$

Collinearity — where two features are linearly associate (high correlation) and they are used to predict target variable

multi - Collinearity — where a feature exhibits a linear relationship with more than two variable.

No multicollinearity. $= 0$

Multicollinearity

* Concern :-

→ It increases Overfitting

→ Affects interpretation

$X_1 = X_2$   $8X_1 + 2X_2 = y$

$X_1 \approx X_2 \, X_3 \qquad y$

* Soln

① VIF and drop feature one by one with high VIF

VIF - Variance Inflation factor.

② RFE — Recursive feature elimination.

$X_1 - X_2 \rightarrow$ correlation / heatmap

$X \sim X_2 X_3 X_4 \rightarrow$ VIF

* VIF is a measure of amount of multicollinearity in regression.

$X_1 \; X_2 \; X_3 \; y$

$X \approx X_2 X_3$

$R_i^2$ — %age variation in Y explained by X.

$$VIF_i = \frac{1}{1 - R_i^2}$$

$X_1 \; X_2 \; X_3 \; X_4 \; X_5$

$VIF_{x_1} \longrightarrow x_1 \simeq (x_2 \, x_3 \, x_4 \, x_5)$    $R square$

$\overset{\downarrow}{\textcircled{y}} \longrightarrow x \longrightarrow \underline{R_i^2}$   $VIF = \dfrac{1}{1-R_{x_1}^2}$

$X_1 X_2 X_3 \quad y$

$Rsquare = R^2 = $ %age variate explained in Y by X

$\checkmark X_1 \sim \boxed{X_2 \, X_3} \xrightarrow{\hspace{3cm}} VIF_{x_1} \sim \cancel{x_i} \text{ and } \boxed{x_2 | x_3} \rightarrow R square$

$\checkmark X_2 \approx \boxed{X_1 \, X_3}$      $\overset{\downarrow}{y} \quad \cancel{\times}$

$\checkmark X_3 \sim \boxed{X_1 \, X_2}$      $VIF = 1/1 - R_i^2$



NO
multicolling.          low multicolinearity        High multicollinery        High
                                                                              multicollinea

$\rightarrow \quad VIF \geq 10$

When features have VIF $\geq 10$ then
drop the feature one by
one

$\rightarrow 10 = \dfrac{1}{\searrow 1 - Rsquar}$

$1 - Rsquar = \dfrac{1}{10}$

$Rsquar = 1 - \dfrac{1}{10} \quad \Rightarrow \quad \dfrac{10-1}{10} = \dfrac{9}{10} = 0.9$

$\begin{cases} X \sim y \\ X_1 \sim (X_2 X_3 X_4) \\ \quad VIF \\ \quad \downarrow \\ \underset{(y)}{\boxed{X_1}} \sim X_2 X_3 X_4 \\ \quad (X) \end{cases}$

$\boxed{X_1} \sim X_2 X_3 X_4$

90% variance in $X_1$ is explained $X_2 X_3 X_4$

| Feature | VIF |
|---------|-----|
| $X_1$ | 12 |
| ✓ $X_2$ | 13 |
| $X_3$ | 8 |
| $X_4$ | 7 |

$VIF \geq 10 \implies$ drop the feature one by one.

* **Two kinds of multicollinearity**

✓ ① Data - based collinearity. — Present in data itself.

ex: latitude | longitude

✓ ② structural multicollinearity — Caused due to new feature from existing feature.

Distance | time | Speed = $\dfrac{Distance}{time}$

* $\underline{1000 \rightarrow \text{features}}$
   ↓

→ R F E — Recursive Feature Elimination

  — It will make a model with all 1000 features
  — Start dropping one by one least important feature.
  → Untill the desired no of feature is achieved.

→ PCA