# What exactly happens in Multicollinearity ?

When multicollinearity is present in a model, it can lead to several issues, including:

1. **Difficulty in identifying the most important predictors:** Due to high correlation between independent variables, it becomes challenging to determine which variable has the most significant impact on the dependent variable.

2. **Inflated Standard errors:** Multicollinearity can lead to larger standard errors for the regression coefficient, which decrease the statistical power and can make it challenging to determine the true relationship between the Independent and dependent variables.

3. **Unstable and unreliable estimates:** The regression coefficient becomes sensitive to small changes in the data, making it difficult to interpret the results accurately.

# Perfect Multicollinearity

Perfect multicollinearity occurs when one independent variable in a multiple regression model is an exact linear combination of one or more other independent variables. In other words, there is an exact linear relationship between the independent variable, making it impossible to uniquely estimate the individual effects of each variable on the dependent variable.

$$\boxed{percent = 10 \times cgpa + error}$$

## Corr linear ↗

$$\boxed{x_1 = a_1 x_2 + a_0 + error}$$

→ **example** ↙ , add diciplne marks

| Cgpa | percent | Lpa |
|------|---------|-----|
| 8.5  | 83      | -   |
| 9.12 | 95      | -   |

## Perfect collinearity ↗

$$\boxed{x_1 = a_1 x_2 + a_0}$$

→ example ↗

| Cgpa | percent | Lpa |
|------|---------|-----|
| 8.5  | 8.5     | -   |
| 9.12 | 91.2    | -   |

$$\boxed{percut = 10 \times Gpa}$$

# Example

| Cgpa | percent | lpa |
|------|---------|-----|
| 8    | 80      | 3   |
| 6    | 60      | 4   |

$$\boxed{lpa = \beta_0 + \beta_1 Cgpa + \beta_2 percut + error}$$

OLS ⌐

$$\beta = (X^T X)^{-1} X^T Y$$

→ design Matrix

$$\begin{bmatrix} 1 & 8 & 80 \\ 1 & 6 & 60 \end{bmatrix}$$

input data

3×2                    2×3

$$\begin{bmatrix} 1 & 1 \\ 8 & 6 \\ 80 & 60 \end{bmatrix}$$ $$\begin{bmatrix} 1 & 8 & 80 \\ 1 & 6 & 60 \end{bmatrix} = \begin{bmatrix} 2 & 14 & 140 \\ 14 & 84 & 8400 \\ 140 & 8400 & 840000 \end{bmatrix}$$

$X^T$                                              $X^T X$

\* before findly Inverse Check Determinate (det)



Det = 2(0) − 14(0) + 140(0) = 0

if det is zero then it is singular matric and we
cant find inverse. and also |β| not find.

## Proofe

\# if our data has multicollinearity
→ unstable coefficient } Problem
→ High S·E

# formula

$$\beta = (X^T X)^{-1} X^T Y$$

$$\text{Var}(\beta) \text{ or } SE(\beta)$$

Underroot of var(β)

$\beta_0 \quad \beta_1 \quad \beta_2$

$SE(\beta_0) \quad SE(\beta_1) \quad SE(\beta_2)$

$$SE(\beta) = \sqrt{\text{diag}(\sigma^2 (X^T X)^{-1}}$$

not ~~~~
→ inverse

$$\beta = (X^T X)^{-1} X^T Y$$

→ perfect → det $(X^T X) = 0$
multicoli

→ Stoony multicoli
det $(X^T X)$ → Very small

inverse = $\dfrac{1}{\det}$ 

very small

inflate → high

# How to Remove multicollinearity

1. **Collect more data**: In some cases, multicollinearity might be a result of a limited sample size. Collecting more data, if possible, can help reduce multicollinearity and improve the stability of the model.

2. **Remove One of the highly correlated variable**: If two or more are highly correlated, consider removing one of them from the model. This step can help eliminate redundancy in the model and reduce multicollinearity. Choose the variable to remove based on domain knowledge, variable importance, or the one with the highest VIF.

3. **Combine Correlated Variable**: If correlated Independent variable represent similar information, consider combining them into a single variable. This combination can be done by average, summing or any other mathematical operations, depending on the content and the nature of the variables—

4. <u>Use partial least square regression (PLS)</u>: PLS is a technique that combines features of both principal component analysis and multiple regression. It Identifies linear combinations of the predictor variables (called latent variable) that have the highest covariance with the response variable, reducing multicollinearity while retaining most of the predictive power.

* First find PCA → (then) → LR