# Multi collinearity

Multi collinearity is a statistical phenomenon that <mark>Multicollinear</mark> variable occur when two or more independent variable in a multiple ~~linear~~ regression model are highly correlated. In other words, these variables exhibit a strong linear relationship, making it difficult to isolate the individual effects of each variable on the dependent variable.

eg:- ①

| cgpa | iq | lpa |
|------|-----|-----|
| 8 | 80 | 8 |

$$iq \uparrow \rightarrow cgpa \uparrow$$

| Corr | → linear regression |

②

| iq | backlogs | lpa |
|----|----------|-----|
| ↑↑ | ↓ ↓ | |

└─corr.

③

| cgpa | dob | lpa |
|------|-----|-----|
| | | |

not corr

$$\overset{\nearrow ??}{\underline{lpa}} = \beta_0 + \beta_1 \overset{\uparrow}{\underline{cgpa}} + \beta_2 \underline{iq}^{\curvearrowright c}$$

How its work?

① let iq is constant

② and if cgpa is increase from 1.

③ Then how much lpa increase.

But in Multicollinearity how its work?

① let assume iq is constant

⑤ and cgpa increase from 1.

③ If cgpa and iq is multicollinearity then iq is also increase from 1 (no constant)

④ then this formula is not working

* When is Multicollinearity bad?

1. Inference:

* Inference focus on understanding the relationship bet^n variable in a model.

eg:- cgpa | iq) lpa

* How much cgpa and iq perform sole to give output data (lpa)

Let assume

| cgpa | iq | lpa |
|------|-----|-----|
| perform 75% | perform 25% | to give output |

(in game)

eg:-

| User activity | ban / unban |
|---------------|-------------|

bullet shots     Running     guns

                                    Hack
if user use auto bulletshot ⟹ Ban
           perform 75%

if user not using Hack ⟹ Unban

# Prediction

* Prediction focuses on using a model to make accurate forecasts or estimate for new, unseen data.

    eg:- learn from old or train data set to predict new unseen data.

---

**✓** ** Multi collinearity doesn't affect the model when you are building a predictive model. But if you are using for inference (find the relationship betn input and output) then multi - collinearity is bad.

$x_1 | x_2 | Y$

## Explaining

In Prediction $\rightarrow$ $\boxed{x_1 = a_0 + a_1 x_2 + \lambda}$ $\rightarrow$ $\boxed{\text{reduced Error}}$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + t$$

$$Y = \beta_0 + \beta_1(a_0 + a_1 x_2 + \lambda) + t$$

$$Y = \beta_0 + \beta_1 a_0 + \beta_1 a_1 x_2 + \beta_1 \lambda + t$$

$$Y = (\beta_0 + \beta_1 a_0) + \underset{x_2}{\beta_1 a_1 x_2} + (\underbrace{\beta_1 \lambda + t})$$

$\underset{\text{constant}}{\swarrow}$   $\qquad \qquad \qquad \longrightarrow$ some erro

$\boxed{x_1 \text{ not in new form}}$

* which means automatically convert into $x_2$ in multicollinearity in prediction

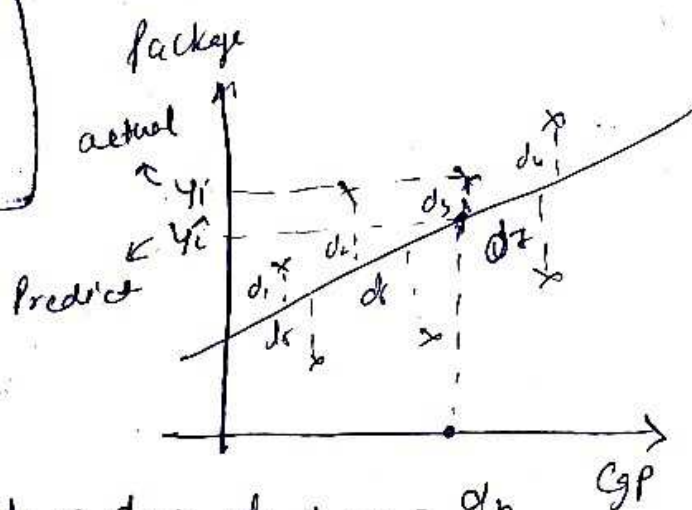# Simple Linear Regression

direct formula

↓

Closed form

$\boxed{\text{OLS}}$ ↙

$$\boxed{m = \dfrac{\sum\limits_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sum\limits_{i=1}^{m} (x_i - \bar{x})^2}}$$

$$\boxed{b = \bar{Y} - m\bar{x}}$$

Non-closed form

$$\boxed{\begin{array}{c}\text{Gradient} \\ \text{Descent}\end{array}}$$

package

actual

$\uparrow y_i$

$\leftarrow \hat{y_i}$

Predict

$d_1$ $d_2$ $d$ $d_3$ $d_4$

CgP

$$E = d_1 + d_2 + d_3 + \cdots d_n$$
$$E = d_1^2 + d_2^2 + d_3^2 + \cdots - d_n^2$$

$$\boxed{E = \sum_{i=1}^{n} d_i^2}$$

$$d_i = (y_i - \hat{y_i})$$

$$\boxed{E = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y_i})^2} \rightarrow \text{Avg Error}$$

$$\boxed{E = \sum_{i=1}^{n} (y_i - \hat{y_i})^2} \rightarrow \text{Total Error}$$

$$\hat{y_i} = mx_i + b$$

$$E(m, b) = \sum_{i=1}^{n} (y_i - mx_i - b)^2$$

## Types of Multicollinearity

1. Structural Multicollinearity :- Multi collinearity create by data scientist.

eg:① Apply One Hot encoding

category

| City |

apply One hot encoding

Delhi
Mumbai
Kolkata

| | Delhi | Mumbai | Kolkata |
|---|---|---|---|
| | 1 | 0 | 0 |
| | 0 | 1 | 0 |
| | 0 | 0 | 1 |

↗ { Perfect Multicollinearity }

$$K = 1 - Delhi - Mumbai$$
$$K = 1 - 1 - 0$$
$$\boxed{K = 0}$$
↳ Kolkata

↳ $\beta \Rightarrow$ we cannot find $\beta$ because Perfect Multicollinearity
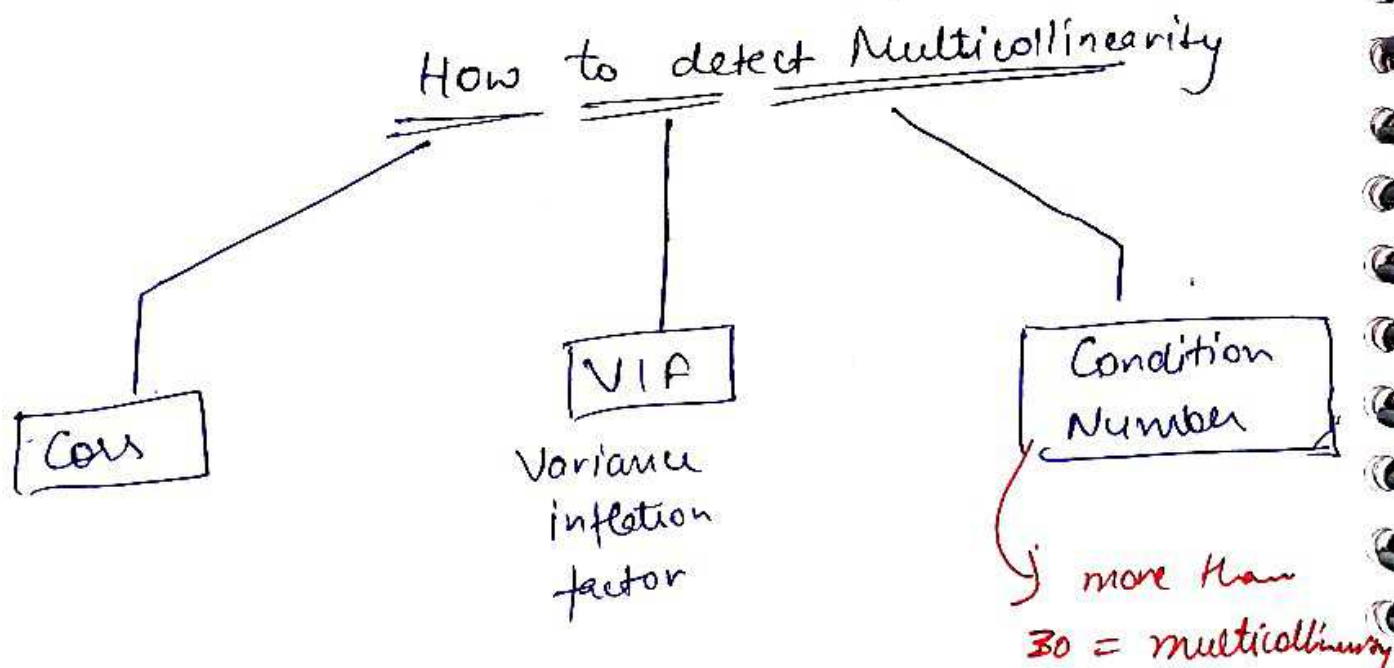
② If we apply polynomial regression

② **Data driven multicollinearity :-** Data - driven multi-collinearity occur when the independent variables in the dataset are highly correlated due to the specific data ~~due~~ being analysed. On this case, the high correlation between the variables is not a result of the way the variables are defined or the model is constructed but rather due to the observed data patterns.

   flat data

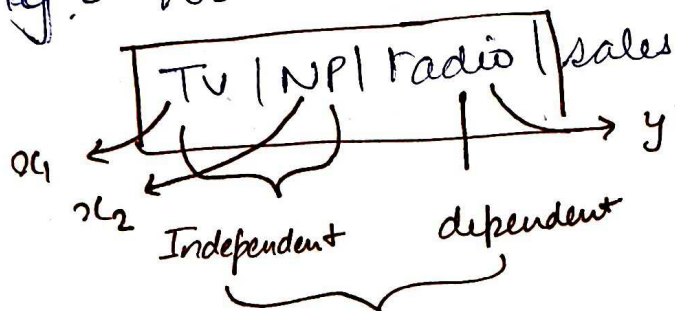| Size | no. of washroom |
|------|-----------------|

↳ data-driven multicollinearity

How to detect Multicollinearity

| Cors |

| VIF |
Variance
inflation
factor

| Condition Number |
↳ more than
30 = multicollinearity

## ① Correlation

↳ using [corr] function

## ② Variance Inflation factor

eg:- We have dataset

| TV | NP | radio | sales |

$x_1$ ← ... $x_2$ ... Independent ... dependent

→ y

assume

y (dependent)

| TV | NP | radio | sales |

$x_1$ ← ... Independent ... → $x_2$

↓ dependent

| TV | NP | radio | sales |

↓ y ... Independent → $x_2$, $x_1$

↓ dependent

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

If linear relation have
then linear regression
bta dega.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

3 input

↓

③ LR (Linear Regression)

↳ R2 Store

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

if VIF value (greater than 5 or 10) depending on the context) then multicollinearity exist. if VIF value is near 1 then multicollinear-ity doesn't exist.

## Code

```
from statsmodels.stats.outliers_influence
              import variance_inflation_factor

vif = []

for i in range(3)
     vif.append (variance_inflation_factor(df.iloc
                                       [:,1:4],i))
```

```
pd.Dataframe ({ 'vif' : vif }, index = df.columns [1:4]).T
```

## 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 highly correlated variable.

③ Combine correlated variable

④ use Partial least squares regression

$\quad\quad$ ↳ PCA → LR

$\quad\quad$ $\boxed{x_1 \quad x_2}$ → $\boxed{x_1' \quad x_2'}$

$\quad\quad\quad\quad\quad$ PCA