# Linear Regression - 04

# Assumptions of Linear Regression

→ Assumption of Linearity

→ No Multi-Colinearity

→ Normality of Residuals $(y - \hat{y})$

→ No Heteroskadasticity

→ No Autocorrelation

# No Multicolinearity

Colinearity ?

$f_1$, $f_2$

if $\boxed{f_2 = \alpha f_1 + \beta}$

$f_1$ & $f_2$ are colinear

age, year

age = 2024 − year

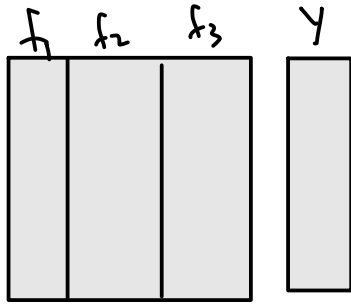Colinearity multiple features = Multi-Colinearity

$f_1$ $f_2$ $f_3$ $f_4$

$$f_2 = \alpha_1 \cdot f_1 + \alpha_3 \cdot f_3 + \alpha_4 \cdot f_4 + \alpha_0$$

$f_2$ is MultiColinearity

Height (cms) | Height (ft.)

No MultiColinearity !!

Q⇒ Why MC is a Problem?

|       | $f_1$ | $f_2$ | $f_3$ |   | $y$ |
|-------|-------|-------|-------|---|-----|

M.C exists
if $x_2 = 1.5 x_1$

Train / Optimized →

$$W^* = [w_1, w_2, w_3], w_0$$

$$w^* = [1, 2, 3], w_0 = 5$$

$x_q = [x_1, x_2, x_3]$

①: $\hat{y} = 1.x_1 + 2.x_2 + 3.x_3 + 5 \longrightarrow w = [1, 2, 3]$

$$\hat{y} = 1.x_1 + 2.(1.5x_1) + 3.x_3 + 5$$

$$\hat{y} = 1 \cdot x_1 + 3 x_1 + 3 x_3 + 5$$

②: $\hat{y} = 4 x_1 + 3 x_3 + 5 \qquad \longrightarrow \omega = [4, 0, 3]$

$x_q = [2, 3, 1]$

① $1 \times 2 + 2 \times 3 + 3 \times 1$

$\hat{y} = 2 + 6 + 3$

$\hat{y} = 11$

② $\hat{y} = 4 \times 2 + 3 \times 0 + 3 \times 1$

$\hat{y} = 8 + 0 + 3$

$\hat{y} = 11$

→ No feature importances

→ No interpretability
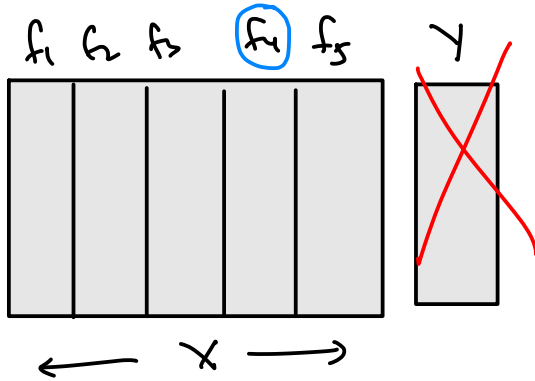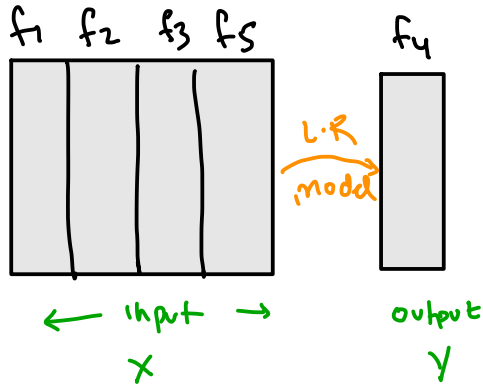
Messed up with weights → Unstability of weights

# VIF [Variance Inflation factor]

$\widehat{f_j}$
$\longleftarrow \omega_j \longrightarrow$

$f_1 \ f_2 \ f_3 \ \boxed{f_4} \ f_5$    $y$

$\longleftarrow X \longrightarrow$

## VIF($f_4$)

$f_1 \ f_2 \ f_3 \ f_5$    $f_4$

L.R
model

$\longleftarrow$ input $\longrightarrow$    output
X    y

$$\hat{f_4} = \omega_1 f_1 + \omega_2 f_2 + \omega_3 f_3 + \omega_5 f_5 + \omega_0$$

$R_2$ Score

if $R_2 = 0.98$ : M.C exists

if $R_2 = 0.18$ : M.C Not exists

$$VIF_j = \frac{1}{1 - R^2_j}$$

if $R_2 \to 1$      $VIF \to \infty$

$R_2 \to 0$      $vif \to 1$

$VIF \to [1, \infty)$

No. M.C      v. large M.C

THUMB RULE :

$Vif > 10$ : v. high m.c

$Vif > 5$ : High m.c

$vif < 5$ : low m.c

$VIF(f_1) = \underline{\quad}$

$VIF(f_2) = \underline{\quad}$

$VIF(f_d) = \underline{\quad}$

Independent features

# Normality of Residuals
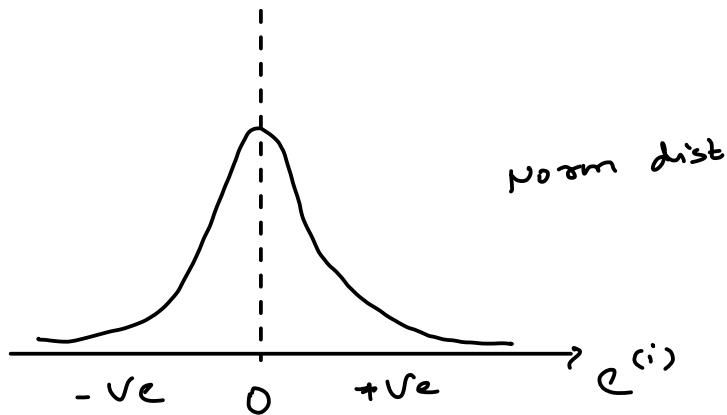
$$\rightarrow y - \hat{y}$$

$$y = \boxed{w^T x + w_0} + \epsilon$$

$$\downarrow$$

$$\hat{y}$$

$$\boxed{\epsilon = y - \hat{y}} \quad \text{residuals.}$$
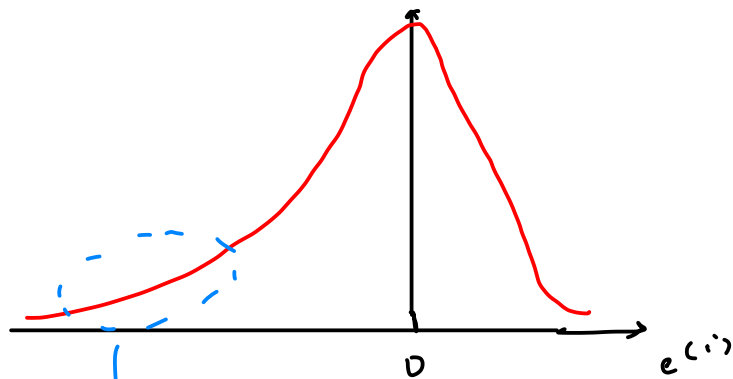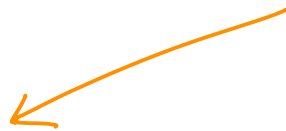
Norm dist
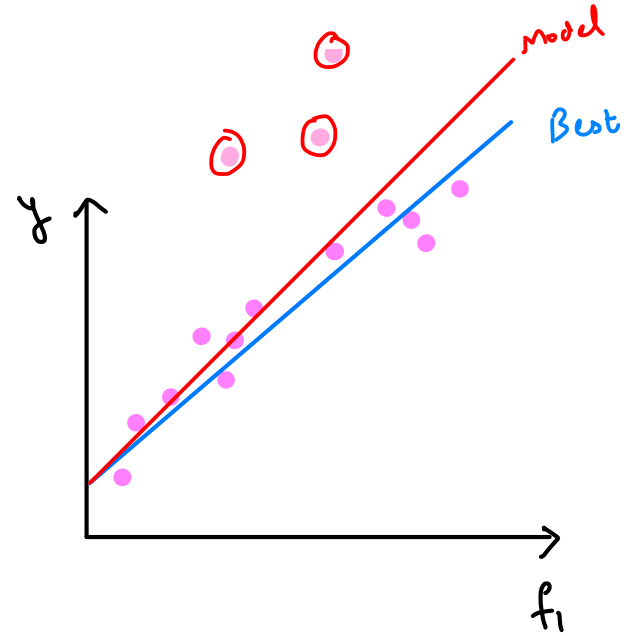
$-ve \qquad 0 \qquad +ve \qquad \rightarrow e^{(i)}$

# Impact of Outliers

I. identify outliers ?

II. Deal outliers !

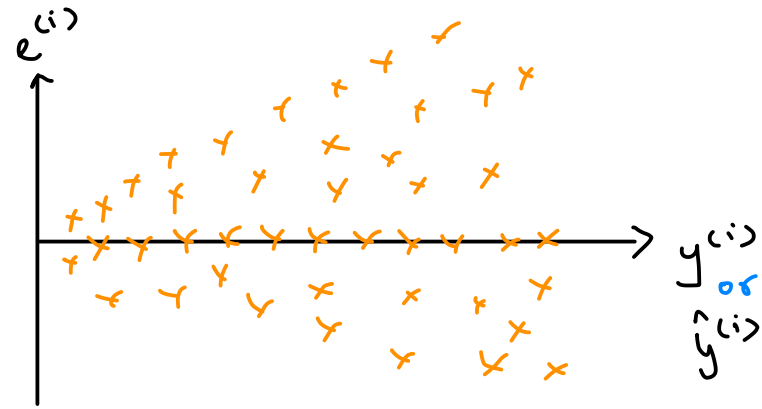**Residual Analysis**

# No Heteroskadasticity

$\epsilon^{(i)}$ vs. $y^{(i)}$

Homoskedasticity

# No Autocorrelation

→ in TimeSeries

| Time | Date | Sales |
|------|------|-------|
| — | - | |
| — | - | |
| — | - | |
| — | - | |
| — | — | □ |
| - | - | |
| - | , | □ |

Correlation with dataPoints itself

→ every dataPoint should be independent

# Quiz time!

🕐 Quiz Ended!

## In linear regression, a high VIF value suggests:

31 users have participated

| | | |
|---|---|---|
| **A** | Heteroskedasticity is present | 10% |
| **B** | A strong linear relationship between the independent and dependent variables. | 29% |
| **C** | The absence of outliers in the dataset. | 0% |
| ✓ **D** | Strong multicollinearity between predictor variables. | 61% |

$(y)$ e-g → Price