Date: 24/08/2024

Session 5: Regression & Classification Model

Topic- SLR, MLR, PLR, Logistic Regression, DT, RF

---

Simple Linear Regression

SLR Equation → $y = mx + c$

$m \downarrow$ Slope  $c \downarrow$ Y intercept

$y = \boxed{\beta_0} + \boxed{\beta_1} \bigotimes \quad \bigcirc X \to 20$

Dependent Variable

Independent Variable (feature)

Y-intercept

feature

Coefficient

## Multiple Linear Regression

- Uses more than one Independent Variable $(x_1, x_2, x_3)$ to predict a dependent variable

| $x_1$ | $x_2$ | $x_3$ | Y | |
|---|---|---|---|---|
| x-train | | | | y-train → Tr |
| x-test | | | | y-test → Test |

$\leftarrow x_1, x_2, x_3 \rightarrow$ Actual data

Equation for MLR

$$Y = \bigcirc\beta_0 + \bigcirc\beta_1 x_1 + \bigcirc\beta_2 x_2 + \bigcirc\beta_3 x_3$$

$\beta_0 = 0.5$

$\beta_1 = 0.7$

$\beta_2 = 2.7 \longrightarrow \bigcirc 2.7 \to$ Most imp feature

$\beta_3 = 1.8$

Predicted - Y

Predicted ← 4.3

Predicted y value

error

$y = \beta_0 + \beta_1 x$
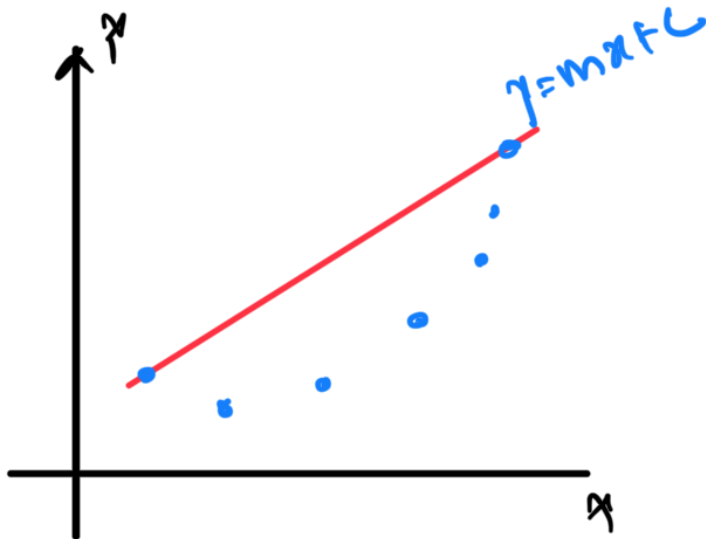
Random error

3

y-intercept

$x=0$    2    4    4.5    $X$

Dependent variable

Independent variable

$x_1, x_2, x_3$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots \epsilon$$

Slop Coefficient

Error term

## Polynomial Regression



$y = mx + c$

SLR

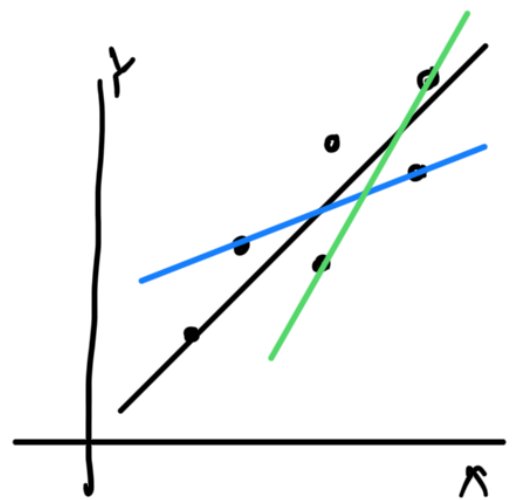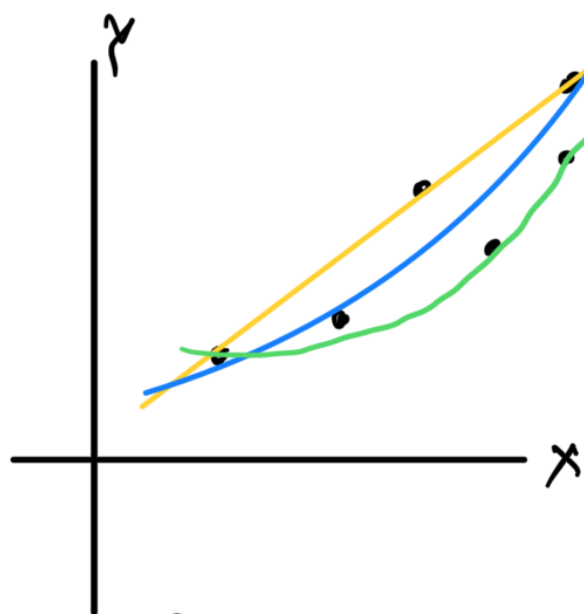Non-linear relationship

## Polynomial Regression

SLR —    $y = \beta_0 + \beta_1 x$

MLR —    $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

PLR —    $y = \beta_0 x^0 + \beta_1 x_1^1 + \beta_2 x_2^2 + \beta_3 x_3^3 + \beta_4 x_4^4$

$x^0 = 1$     degree $= 1$

degree = 2

degree = 4

d=2

d=3

$Y$

$X$

$Y$

$x$

$\underline{\underline{\text{Best fit line}}}$ — Minimize the error between

predicted value and actual value

$$\underline{\underline{\text{Error}}} \leq \in \begin{cases} y = \text{Actual value } (y.\text{test}) \\ \hat{y} = \text{Predicted value ( Model output)} \end{cases}$$

— draw scatterplot & identify the best fit line

## 1. R - Square Method

- statistical method that determines the goodness of best fit

- measure the strength of the relationship between the DV and IV
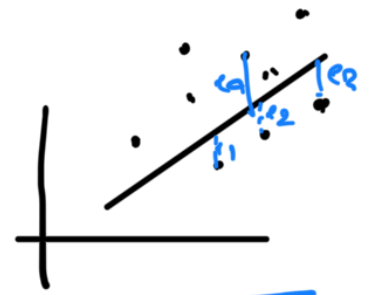
- Calculated 0 — 100%

## 2. Residuals (Regression Error)

- Error in regression represents the difference of the observed data point from the predicted data point in regression line.

Residuals = Actual $y$ ($y_i$) − Predicted $y$ ($\hat{y_i}$)

## 3. Root Mean Square Error (RMSE)

- RMSE represents the standard deviation of the residuals. It gives an estimate of the spread of observed data points across the predicted regression line
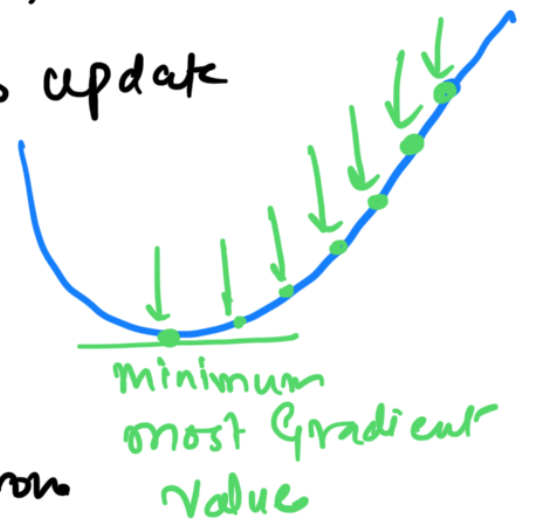
$$RMSE = \sqrt{\frac{e_1^2 + e^2 + e^3 \dots}{n}}$$

## To improve the performance of Regression

Technique — Gradient Descent

Purpose - Minimize the MSE by calculating the gradient of the cost fⁿ (Eqⁿ)

process - Iterative approach to update the coefficients to reduce the cost function

Minimum
most Gradient
value

## Model Performance → ↑ improve

- Goodness of Fit - (Best fit line)
- R - Square → Measure the strength of the relationship betⁿ DV and IV
- Values ranges from 0 − 100%
- Higher values indicates better model fit.
- It is also called as Coefficient of determinⁿ

# Assumption of Linear Regression

1. **Linear Relationship** — Assume a linear relationship between features and the target.

2. **Multicollinearity** — Assume little or no multicollinearity between features

3. **Homoscedasticity** — Assume that the error term is the same for all values of the IV

4. **Normal distribution of Error Term** — Assume error term follow a normal distribution.

5. **No Autocorrelation**: Assume no correlation in error term.
It reduces model accuracy

## Additional Regression Model

1. Ridge Model → L1 Regularization

2. Lasso Model → L2 Regularization

3. Elastic Net Model → (L1 + L2) Regularization

## Regularization Techniques —

$$MSE = \frac{1}{N} \sum_{i=1}^{n} (y_i - \hat{y}_i)$$

Actual $y_i$    Predicted $y_i$

Regulariza t

$$J(\theta) = \frac{1}{N} \sum_{i=1}^{n} (y_i - \hat{y}_i) + \boxed{\lambda} \boxed{\sum_{j=1}^{n} \theta_j^2}$$

Regularisation Term

Regulization Parameter

## Lasso Regression

$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y - \hat{y})^2 + \lambda \sum_{j=1}^{n} |w_j|$$

## Ridge Regression

$$J(\theta) = \frac{1}{n} \underbrace{\sum_{i=1}^{n} (y - \hat{y})^2}_{\text{Loss function}} + \boxed{\lambda} \sum_{j=1}^{n} (|w_j|)^2$$