

Linear Regression { Regressive Prediction model }

- This Model gives o/p in Continuous Value format i.e. $[-\infty$ to $\infty]$ Supervised learning model.

(Unknown Value)

$$y = f(x)$$

Straight line method

(or) Ordinary least square method (OLS)

o/p Variable
(Dependent Variable)

i/p Variable

(Independent Variable)

(Known Value)

(Explanatory Variable)

↓ Regressive Prediction Value
 $[-\infty$ to $\infty]$

This method is used to cal. the Unknown Value based on the Known Value

"Assumes linear relation b/w independent & dependent Variable"

Regression

Statistical tech. used to model the relationship b/w the dependent variable

① the Independent Variable (one/more)

• Based on the Independent Variable we can Predict the dependent Variable.

m & c value calculated based on Historical data

Simple

Linear Regression

(Only one Independent Variable)

→ Straight line Used as Best fit line Prediction

Multiple

Linear Regression

(two or more Independent Variable) i/p Variable

$$y = m_1x_1 + m_2x_2 + \dots + m_nx_n$$

Predicted o/p

$$y = mx + c$$

↑ ↑
Slope Intercept

Predicted o/p based on known value (X)

Goal of Linear Regression \div find Best fit line which minimize the error b/w the predicted o/p & actual o/p based on the historical data (Training dataset)

$$\text{Error (Cost function)} = y - \hat{y}$$

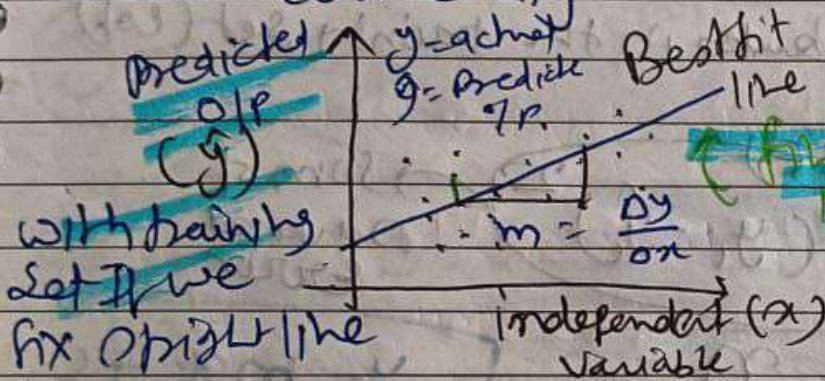
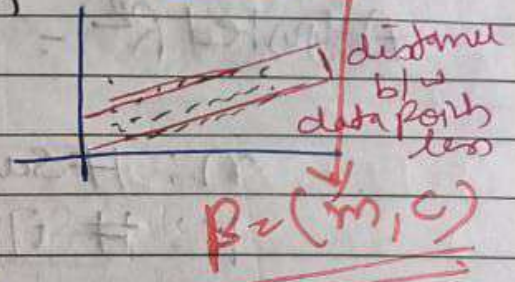
actual o/p predicted o/p

$$(X^T X) \beta = (X^T Y)$$

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

Assumption in the linear Regression

- Linearity ($x \propto y$ or $x \propto \frac{1}{y}$)
- Homoscedasticity
- Multicollinearity



Best fit line use Performance matrix

Performance Matrices \div

- ① R^2 measure
- ② Adjusted R^2 measure
- ③ Mean Square Error (MSE)
- ④ Root Mean Square Error (RMSE)
- ⑤ Mean Absolute Error (MAE)
- ⑥ Mean Absolute Percentage Error (MAPE)

Penalize large errors more than smaller ones

R^2 Measure \div

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

n = no. of samples

Predicted o/p

RSS

TSS

mean value of actual o/p

Variance in the dependent variable that is predictable from independent variable

R^2 value will be measured [0 to 1]
 R^2 value is "1" means no error in the model

No error in the model (100% accuracy)
 R^2 close to zero then more errors present in the model - error is high

Re-training Reg. $\beta (X^T X) = (X^T Y)$

Adjusted R^2 measure

$$\text{Adjusted } R^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - k - 1}$$

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

n : # samples in training set (rows)

k : # of attributes in the training set (col)

Mean Square Error (MSE)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Sum of squared error

$$\text{MSE} = \frac{\text{SSE}}{n}$$

of samples n in training set

④ RMSE (Root Mean Square Error)

$$\text{RMSE} = \sqrt{\text{MSE}}$$

$$y = mx + c$$

$$\textcircled{1} \hat{y} = m\bar{x} + c$$

$$\textcircled{2} m = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

$$\sum XY = m \sum X^2 + c \sum X$$

⑤ Mean Absolute Error

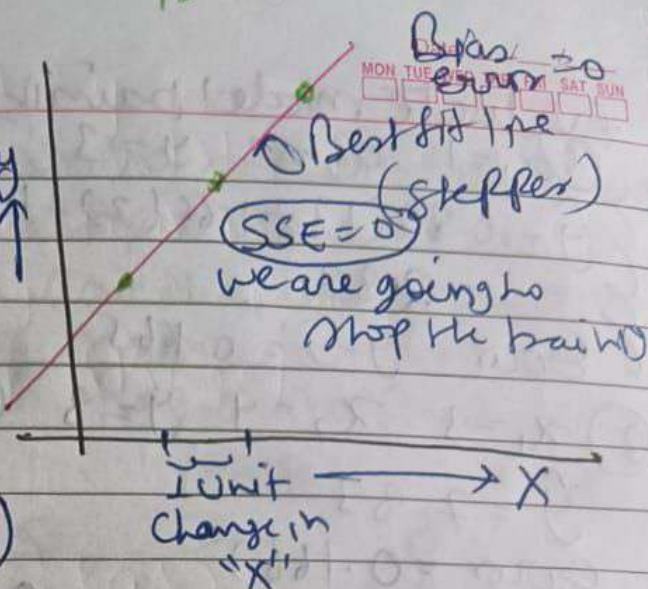
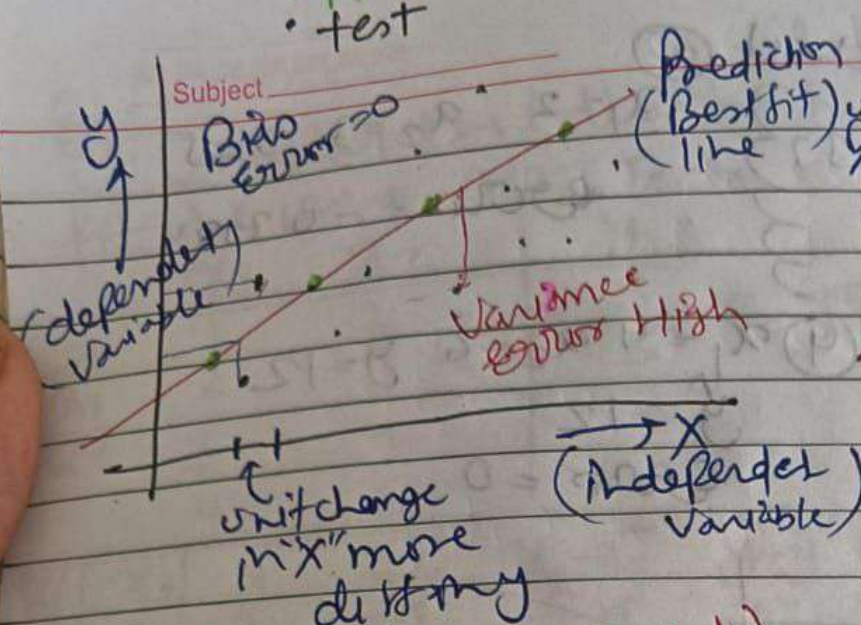
$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

⑥ MAPE (Mean Absolute Percentage Error)

$$\text{MAPE} \times 100$$

• train
• test

• Train



(Overfitting Model)

Regularization Techniques → Introduce some error in the model to improve the performance.

Matrix the model undergoes Overfitting.

Ridge Regularization

Lasso Regularization

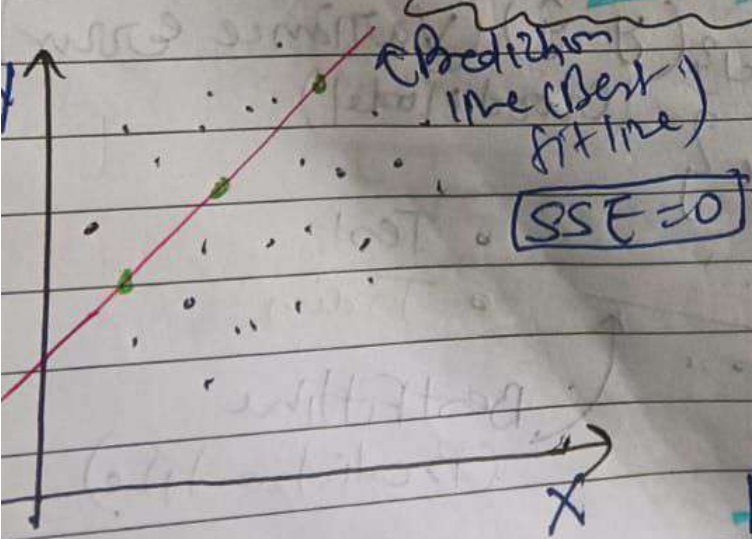
(L-2 Reg)
Ridge Regularization : error is $\lambda (\text{slope})^2$

$\lambda: 0 - \infty$

Error is added to the SSE Prediction line

$$\text{Model} = \text{SSE} + \lambda (\text{slope})^2$$

Circle Boundary



$$\text{SSE} + \lambda (\text{slope})^2$$

$$0 + \lambda (\infty) (\text{slope})^2$$

If SSE again zero Model back on "overfitting"

if $\lambda \rightarrow \infty$ then Prediction line is close to 'x' axis i.e origin.

→ L1 Regularization (like slope)
lasso Regression = lasso penalty + error is $\lambda |\text{slope}|$ (0-∞)

Add lasso penalty to the SSE, to improve the model accuracy when the model undergoes overfitting

Lasso Penalty = SSE + $\lambda |\text{slope}|$
Model Error

Major diff. b/w Ridge & lasso

Ridge

In this tech, new Prediction line (Slope) close to origin (X-axis) but it never reach to origin

$\lambda (\text{slope})^2$

lasso

In this tech, new prediction line (Slope) is equal to origin (X-axis)

$\lambda |\text{slope}|$
 Slope of any attribute is 0 then we remove that attribute
feature selection

Simple Linear Regres.

$\hat{y} = m_1 x_1 + c$

Lasso Penalty = SSE + $\lambda |m_1|$

Multiple Linear Reg.

$\hat{y} = m_1 x_1 + m_2 x_2 + \dots + m_n x_n$

Lasso Penalty = SSE + $\lambda [|m_1| + |m_2| + \dots + |m_n|]$

Eg. Consider the following multiple linear Reg.
 $\hat{y} = m_1 x_1 + m_2 x_2 + m_3 x_3 + m_4 x_4 + c$

Suppose

lasso penalty = SSE + $\lambda [|m_1| + |m_2| + |m_3| + |m_4|]$

Assume m_2 & m_4 are zero then realize that attributes 2 & 4 are not so important (less) in model

So delete 2 & 4 attributes from training dataset
 $\hat{y} = m_1 x_1 + m_3 x_3$