# Recap of Previous Lecture

Topic

Topic

Topic

Topic

Topic

- OLS → overfit
  → Multicollinearity → unstable model
  ↳ means Min RSS
  → Solution Ridge Reg.
  Regularisation

# Topics to be Covered

Topic

**Ridge regression**

Topic

Topic

Topic

Topic

Ridge Regression Final expression

## Shrinkage Methods : Ridge Regression

❖ Ridge regression is a regularisation techniques...

updated loss fxn.

Hyperparameter

$$d = \left[ \sum_{i=1}^{N} \left( y_i - \hat{y}_i \right)^2 + \frac{\lambda}{2} \sum_{i=1}^{D} \beta_i^2 \right]$$

$\beta_0$ not included ??

model LR wala hai $\left( \hat{y}_i = \beta_0 + \beta_1 x^1 + \beta_2 x^2 + \cdots \right)$

→ we want to min the loss fxn, to find best $\beta$'s.

## Shrinkage Methods : Ridge Regression

❖ "In regularization technique, we reduce the magnitude of the features by keeping the same number of features.

❖ This helps in ....

(i) we are not reducing / limiting the No of dimension.
for example if data has 100 dimensions
then we train model with 100 dimensions

(ii) In Regularisation we put limit on $\beta$'s of the dimension.

## Shrinkage Methods : Ridge Regression

❖ **Ridge regression shrinks the regression coefficients by imposing a penalty on their size.**

❖ **The ridge coefficients minimize a penalized residual sum of squares of the weights.**

$\longrightarrow$ we add $\dfrac{\lambda}{2} \sum\limits_{i=1}^{D} \beta_i^2$ in the loss fxn.

The loss function are updated

## Shrinkage Methods : Ridge Regression

done

The loss function are updated

## Shrinkage Methods : Ridge Regression

The main reason for not regularizing the intercept term is that it represents the mean value of the target variable when all the features are zero. Regularizing the intercept can lead to shifting this mean value away from its natural value, which might not be desirable in many cases.

Why the $w_0$ term is not included in regularisation ..

* Why $\beta_0$ is not included
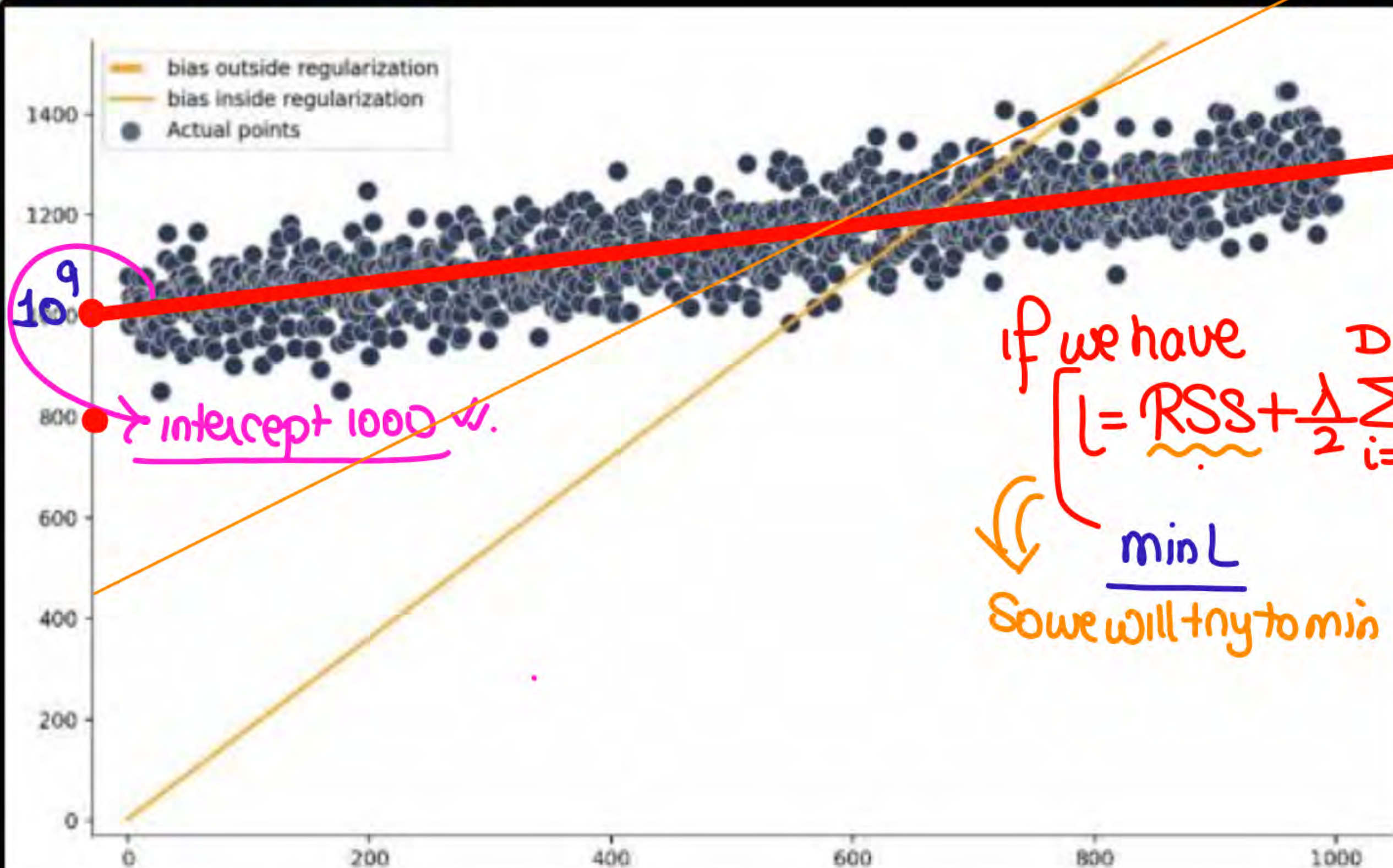* Why we donot want $\beta_0$ to be minimize ??

we always want values of $\beta$'s sha be
within some limit

bcoz if any dimension get large $\beta$, then that
dimension dominate model ⇒ unstable model.

bias outside regularization
bias inside regularization
● Actual points

$10^9$

intercept 1000 w.

if we have

$$L = RSS + \frac{\lambda}{2} \sum_{i=0}^{D} \beta_i^2$$

$\min L$

Some will try to min $\beta_0$

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

$$\underline{\beta_0} = \bar{y} - \beta_1 \bar{x}$$

Similarly

$$y = \beta_0 + \beta_1 x^1 + \beta_2 x^2 + \text{ --- }$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x^1} - \beta_2 \bar{x^2} \text{ --- } )$$

* Value of $\beta_0$ depend on mean location of data

* $\beta_0$ shd not be minimized

# Complete analysis of RR

we have the data $\Rightarrow$ Step1 Remove the $\beta_0$ from analysis.

So Create Centred data

Now $x^1$.    $x^2$.    $y$.

| $x^1$ | $x^2$ | $y$ |
|---|---|---|
| | | |

$\hookrightarrow \overline{x^1}$  $\hookrightarrow \overline{x^2}$  $\hookrightarrow \overline{y}$

| $x^1 - \overline{x^1}$ | $x^2 - \overline{x^2}$ | $y - \overline{y}$ |
|---|---|---|
| | | |

Now apply RR $\Rightarrow$ 2D data.

$$\mathcal{L} = \sum_{i=1}^{N} (y_i - \hat{y}_i)^2 + \frac{\lambda}{2} \sum_{i=1}^{2} \beta_i^2$$

N data points
and 2 dimension

$$\left( \hat{y}_i = \beta_1 x_i^1 + \beta_2 x_i^2 \right)$$

$$\mathcal{L} = \sum_{i=1}^{N} \left( y_i - \left( \beta_1 x_i^1 + \beta_2 x_i^2 \right) \right)^2 + \frac{\lambda}{2} \sum_{i=1}^{2} \beta_i^2$$

min L $\Rightarrow$

$$\frac{\partial L}{\partial \beta_1} = -2 \sum_{i=1}^{N} x_i^1 (y_i - \beta_1 x_i^1 - \beta_2 x_i^2) + \lambda \beta_1 = 0$$

$$\frac{\partial L}{\partial \beta_2} = -2 \sum_{i=1}^{N} x_i^2 (y_i - \beta_1 x_i^1 - \beta_2 x_i^2) + \lambda \beta_2 = 0$$

$$\min L \Rightarrow \frac{\partial L}{\partial \beta_1} = -2 \sum_{i=1}^{N} x_i^1 (y_i - \beta_1 x_i^1 - \beta_2 x_i^2) + \lambda \beta_1 = 0$$

$$\frac{\partial L}{\partial \beta_2} = -2 \sum_{i=1}^{N} x_i^2 (y_i - \beta_1 x_i^1 - \beta_2 x_i^2) + \lambda \beta_2 = 0$$

$$\frac{\partial L}{\partial \beta} = \begin{bmatrix} \partial L/\partial \beta_1 \\ \partial L/\partial \beta_2 \end{bmatrix} = -2 \begin{bmatrix} \sum x_i^1 y_i - \beta_1 \sum (x_i^1)^2 - \beta_2 \sum x_i^1 x_i^2 \\ \sum x_i^2 y_i - \beta_1 \sum x_i^2 x_i^1 - \beta_2 \sum (x_i^2)^2 \end{bmatrix} + \lambda \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = 0$$

$$= -2 \left[ \begin{bmatrix} \sum x_i^1 y_i \\ \sum x_i^2 y_i \end{bmatrix} - \begin{bmatrix} \sum (x_i^1)^2 & \sum x_i^1 x_i^2 \\ \sum x_i^1 x_i^2 & \sum (x_i^2)^2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \right] + \lambda \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = 0$$

$$\frac{\partial L}{\partial \beta} = \begin{bmatrix} \partial y / \partial \beta_1 \\ \partial y / \partial \beta_2 \end{bmatrix} = -2 \begin{bmatrix} \sum x_i^1 y_i - \beta_1 \sum (x_i^1)^2 - \beta_2 \sum x_i^1 x_i^2 \\ \sum x_i^2 y_i - \beta_1 \sum x_i^2 x_i^1 - \beta_2 \sum (x_i^2)^2 \end{bmatrix} + \lambda \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = 0$$

$\nearrow X^T Y$

$$= -2 \left[ \begin{bmatrix} \sum x_i^1 y_i \\ \sum x_i^2 y_i \end{bmatrix} - \begin{bmatrix} \sum (x_i^1)^2 & \sum x_i^1 x_i^2 \\ \sum x_i^1 x_i^2 & \sum (x_i^2)^2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \right] + \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = 0$$

$\searrow X^T X$

$\hookrightarrow \beta$

$$X = \begin{bmatrix} x_1^1 & x_1^2 \\ x_2^1 & x_2^2 \\ x_3^1 & x_3^2 \\ \vdots & \vdots \end{bmatrix} \quad X^T = \begin{bmatrix} x_1^1 & x_2^1 & x_3^1 & \cdots \\ x_1^2 & x_2^2 & x_3^2 & \cdots \end{bmatrix}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \end{bmatrix}$$

$$\frac{\partial L}{\partial \beta} = -2\left[ X^T Y - (X^T X)\beta \right] + \lambda I \beta = 0$$

$$= -2X^T Y + 2(X^T X)\beta + \lambda I \beta = 0$$

$$\Rightarrow (X^T X)\beta + \frac{\lambda}{2}I\beta = X^T Y$$

$$\Rightarrow \left( X^T X + \frac{\lambda}{2}I \right) \beta = X^T Y$$

$$\beta = \left( X^T X + \frac{\lambda}{2}I \right)^{-1} X^T Y.$$

① Centre the data

⑪ we will get $\beta_1, \beta_2$

Centred data ka model
$$y = \beta_1 \underline{x'} + \beta_2 \underline{x^2}$$

Centred

Original model will be

$$y = \beta_0 + \beta_1 \underline{x^1} + \beta_2 \underline{x^2}$$

Original Values

$$\beta_0 = \bar{y} - \beta_1 \bar{x}^1 - \beta_2 \bar{x}^2$$

## 2D data

| $x^1$ | $x^2$ | $y$ |
|-------|-------|-----|
| 5 | 8 | 16 |
| 7 | 10 | 26 |
| 9 | 12 | 30 |
| 7 | 14 | 40 |

$\overline{x^1} = 7 \quad \overline{x^2} = 11 \quad \overline{y} = 28$

Find $\beta_0, \beta_1, \beta_2$.

Ridge Reg $\lambda = 2$.

| $x^1 - \overline{x^1}$ | $x^2 - \overline{x^2}$ | $y - \overline{y}$ |
|-------|-------|-----|
| -2 | -3 | -12 |
| 0 | -1 | -2 |
| 2 | 1 | 2 |
| 0 | 3 | 12 |

$$\left(X^TX + \frac{\lambda}{2}I\right) = \begin{bmatrix} 9 & 8 \\ 8 & 21 \end{bmatrix}$$

$$X = \begin{bmatrix} -2 & -3 \\ 0 & -1 \\ 2 & 1 \\ 0 & 3 \end{bmatrix} \quad X^T = \begin{bmatrix} -2 & 0 & 2 & 0 \\ -3 & -1 & 1 & 3 \end{bmatrix}$$

$$(X^TX) = \begin{bmatrix} 8 & 8 \\ 8 & 20 \end{bmatrix}$$

$$X^TY = \begin{bmatrix} 28 \\ 76 \end{bmatrix}$$

$$\beta = \left(X^TX + \frac{\lambda}{2}I\right)^{-1}X^TY$$

$$= \frac{1}{125}\begin{bmatrix} 21 & -8 \\ -8 & 9 \end{bmatrix}\begin{bmatrix} 28 \\ 76 \end{bmatrix} \Rightarrow \begin{bmatrix} 20/125 \\ 460/125 \end{bmatrix}$$

$$\checkmark \beta_1 = -\cdot 16$$

$$\checkmark \beta_2 = 3\cdot 68$$

$$\overline{y} - \beta_1 \overline{x}_1 - \beta_2 \overline{x}^2 = \beta_0$$

$$\checkmark \beta_0 = -11\cdot 36$$

- if RR
- $L = RSS + \lambda \sum_{i=1}^{D} \beta_i^2$

$$\left( X^T X + \lambda I \right)^{-1} X^T Y$$

❖ Here $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage:

- we always keep $\lambda \geq 0$

$$\alpha = \left\{ \sum_{i=1}^{N} \left( y_i - \hat{y}_i \right)^2 + \frac{\lambda}{2} \sum_{i=1}^{D} \beta_i^2 \right\}$$

- $\lambda = 0 \Rightarrow$ Same as L.R

- $\lambda = $ Bahut Badi

$\lambda :$ hyper parameter

$\lambda \sim\sim\sim\sim\sim$

Control Karta hai.

loss $\mathcal{F}$xn has 2 terms

RSS + Penalty.

- $\lambda$ Control that which term is more Imp

- $\lambda$ Large Penalty term imp

- we always keep $\lambda \geq 0$

$$\alpha = \left\{ \sum_{i=1}^{N} (y_i - \hat{y_i})^2 + \frac{\lambda}{2} \sum_{i=1}^{D} \beta_i^2 \right\}$$

- $\lambda = 0 \Rightarrow$ same as L.R
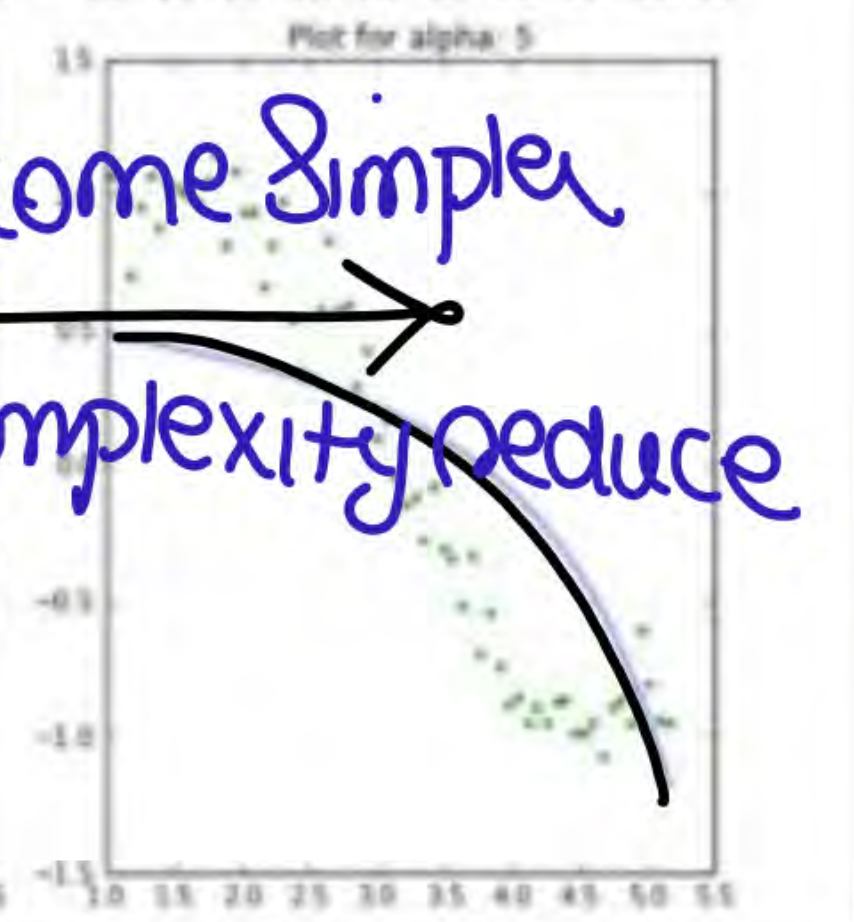
- $\lambda =$ Bahut Badi

$\lambda$: hyperparameter

$\lambda$ ~~~~~~~

Control Karta hai.

- $\lambda$ is very small then RSS shd be zero or very small even if $\sum \beta^2$ is large

loss fxn has 2 terms
RSS + Penalty.

- $\lambda$ Control that which term is more Imp

- $\lambda$ large Penalty term imp
  → matlal chahe RSS is large, we need $\sum \beta^2$ to be v. small

# λ Ko bachane Se Kya hoga.

$\lambda = 0$ small

RSS Pyara hai

RSS → 0

$\sum \beta^2$ Can be large

*Same as Linear Reg. overfit model.*

How to find best β

hyperparameter tuning

Cross Validation

λ large

RSS kuch hoJae

$\sum \beta^2$ very very Small

*Underfit model*

- we always keep $\lambda \geqslant 0$

$$\alpha = \left\{ \sum_{i=1}^{N} \left( y_i - \hat{y}_i \right)^2 + \frac{\lambda}{2} \sum_{i=1}^{D} \beta_i^2 \right\}$$

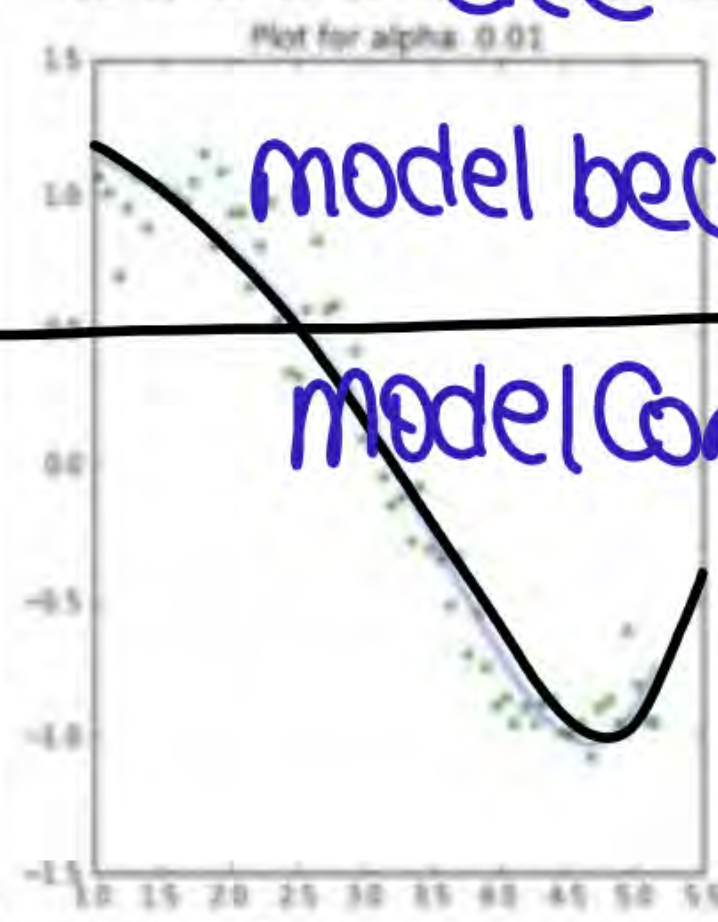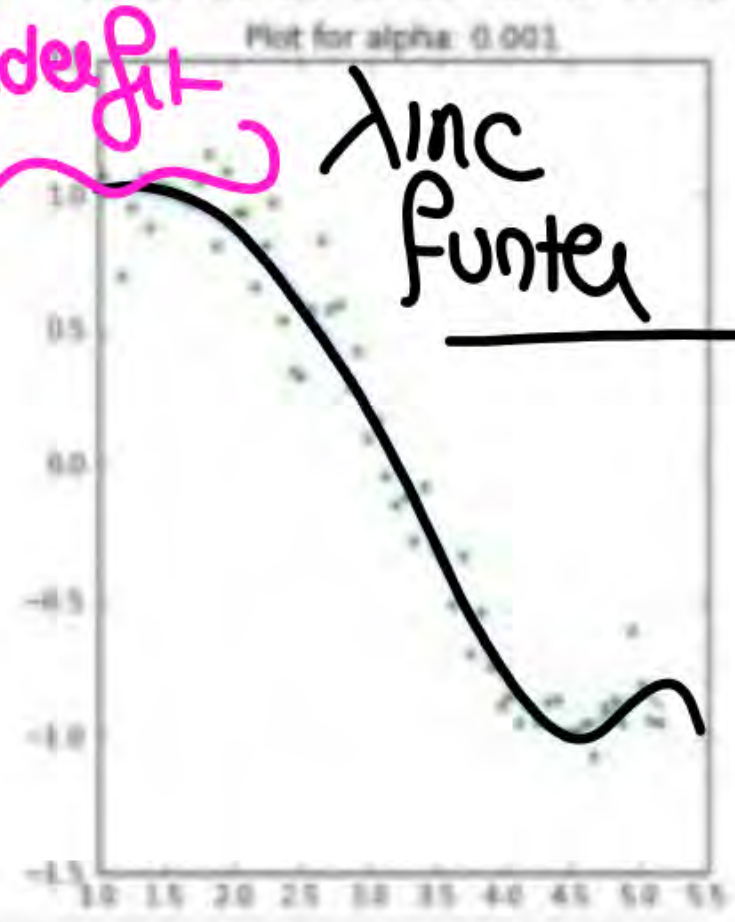- $\lambda = 0 \Rightarrow$ Same as L.R
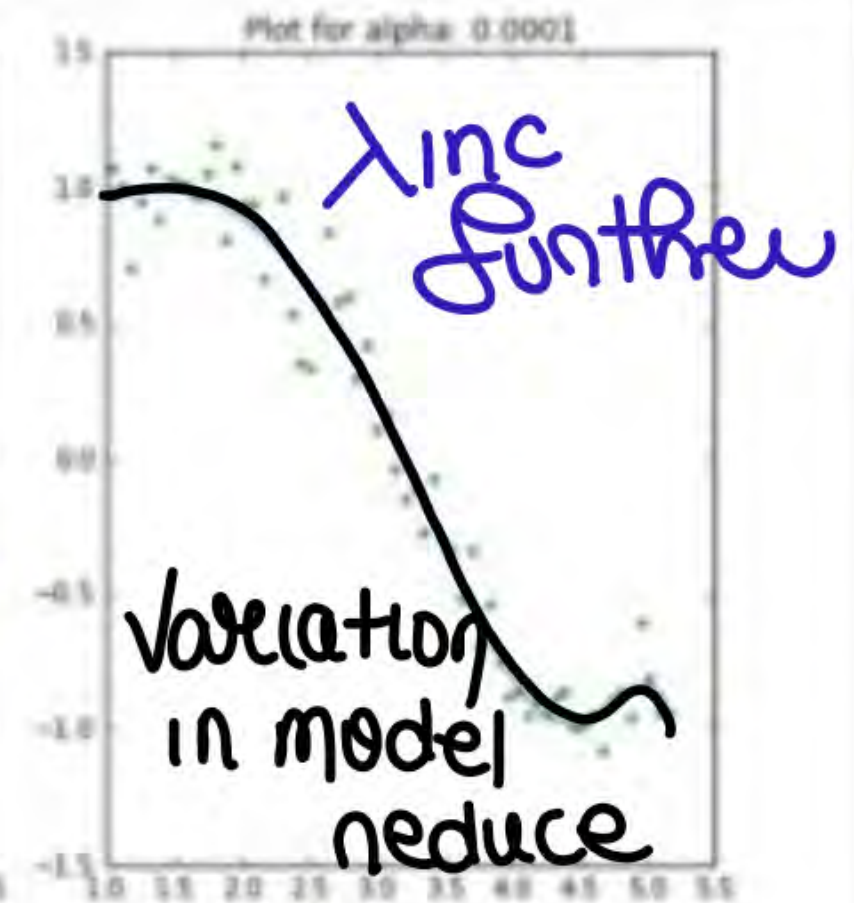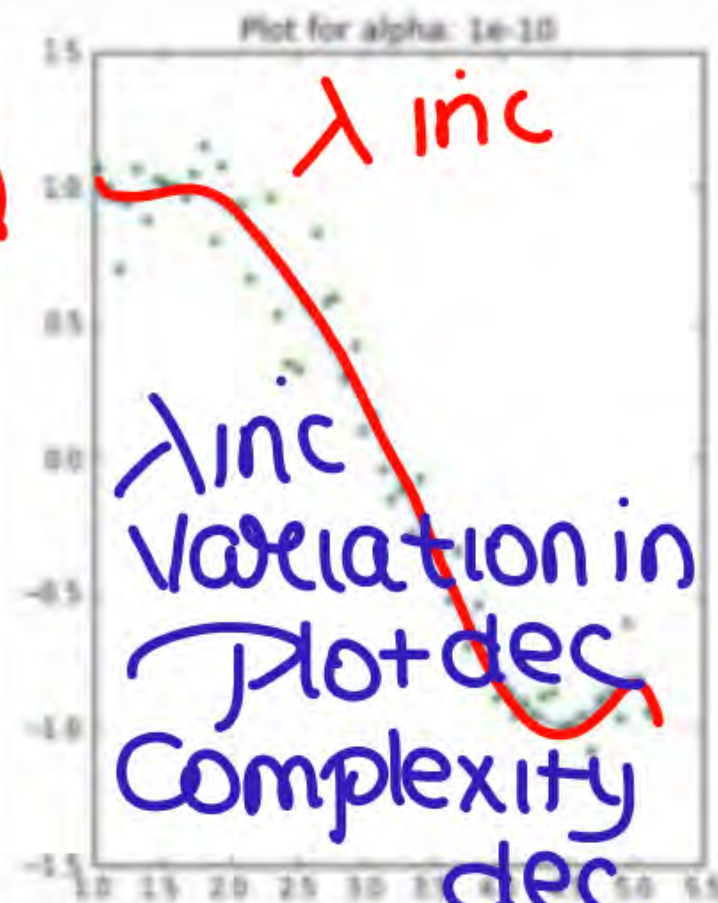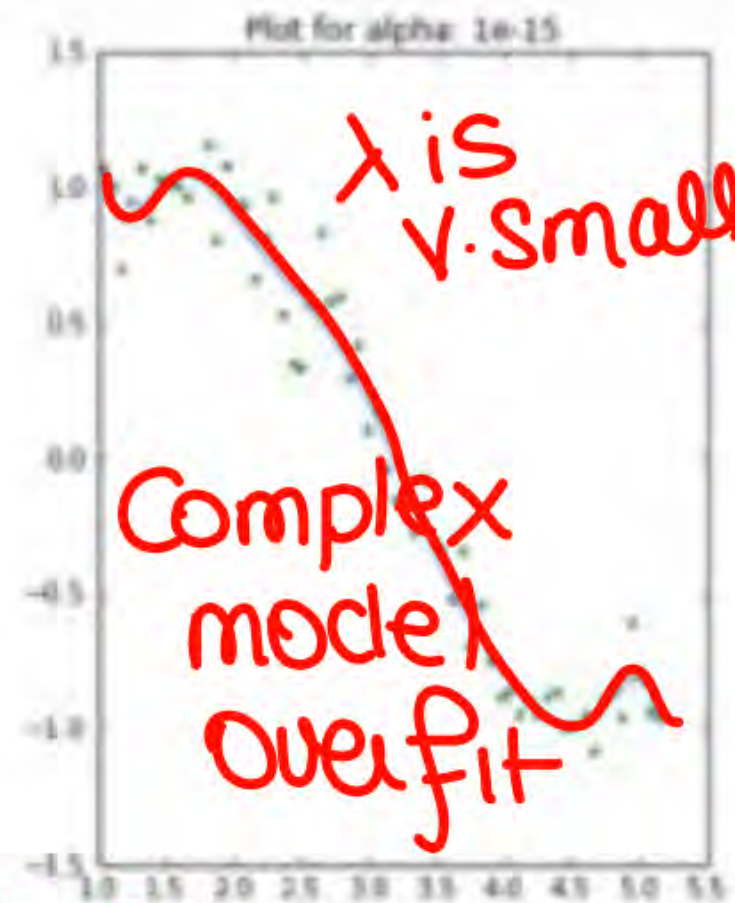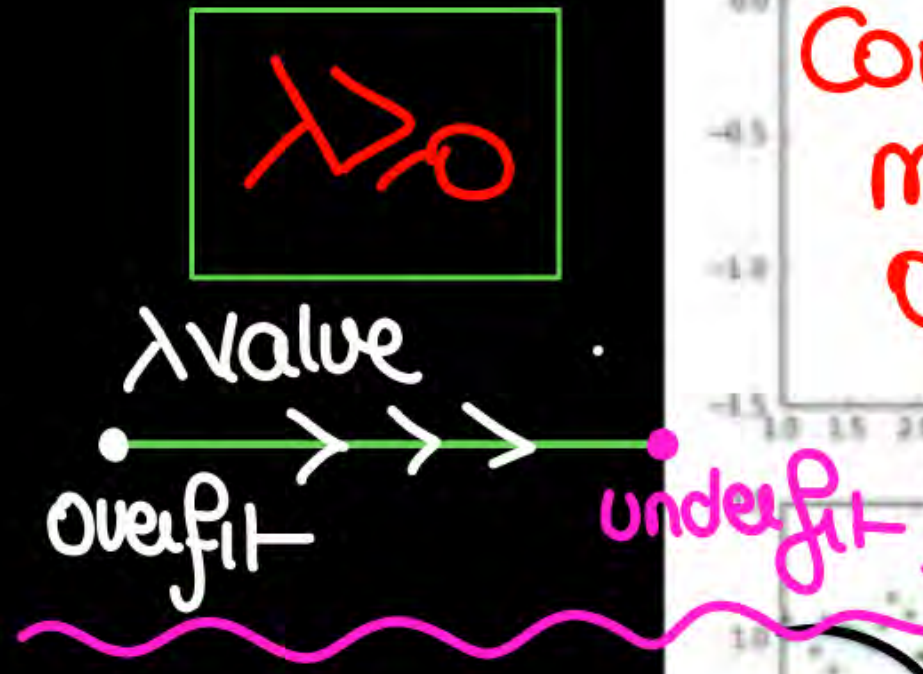- $\lambda =$ Bahut Badi

$\lambda$: hyper parameter

$\lambda$ ~~~~~~~~~~
Control Karta hai.

what if $\lambda$ is -ve

- if $\lambda$ is -ve
- minimize loss fxn.
- underfit model
- -ve $\lambda$ never used.

λ > 0

λ value → → → overfit ... underfit

Plot for alpha: 1e-15
λ is v. small
Complex model overfit

Plot for alpha: 1e-10
λ inc
λ inc Variation in plot dec Complexity dec

Plot for alpha: 0.0001
λ inc further
Variation in model reduce

Plot for alpha: 0.001
λ inc furter

Plot for alpha: 0.01
model become simpler
model Complexity reduce

Plot for alpha: 5

How to find best $\lambda$

many algo we have

In many algo we have hyperparameter.

Inki value KO hit and trail karke best value choose karte hai

they are used while training the model.

Parameters $\Rightarrow$ Jo hamari final model ki equation ka part bante hein

Values of these are found after training process ends.

THANK - YOU