

# Data Science and Artificial Intelligence

## Machine Learning



**Regression**

**Lecture No. 9**



**By- SIDDHARTH SABHARWAL SIR**



# Recap of Previous Lecture



Topic

Ridge Regression

Topic

$\beta$ 's expression

Topic

Topic

Topic



# Topics to be Covered



Topic

What is  $\lambda$

Topic

Effect of  $\lambda$ , how to find best value of  $\lambda$

Topic

Lasso Regularisation (only v. brief)

Topic

Questions.

Topic



**THINK BIG.  
TRUST  
YOURSELF  
AND MAKE  
IT HAPPEN**

apne/  
Teachers!

Think big  
AIR 1.





## Ridge Regression Final expression

$$\bullet \Rightarrow \beta = (X^T X + \lambda I)^{-1} X^T Y$$

→ we use Centred data here

$$\rightarrow y = y - \bar{y}, x^i = x^i - \bar{x}^i$$

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

→  $\beta$ 's obtained from the Centred data is valid in original data.



## Ridge Regression Final expression

- $\beta_0$  is not included in Reg term,

- L2 Reg. term

- $L \Rightarrow \min \left( \frac{1}{2} \sum (y_i - \hat{y}_i)^2 + \frac{\lambda}{2} \sum \beta_i^2 \right)$

Constant  
we can find best  $\lambda$ .



$X = \begin{bmatrix} \end{bmatrix}_{n \times 1}$ ,  $n$ : Number of unknowns

- A is matrix of size  $(n \times n)$

- The basic Linear Algebra  $\Rightarrow$

$$Ax = b$$

## Unique Solution

No solution

$\infty$  many Solution

We create a Augmented matrix

$$\begin{bmatrix} A & B \\ \vdots & \vdots \end{bmatrix}$$

⇒ Solution ⇒  $X = A^{-1}B$

1. Rank of Augmented matrix = Rank of  $A = n \rightarrow$  Unique sol



Q. What will happen if  $X^T X^{-1} \Rightarrow$  non Inv

So in linear regression

$$\underbrace{(X^T X)}_A \beta = \underbrace{X^T Y}_B$$

$$\beta = \underbrace{(X^T X)^{-1}}_{\text{non Inv}} X^T Y$$

$$X^T Y = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}_{D+1 \times 1}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_D \end{bmatrix}_{D+1 \times 1}$$

$$(X^T X) = \begin{bmatrix} \vdots & \vdots \end{bmatrix}_{D+1 \times D+1}$$

•  $(X^T X)^{-1}$  will have dimension of  $(D+1 \times D+1)$

$\rightarrow X^T X$  will be non Inv if  $|X^T X| = 0$   
OR Rank of  $X^T X < (D+1)$



Now Augmented matrix in this case

$$\left[ \begin{array}{c|c} (X^T X) & X^T Y \\ \hline A & B \end{array} \right]$$

If  $X^T X$  is non inv then Rank of A will be  $< D+1$ , and in that case we can have 2 Situation

→ i.e No Solution ✓

→ i.e  $\infty$  many Solutions ✓

$$\underbrace{\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}}_A \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \underbrace{\begin{bmatrix} j \\ k \\ l \end{bmatrix}}_B$$

$$\text{Aug} = \begin{bmatrix} a & b & c & j \\ d & e & f & k \\ g & h & i & l \end{bmatrix} \quad \text{3x4}$$

If Rank of  $A=2$ , then Rank of Aug  $\Rightarrow$  Can be  $3/2$ .









# Ridge Regression



## Shrinkage Methods : Ridge Regression

### ❖ Solution to this ridge regression problem

$$\text{So loss function} \Rightarrow \min \left( \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \frac{\lambda}{2} \sum_{i=1}^D \beta_i^2 \right)$$

•  $\lambda=0 \Rightarrow$  this becomes basic LR  $\Rightarrow$  Overfitting

• if  $\lambda=\infty$  then Regularisation terms become dominant, the algorithm try to make all

$\beta=0$   
 $\rightarrow$  underfitting  
 $\rightarrow$  high training and testing error.

$\rightarrow$  V. low training error  
 $\rightarrow$  i.e V. low Bias  
 $\rightarrow$  V. high testing error  
 $\rightarrow$  i.e V. high Variance





# Ridge Regression



## Shrinkage Methods : Ridge Regression

❖ Solution to this ridge regression problem

done





# Ridge Regression



## ❖ How to choose the best $\lambda$ .

- In basic linear regression the complexity was very less
- we solve the algorithm just once to find  $\beta_0$

Cross  
Validation

- But in RR the algorithm has to run for various  $\lambda$ 's and then we find the best  $\lambda$  to be used in the algorithm.
- So RR has much more complexity than LR



# Ridge Regression



## ❖ How to choose the best $\lambda$ .



- So data is broken into training & testing data
- Testing data is not available to us
- we only have the training data.

**Cross  
Validation**



To find best  $\lambda \Rightarrow$

So we break the Training data into 2 parts  $\Rightarrow$   
Training and Validation parts.

- So we train the model on training part for various  $\lambda$ 's and then test the model on the validation part.
- The  $\lambda$  that give the best performance on validation part will be chosen

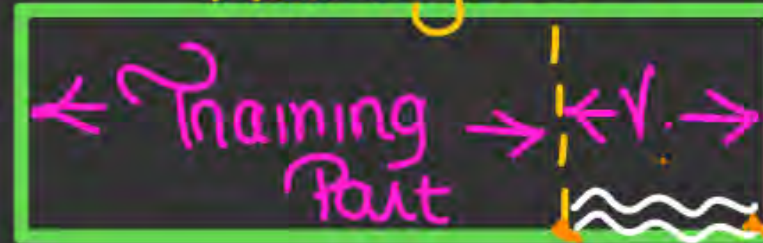




data



Training data



Validation part

$\lambda$

- 1
- 2
- 3
- 4
- 5
- 6
- ...

Solve the loss  
fxn for each  $\lambda$   
On the training  
Part

$\beta$   
 $\beta$   
 $\beta$

Now work on Validation  
Part  $(x \cdot \beta) = \hat{y}$

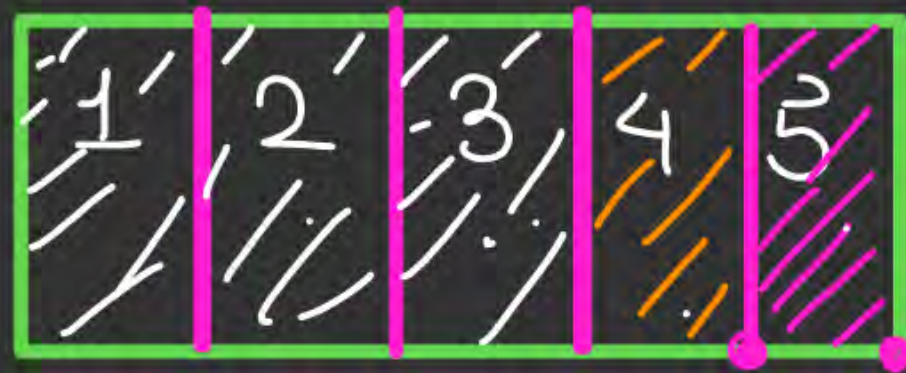
$\sum |y - \hat{y}| \Rightarrow$  error on  
Validation part

- So for various  $\lambda$ 's we train the model and then test it on the validation and then simply find best  $\lambda$

• But in the above process we are actually finding best  $\lambda$  that suits our validation part.



So we do cross validation



Step 3  $\Rightarrow$

Step 4  $\Rightarrow$

Step 5  $\Rightarrow$

5-fold cross validation

Training data is broken into 5 parts

Step 1  $\Rightarrow$  1, 2, 3, 4  $\Rightarrow$  Training Part  
5  $\Rightarrow$  Validation part

$\rightarrow$  find best  $\lambda$ .

Step 2  $\Rightarrow$  1, 2, 3, 5  $\Rightarrow$  Training Part  
4  $\Rightarrow$  Validation  
 $\rightarrow$  find best

○ The best  $\lambda$  for overall data  $\Rightarrow$  avg of all

$\hookrightarrow$  This process is called  $K$  fold cross validation.





## Ridge Regression

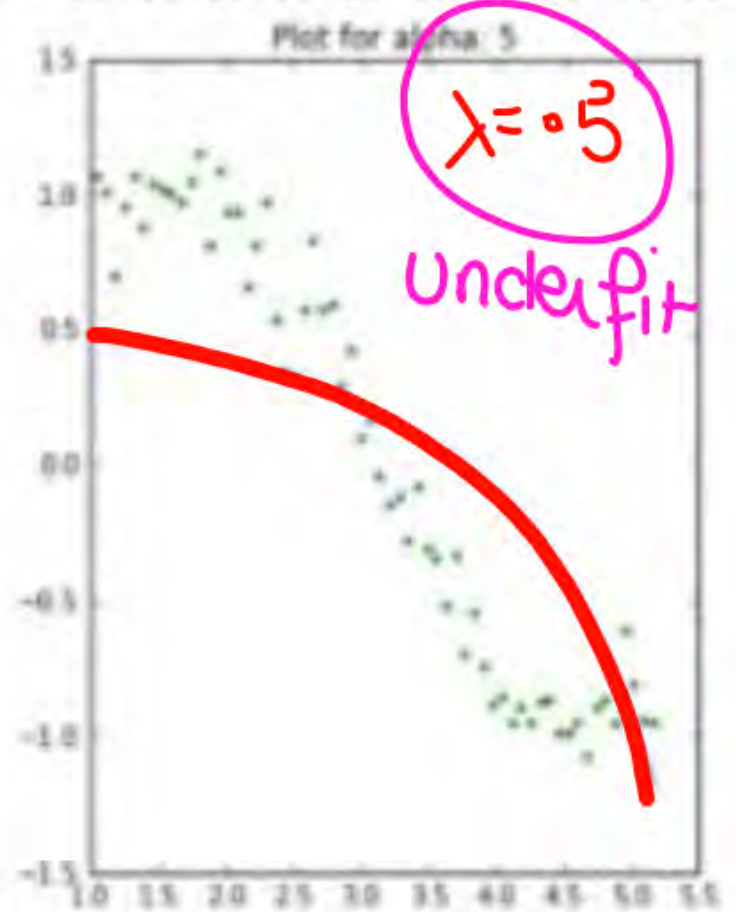
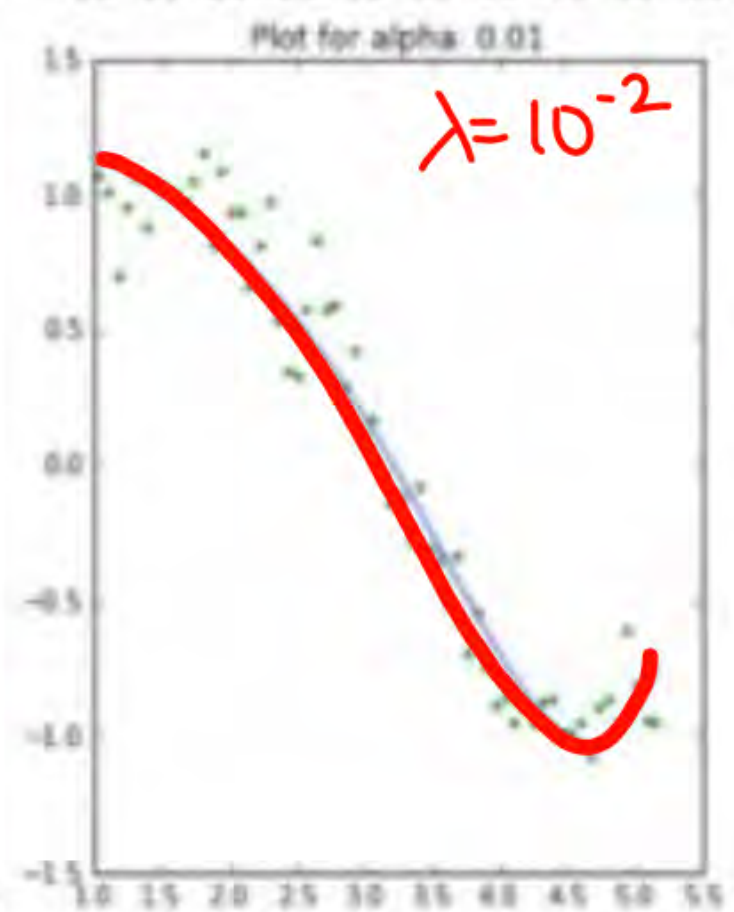
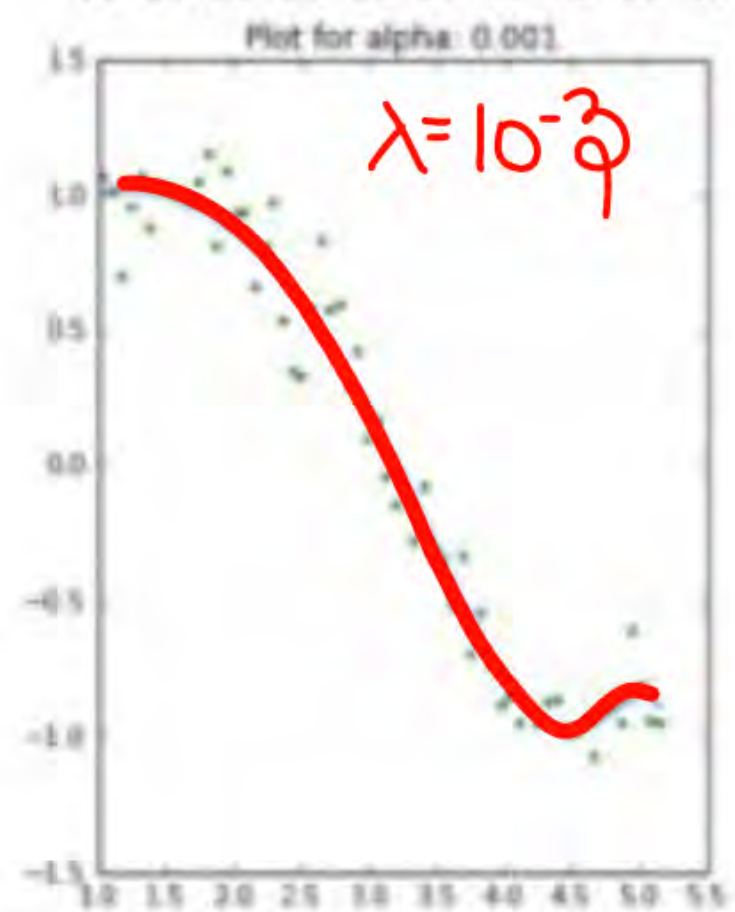
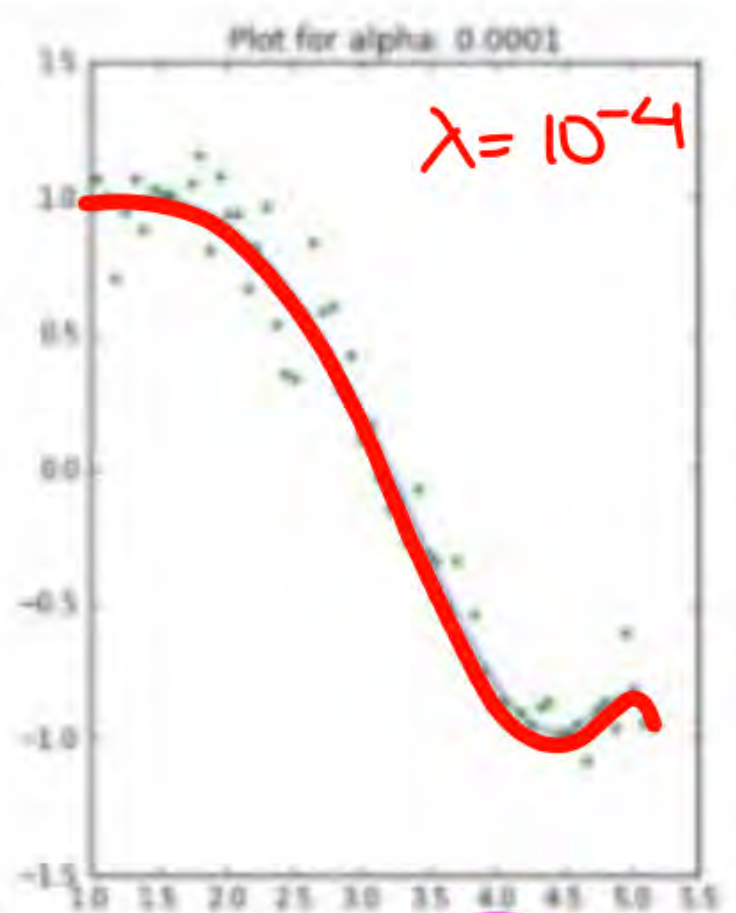
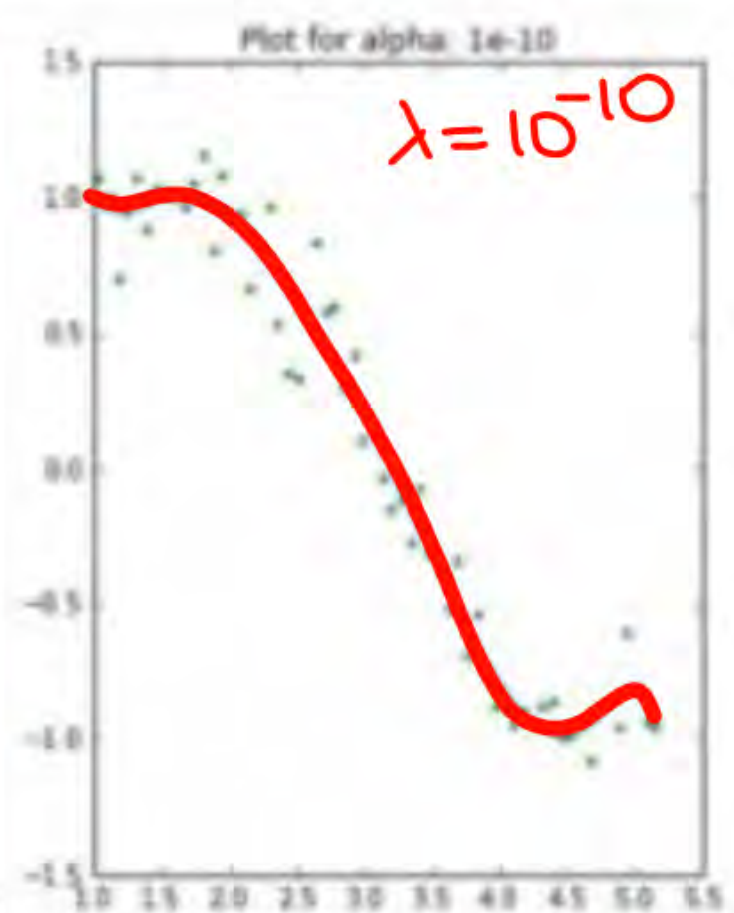
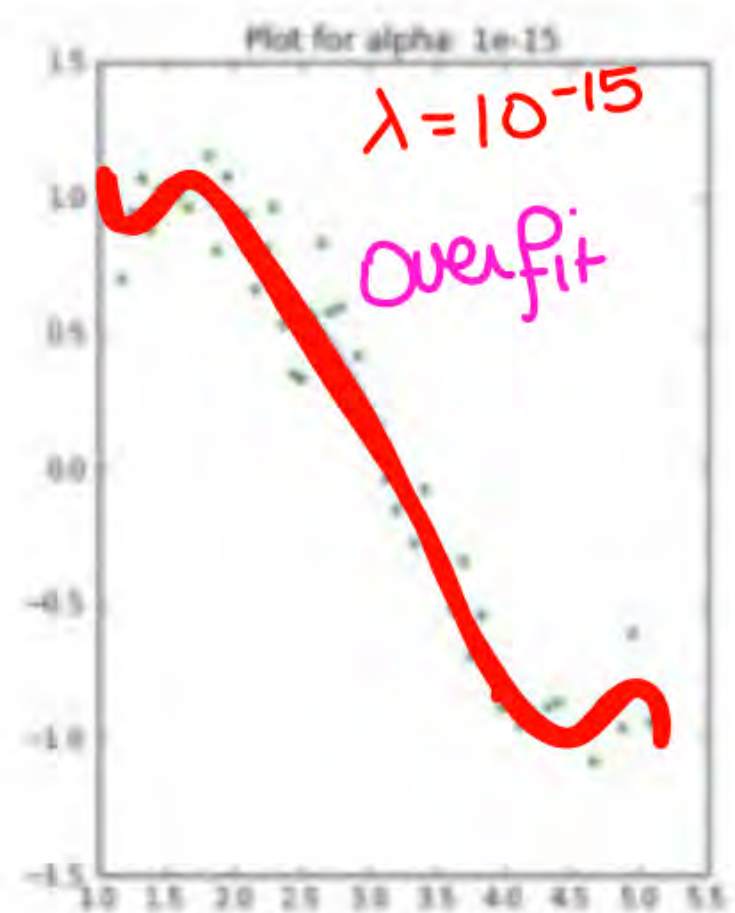


❖ Can  $\lambda$  be negative.

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

\* if  $\lambda$  is negative  $\Rightarrow$  then  $\beta \rightarrow \infty$

\* so  $\lambda$  will never be -ve.





So as  $\lambda$  inc

→ model Complexity reduce

→ Training error inc / becoz overfitting reduce

→ model generalise better

→ Testing error reduce

But large  $\lambda \Rightarrow$  underfitting



## Ridge Regression



### Ridge Regression – lets practise

Ridge Regression is a regularization technique used in linear regression to:

- A) Increase model complexity.
- B) Reduce model complexity and prevent overfitting.
- C) Make the model fit the training data perfectly.
- D) Enhance the interpretability of the model.





## Ridge Regression



### Ridge Regression – lets practise

In Ridge Regression, the penalty term added to the cost function is based on:

- A) The absolute values of the regression coefficients.
- B) The square of the regression coefficients.
- C) The number of features.
- D) The dependent variable.



## Ridge Regression



### Ridge Regression – lets practise

What happens to the magnitude of regression coefficients in Ridge Regression compared to ordinary linear regression?

- A) They become larger.
- B) They become smaller.
- C) They stay the same.
- D) It depends on the dataset.





## Ridge Regression



### Ridge Regression – lets practise

Ridge Regression is particularly useful when:

- A) There is no multicollinearity among the independent variables.
- B) There is a high degree of multicollinearity among the independent variables.
- C) The model needs to fit the training data perfectly.
- D) The dataset has very few observations.



## Ridge Regression



### Ridge Regression – lets practise

Which of the following values of  $\lambda$  (lambda) in Ridge Regression would lead to the strongest regularization effect?

- A)  $\lambda = 0$
- B)  $\lambda = 1$
- C)  $\lambda = 10$
- D)  $\lambda = \infty$





## Ridge Regression



### Ridge Regression – lets practise

Ridge Regression can help prevent overfitting, but what is the trade-off?

- A) Increased model interpretability.
- B) Increased computational complexity.
- C) Reduced accuracy on the training data.
- D) Smaller training dataset size.



## Ridge Regression



### Ridge Regression – lets practise

In Ridge Regression, what is the effect of increasing  $\lambda$  (lambda) on the bias and variance of the model?

- A) Increases bias, decreases variance.
- B) Decreases bias, increases variance.
- C) Increases both bias and variance.
- D) Decreases both bias and variance.





### Ridge Regression – lets practise

In Ridge Regression, the penalty term added to the cost function is based on the L2 norm (Euclidean norm) of the regression coefficients. If the sum of squared regression coefficients (L2 norm) is 50 and the value of  $\lambda$  (lambda) is 3, what is the modified penalty term in the Ridge Regression cost function?

- a) 150
- b) 135
- c) 123
- d) 578



## Ridge Regression



### Ridge Regression – lets practise

In Ridge Regression, the penalty term added to the cost function is based on the L2 norm (Euclidean norm) of the regression coefficients. If the sum of squared regression coefficients (L2 norm) is 50 and the value of  $\lambda$  (lambda) is 3, what is the modified penalty term in the Ridge Regression cost function?

- a)150
- b)135
- c)123
- d)578





## Ridge Regression



### Ridge Regression – lets practise

In a ridge regression model, the original sum of squared residuals is 60. If the regularization parameter  $\lambda$  is set to 0.4, and the sum of squared residuals after ridge regression becomes 50, what is the proportion of variance explained by the model?



# Ridge Regression



Ridge Regression – The constraint representation of the equation...

↳ The loss function of RR  $\Rightarrow \alpha \Rightarrow \min \left( \sum_{i=1}^N \frac{(y_i - \hat{y}_i)^2}{2} + \frac{\lambda}{2} \sum_{i=1}^p \beta_i^2 \right)$

Can be written as a Constraint minimization problem

$$\Rightarrow \min \left( \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \right)$$

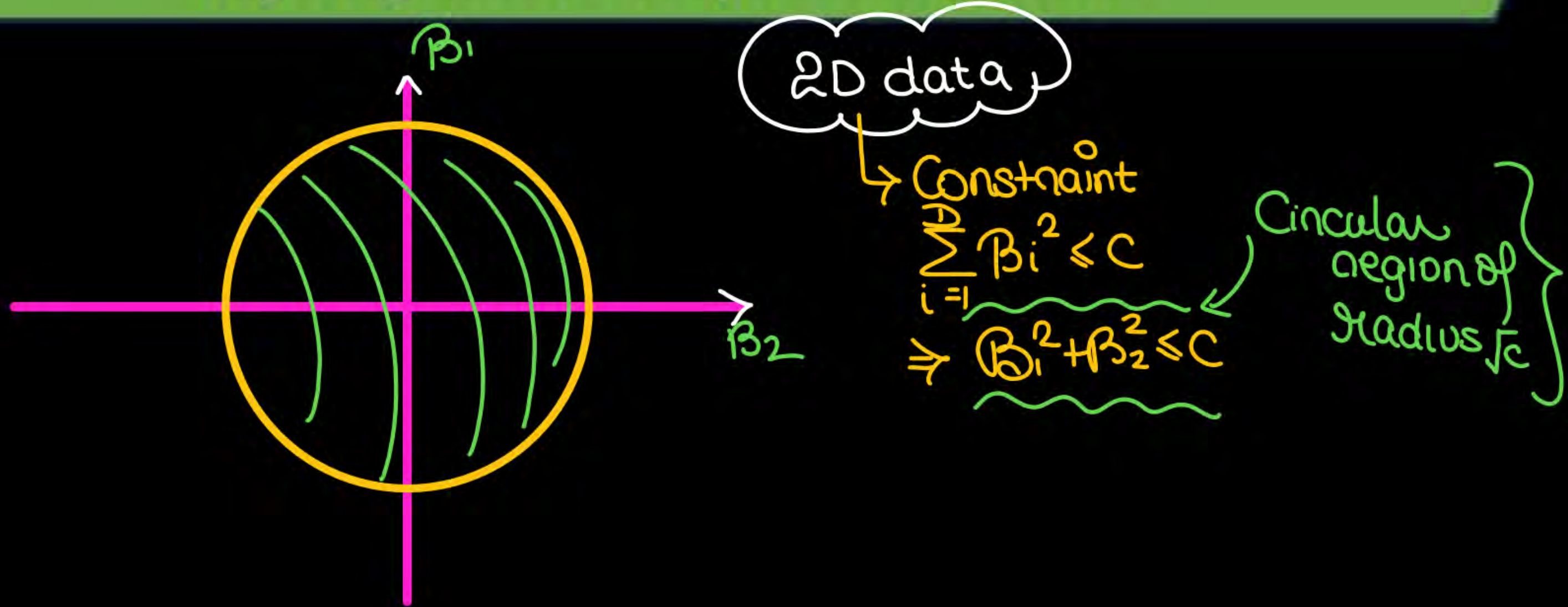
$$\text{Such that } \sum_{i=1}^p \beta_i^2 \leq c$$

- How to find best 'c'  
 $\Rightarrow$  By Cross Validation





## Ridge Regression – Shape of the constraint





# Ridge Regression



## Ridge Regression – lets practise

Numerical

2D data

$x^1$	$x^2$	$y$
5.	6.	10
8.	9	15
15.	10	20
4.	5	9

Centering of data

$$\frac{1}{5} \Rightarrow 13.5$$

$$\frac{1}{5} \Rightarrow 8.$$

$$\frac{1}{5} \Rightarrow 7.5$$

$x^1$	$x^2$	$y$
-3	-1.5	-3.5
0	1.5	1.5
7	2.5	6.5
-4	-2.5	-4.5





# Ridge Regression



## Ridge Regression – lets practise

$$X \Rightarrow \begin{bmatrix} -3 & -1.5 \\ 0 & 1.5 \\ 7 & 2.5 \\ -4 & -2.5 \end{bmatrix}$$

$$(X^T X) = \begin{bmatrix} -3 & 0 & 7 & -4 \\ -1.5 & 1.5 & 2.5 & -2.5 \end{bmatrix} \begin{bmatrix} -3 & -1.5 \\ 0 & 1.5 \\ 7 & 2.5 \\ -4 & -2.5 \end{bmatrix}$$

$$(X^T X) = \begin{pmatrix} 74 & 32 \\ 32 & 17 \end{pmatrix}, X^T Y \Rightarrow \begin{bmatrix} -3 & 0 & 7 & -4 \\ -1.5 & 1.5 & 2.5 & -2.5 \end{bmatrix} \begin{bmatrix} -3.5 \\ 1.5 \\ 6.5 \\ -4.5 \end{bmatrix}$$
$$= \begin{bmatrix} 74 \\ 35 \end{bmatrix}$$

So let  $\lambda = 1$

$$\text{So } \beta = (X^T X + \lambda I)^{-1} (X^T Y)$$

$$\Rightarrow \frac{1}{326} \begin{pmatrix} 18 & -32 \\ -32 & 75 \end{pmatrix} \begin{pmatrix} 74 \\ 35 \end{pmatrix}$$

$$\Rightarrow \frac{1}{326} \begin{pmatrix} 212 \\ 257 \end{pmatrix}$$

$$X^T X + \lambda I \Rightarrow \begin{pmatrix} 74 & 32 \\ 32 & 17 \end{pmatrix} + \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} 75 & 32 \\ 32 & 18 \end{pmatrix}$$

$$\Rightarrow \frac{1}{326} \begin{pmatrix} 18 & -32 \\ -32 & 75 \end{pmatrix}$$



$$\text{So } \beta_1 = \frac{212}{326}, \beta_2 = \frac{257}{326}$$

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

$$\Rightarrow 13.5 - 8 \times \frac{212}{326} - 7.5 \times \frac{257}{326}$$

$$\beta_0 \Rightarrow 2.384$$



## Ridge Regression



### What is Lasso Regularisation

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

- Square term in RR  $\rightarrow$  L2 Regularisation
- In lasso  $\rightarrow$  Single power term  $\rightarrow$  L1 Regularisation.





## Ridge Regression

## Lasso Vs Ridge Regression



Parameter	Ridge Regression	Lasso Regression
<b>Regularization Type</b>	L2 regularization: adds a penalty equal to the square of the magnitude of coefficients.	L1 regularization: adds a penalty equal to the absolute value of the magnitude of coefficients.
<b>Primary Objective</b>	To shrink the coefficients towards zero to reduce model complexity and multicollinearity.	To shrink some coefficients towards zero for both variable reduction and model simplification.
<b>Feature Selection</b>	Does not perform feature selection: all features are included in the model, but their impact is minimized.	Performs feature selection: can completely eliminate some features by setting their coefficients to zero.
<b>Coefficient Shrinkage</b>	Coefficients are shrunk towards zero but not exactly to zero.	Coefficients can be shrunk to exactly zero, effectively eliminating some variables.
<b>Suitability</b>	Suitable in situations where all features are relevant, and there is multicollinearity.	Suitable when the number of predictors is high and there is a need to identify the most significant features.
<b>Bias and Variance</b>	Introduces bias but reduces variance.	Introduces bias but reduces variance, potentially more than Ridge due to feature elimination.
<b>Interpretability</b>	Less interpretable in the presence of many features as none are eliminated.	More interpretable due to feature elimination, focusing on significant predictors only.
<b>Sensitivity to <math>\lambda</math></b>	Gradual change in coefficients as the penalty parameter $\lambda$ changes.	Sharp thresholding effect where coefficients can abruptly become zero as $\lambda$ changes.
<b>Model Complexity</b>	Generally results in a more complex model compared to Lasso.	This leads to a simpler model, especially when irrelevant features are abundant.



## 2 mins Summary



Topic

Topic

Topic

Topic

Topic



**THANK - YOU**