

Data Science and Artificial Intelligence

Machine Learning



Regression

Lecture No. 8

By- SIDDHARTH SABHARWAL SIR

GATE WALLAH

Recap of Previous Lecture



Topic

Assumption in LR

Topic

P.W

Topic

Topic

Topic

Topics to be Covered



Topic

Sol. uploaded

Topic

Adv and dis adv of LR

Topic

Space and time Complexity of LR

Topic

Ridge Regression

Topic



**Be the change
that you wish to
see in the world.**

— MAHATMA GANDI



• advantage of LR } \Rightarrow

① Simplest algorithm

② High interpretability

$$y = 3 + 2x' + 3x^2 + 10x^3$$

• from the eq. of model

• we can interpret

* importance of each dimension

* It is possible to interpret the plane created from model



Problems in LR ...

- disadvantage in LR →
- ① it works only when data is linearly related
 - ② data must not have multicollinearity
 - ③ data must not have heteroscedasticity
 - ④ LR produces unstable model
- ⑤ algorithm is severely affected by outlier
- (Note: An orange line connects the 'V' in the circled 5 to the circled 2, 3, and 4.)*



Linear Regression



Advantages of Simple Linear Regression:

- Simplicity and ease of interpretation. ✓
- Transparent modeling with clear coefficient interpretations. ✓
- Computational efficiency, suitable for large datasets. → Simple algo
- A baseline model for assessing feature significance. → Coefficient of dimension tell significance of dimension.
- Effective when the relationship between variables is linear.

Disadvantages of Simple Linear Regression:

- ✓ Limited to linear relationships, may perform poorly for nonlinear data.
- ✓ Sensitive to outliers, leading to parameter influence.
- ✓ Prone to underfitting when facing complex relationships. underfit if data is NL.
- ✓ Assumptions of independent and normally distributed errors are critical.
- Suitable only when one independent variable is involved in the analysis.

- Assumptions of independent and normally distributed errors are critical.
- Suitable only when one independent variable is involved in the analysis.

Homoscedasticity

Gaussian noise



- LR has one more assumption
 - that data points must be independent to each other.



Linear Regression



Advantage & Disadvantage of Linear Regression

PW

Advantages	Disadvantages
Linear Regression is simple to implement and easier to interpret the output coefficients.	On the other hand in linear regression technique outliers can have huge effects on the regression and boundaries are linear in this technique.
When you know the relationship between the independent and dependent variable have a linear relationship, this algorithm is the best to use because of it's less complexity compared to other algorithms.	Diversely, linear regression assumes a linear relationship between dependent and independent variables. That means it assumes that there is a straight-line relationship between them. It assumes independence between attributes.
Linear Regression is susceptible to over-fitting but it can be avoided using some dimensionality reduction techniques, regularization (L1 and L2) techniques and cross-validation.	But then linear regression also looks at a relationship between the mean of the dependent variables and the independent variables. Just as the mean is not a complete description of a single variable, linear regression is not a complete description of relationships among variables.



Linear Regression



Why Linear Regression is Important

The interpretability of linear regression is a notable strength.

The model's equation provides clear coefficients that elucidate the impact of each independent variable on the dependent variable, facilitating a deeper understanding of the underlying dynamics.

Its simplicity is a virtue, as linear regression is transparent, easy to implement, and serves as a foundational concept for more complex algorithms.



Why LR is unstable model...

$|A| = 0.002$ $A = \begin{bmatrix} 1 & 2 \\ 1.000 & 2.002 \end{bmatrix} \Rightarrow \text{invertible } A^{vv}.$
 $\rightarrow |A| = \text{very small} = \text{close to } 0$

$\sim A^{-1} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ So terms in A^{-1} matrix will be
v. large.

$|A| = 0.001$ $\begin{bmatrix} 1 & 2 \\ 1.000 & 2.001 \end{bmatrix}$

$$X = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 4 \\ 1 & 3 & 6 \\ 1 & 4 & 8 \end{bmatrix}$$

if X has multicollinearity
 $x^2 = 2x^1$

$$X^T X = \begin{bmatrix} \text{2 on more than 2 rows equal} \end{bmatrix}$$

- $|X^T X| = 0$
- $X^T X \Rightarrow$ non invertible

$$X = \begin{bmatrix} 1 & 1 & 0.9999 \\ 1 & 2 & 1.9889 \\ 1 & 3 & 3.0002 \\ 1 & 4 & 4.0001 \end{bmatrix}$$

$$(X^T X) = \begin{bmatrix} \text{Row exactly same nahi} \\ \text{Row exactly same nahi} \end{bmatrix}$$

$$|X^T X| \neq 0, \beta = (X^T X)^{-1} X^T y$$

Slight change
in data

$$X = \begin{bmatrix} 1 & 1 & 0.9999 \\ 1 & 2 & 1.9889 \\ 1 & 3 & 3.0002 \\ 1 & 4 & 4.0001 \end{bmatrix}$$

$(X^T X)^{-1} X^T Y$
we get large
 β 's

$$X = \begin{bmatrix} 1 & 1 & 0.9989 \\ 1 & 2 & 1.9998 \\ 1 & 3 & 3.0001 \\ 1 & 4 & 4.0002 \end{bmatrix}$$

$\beta = (X^T X)^{-1} X^T Y$
 β 's will be of
large value
But

β 's value will be totally
different.

- $(X^T X)^{-1} X^T Y$

Since $(X^T X)^{-1}$, all the terms will be vry. large.

$\beta \Rightarrow$ we will get very large β values.

→ So in LR due to multicollinearity we may get large β 's.

- LR give unstable models



Slight change in data

Produces Completely new model



Space and Time Complexity of Linear Regression

Assumptions:

n = number of training examples, m = number of features, n' = number of support vectors,
 k = number of neighbors, k' = number of trees

- **Linear Regression**

- Train Time Complexity = $O(n \cdot m^2 + m^3)$
- Test Time Complexity = $O(m)$
- Space Complexity = $O(m)$

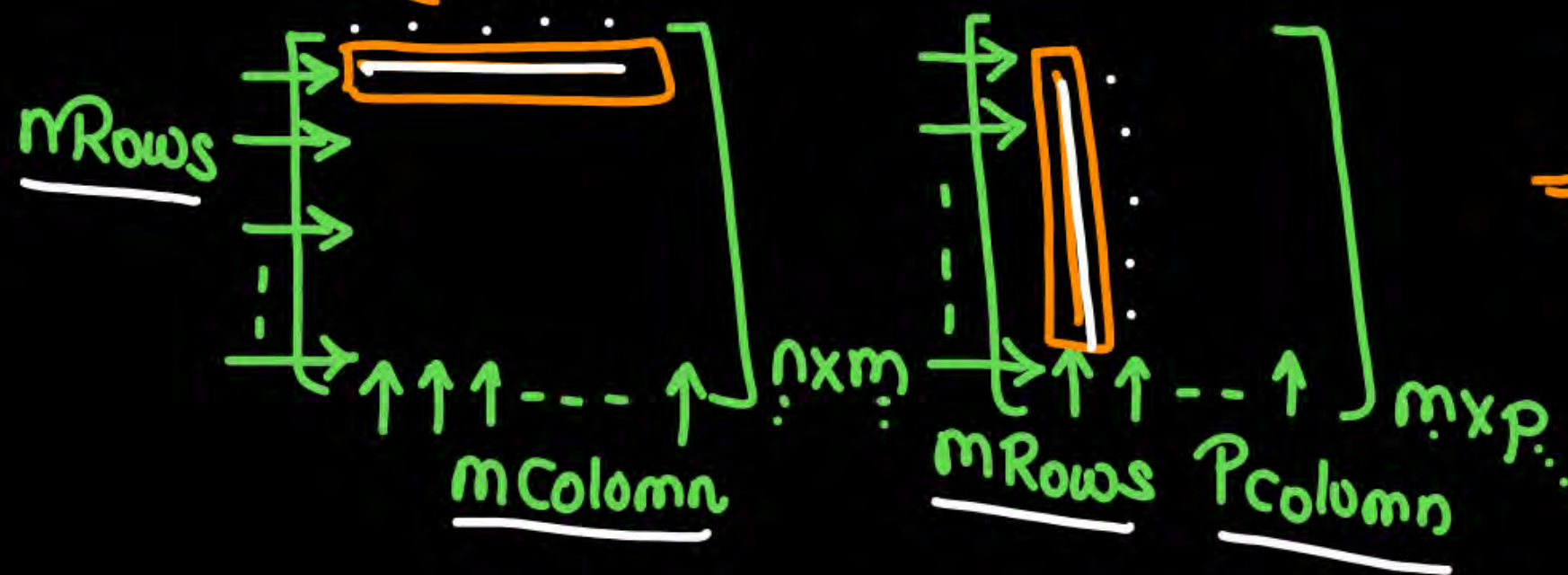


Linear Regression

Space and Time Complexity of Linear Regression

- matrix multiplication

$$(A)_{n \times m} (B)_{m \times p} = (C)_{n \times p}$$



For each value m -times mult,
($m-1$)-times addition
 $n \times p$ number of values
 $(n \times p \times m)$ number of mult
 $n \times p \times (m-1)$ number of add.

Inverse of a matrix

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}_{3 \times 3}$$

$$\begin{bmatrix} C_{11} & C_{12} & C_{13} \\ C_{21} & C_{22} & C_{23} \\ C_{31} & C_{32} & C_{33} \end{bmatrix}$$

$$A^{-1} = \frac{1}{|A|} \begin{bmatrix} \text{Adj } A \end{bmatrix}$$

$$C_{11} = (e_i - f_h)(-1)^{i+h}$$

$$C_{12} = (-1)^{1+2} (d_i - f_g)$$

To find Cofactor = 3 Mult

9 Cofactors \Rightarrow 27 Mult

$\rightarrow (\text{Cof Matrix})^T$

* So for $(3 \times 3)^{-1} \Rightarrow$ we need 27 operation

N : No of data points
 D : dimension of data

$K \times K$ matrix Ka inv \Rightarrow Number of Calculation $\Rightarrow O(K^3)$.

In LR \Rightarrow

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

$(X^T)_{(D+1) \times N} (Y)_{N \times 1} \Rightarrow$ To create $X^T Y$
 $(D+1) \times 1$

$$X = \begin{bmatrix} \quad \end{bmatrix}_{N \times (D+1)}$$

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

To create $X^T X \Rightarrow (D+1)^2 \times N$

\Rightarrow Inv $\Rightarrow (D+1)^3$

$(X^T X)^{-1}_{(D+1) \times (D+1)} (X^T Y)_{(D+1) \times 1} \Rightarrow$ Total op. $(D+1)^3$.

N : No of data points
 D : dimension of data

In LR \Rightarrow

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

$$(X^T)_{(D+1) \times N} (Y)_{N \times 1} \Rightarrow \text{To create } X^T Y \quad (D+1) \times 1$$

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

$$\text{To create } X^T X \Rightarrow (D+1)^2 \times N$$

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

$$\Rightarrow \text{Inv} \Rightarrow (D+1)^3$$

$$(X^T X)^{-1}_{(D+1) \times (D+1)} (X^T Y)_{(D+1) \times 1} \Rightarrow \text{Total op. } (D+1)^2$$

Total No of Calculation $\Rightarrow (D+1)^2 \times N + (D+1)^3 + N(D+1) + (D+1)^2$
 This is order of No of Calci \Rightarrow
 Not exact No of Calci $\left[N(D+1)^2 + (D+1)^3 \right]$

$$K = D + 1$$

• Training time \Rightarrow

\rightarrow Time taken to produce β_0 .
if 1 unit time is needed for
a multiplication

we need $O(NK^2 + K^3)$ number of mult.

So training time $\Rightarrow O(NK^2 + K^3)$

data \rightarrow

LR
machine

\rightarrow

@end of
training
we get.
model, β values

Training Time
Complexity of LR
 $\Rightarrow O(NK^2 + K^3)$

after raining



we only need to store \Rightarrow Only β 's.

$$D+1 = K$$

Space Complexity.

$\rightarrow (D+1)$
Spaces.

$\rightarrow (D+1)$ Beta values

How testing is done

we get a new x value

$$[x_t^1 \ x_t^2 \ x_t^3 \ x_t^4 \ \dots \ x_t^D]$$

To find \hat{y} , testing

$$\begin{bmatrix} 1 & x_t^1 & x_t^2 & \dots & x_t^D \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix} \Rightarrow \hat{y}$$

@ Testing $(D+1)$ multiplication

Order of Complexity in
Testing = $O(K)$



Ridge Regression



Problem with the least square error ??

Problem in OLS
 $\min \text{RSS}$

① Multicollinear data \Rightarrow produce unstable model

② $\text{LR} = \min \sum_{i=1}^N (y_i - \hat{y}_i)^2$ very large β .

\rightarrow we $\min \text{RSS} \Rightarrow$ thus algo want to make $\text{RSS} = 0$.

\rightarrow So algo has overfitting tendency.

@ The method to solve the problem

@ Regularisation Ridge Regression

we change the loss function

updated
loss f_{xn}

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

RSS

= we will minimize this loss f_{xn}

Ex.

Same data. $\xrightarrow{\text{LR}}$

$$x^2 = 100x^1$$

x^1	x^2	y
-	-	-
-	-	-
-	-	-
-	-	-

$$3x^2 + 1x^1 + 100 = 0$$

Best model

$$301x^1 + 0x^2 + 100 = 0$$

$$3.01x^2 + 0x^1 + 100 = 0$$

all 3 perform
Similar
Similar MSE on data

Ex.

Same data $\xrightarrow{\text{LR}}$

$$x^2 = 100x^1$$

x^1	x^2	y
-	-	-
-	-	-
-	-	-
-	-	-

$\downarrow \text{LR}$

~~$3x^2 + 1x^1 + 100 = 0$~~

Best model

~~$301x^1 + 0x^2 + 100 = 0$~~

$3.01x^2 + 0x^1 + 100 = 0$ ✓

① Solution of Ridge Regression remove problems of large β

② The algorithm will make β 's of irrelevant dimension = 0 or close to zero.



Ridge Regression



Problem with the least square error ??

- Not all the dimensions are equally usefull





Ridge Regression



Problem with the least square error ??

- OLS may lead to unstable model ✓

$\min \text{RSS}$



Ridge Regression



Problem with the least square error ??

Let's Summarize : The problem in OLS or minimizing the RSS are as follows

1. It may lead to Overfitting.
2. No boundation on values of Betas may lead to unstable model.
3. The problem of multicollinearity.

So what is the solution

Regularisation



Ridge Regression



Shrinkage Methods : Ridge Regression

❖ Ridge regression is a regularisation techniques...



Ridge Regression



Shrinkage Methods : Ridge Regression

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



Ridge Regression



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.

The loss
function are
updated



Ridge Regression



Shrinkage Methods : Ridge Regression

The loss
function are
updated



Ridge Regression



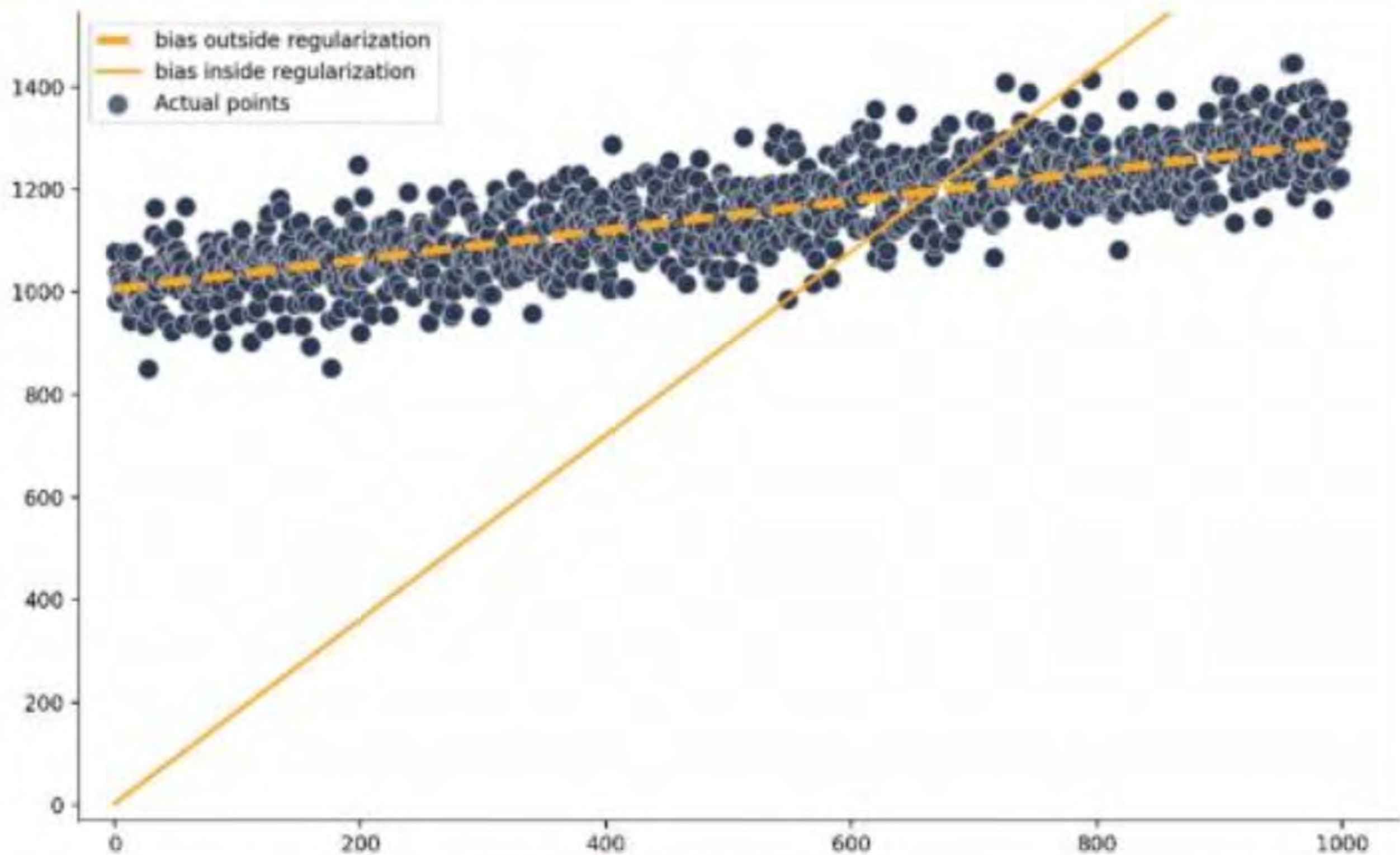
Shrinkage Methods : Ridge Regression

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 bias term is not included in regularisation ..



Ridge Regression



This GIF has been sourced from the author's website



Ridge Regression



Shrinkage Methods : Ridge Regression

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

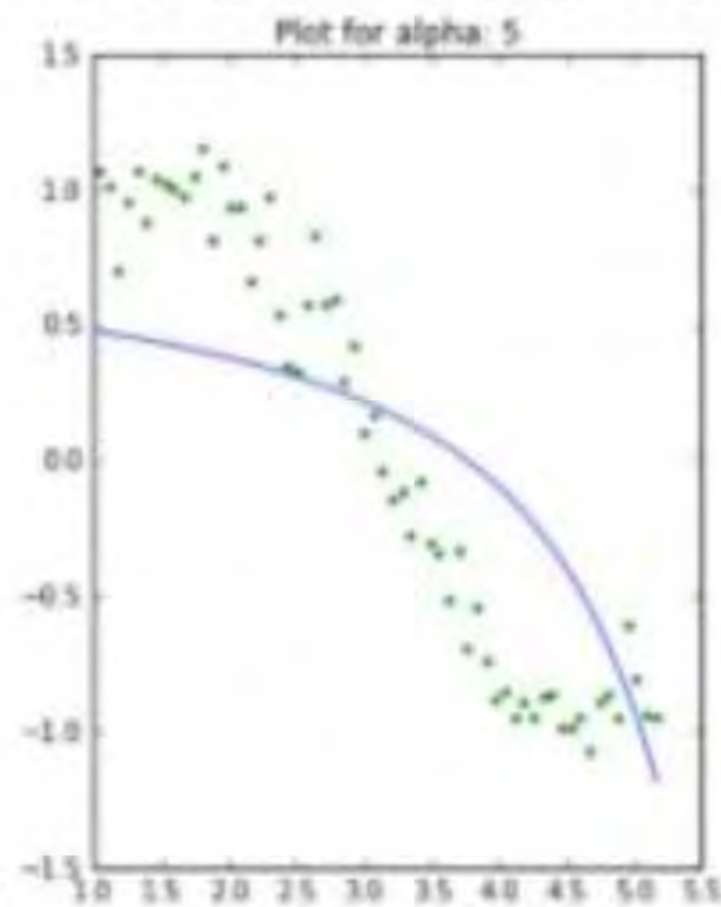
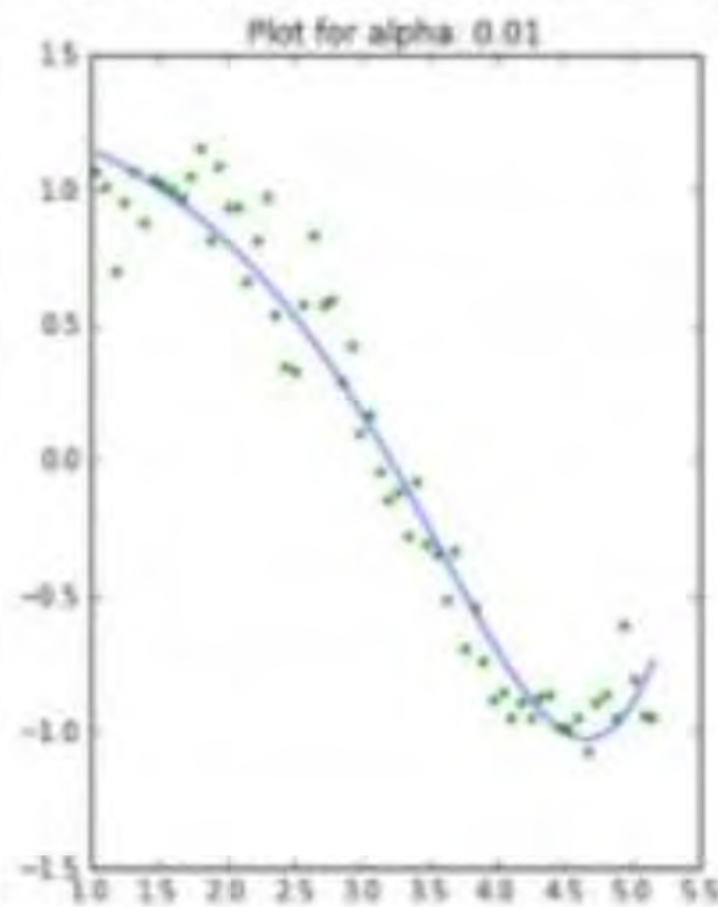
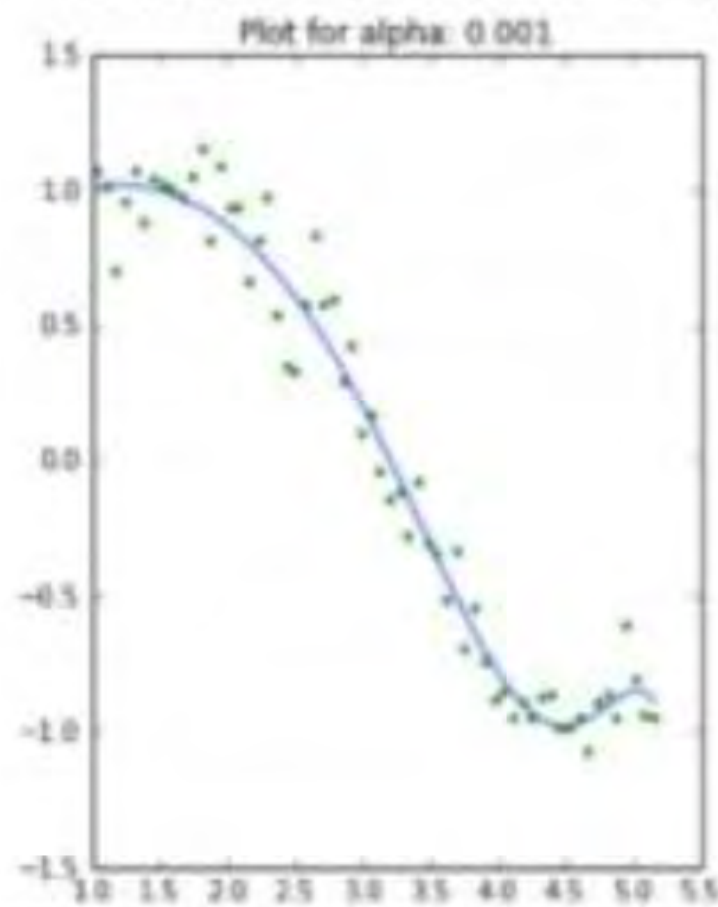
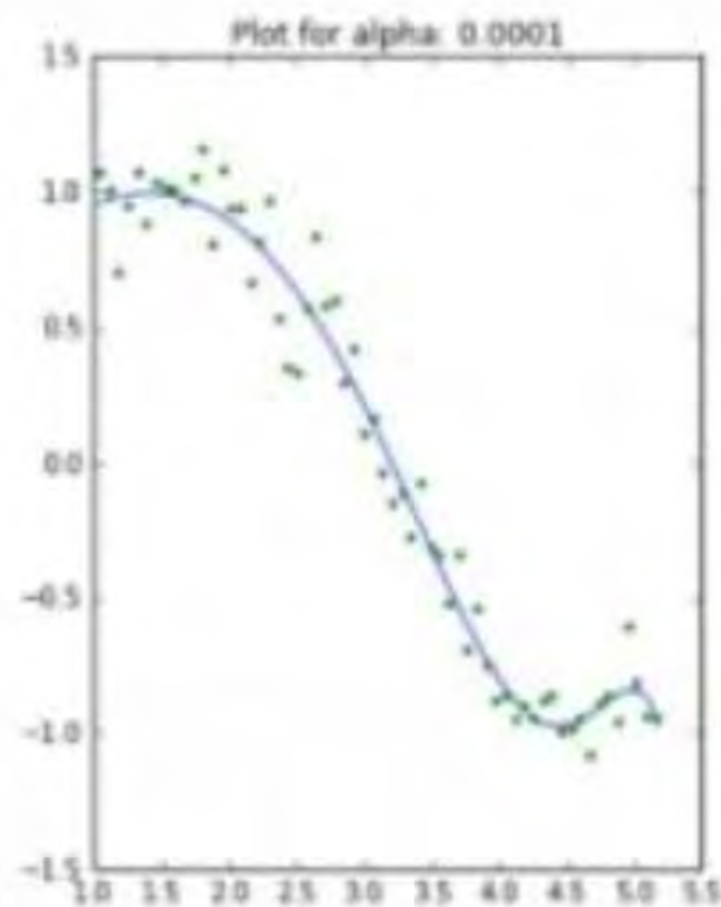
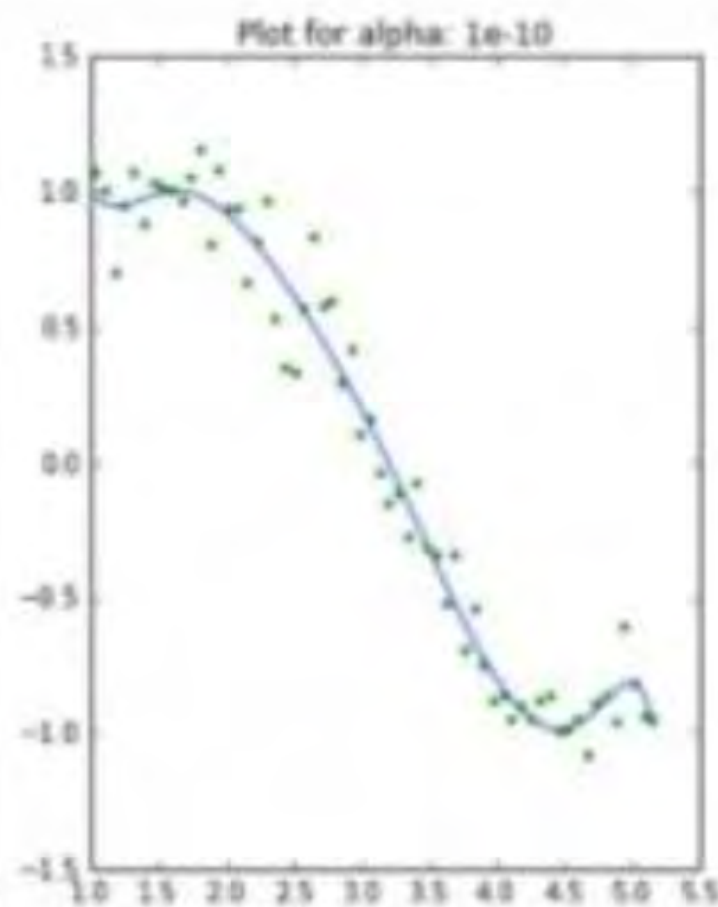
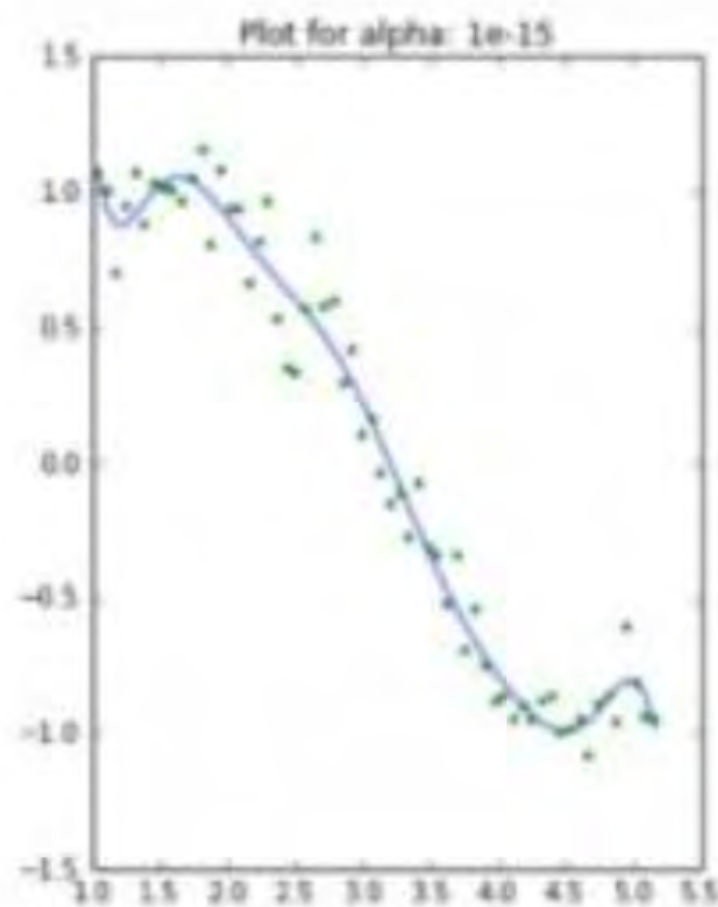


Ridge Regression



Shrinkage Methods : Ridge Regression

❖ Here λ is very important control parameter:





Ridge Regression



Shrinkage Methods : Ridge Regression

❖ Lets find the solution to this ridge regression problem



Ridge Regression



Shrinkage Methods : Ridge Regression

❖ How to find λ (can this be negative?)



Linear 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.



Linear 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 – 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.

THANK - YOU