

Data Science and Artificial Intelligence

Machine Learning



Regression

Lecture No. 10

By- SIDDHARTH SABHARWAL SIR



Recap of Previous Lecture



Topic

Topic

Topic

Topic

Topic

• Ridge Regression
expression
effect of λ

Topics to be Covered



Topic

How to find best λ

Topic

Constraint view of RR

Questions.

Topic

Lasso \rightarrow brief intro

Topic

Classification

Topic

Parade

"NOTHING IS
IMPOSSIBLE.
THE WORD
ITSELF SAYS
'I'M POSSIBLE!'"
— AUDREY HEPBURN

About the Faculty

- AIR 1 GATE 2021, 2023 (ECE).
- AIR 3 ESE 2015 ECE.
- M.Tech from IIT Delhi in VLSI.
- Published 2 papers in field of AI-ML.
- Paper 1 : Feature Selection through Minimization of the VC dimension.
- Paper 2 : Learning a hyperplane regressor through a tight bound on the VC dimension.



By- SIDDHARTH SABHARWAL SIR

Join Our Telegram Channel!

Stay updated, get instant alerts & connect with the community



<https://t.me/siddharthsirPW>



By- SIDDHARTH SABHARWAL SIR



Ridge Regression Final expression



Ridge Regression



Ridge Regression – lets practise

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

- Prevent overfitting, → reduce model Complexity
- ☒ A) Increase model complexity.
 - ☒ B) Reduce model complexity and prevent overfitting.
 - ☒ C) Make the model fit the training data perfectly. → overfit
 - ☒ 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:

$$\text{Loss}_{f_{Xn}} = \text{RSS} + \text{penalty term} \longrightarrow \sum \beta_i^2$$

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

Absolute value $\Rightarrow |\beta_i|$



Ridge Regression



Ridge Regression – lets practise

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

β Ki values.

RR \Rightarrow try to reduce β 's.

- ☒ 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:

β value small
dimensions.

Linear R = give stable model

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. → overfit

✗ D) The dataset has very few observations.

→ LR give unstable model

RR = Solve Multi
Collinearity



Ridge Regression



Ridge Regression – lets practise

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

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

Strongest penalty term.

→ λ that give v. high importance to penalty term

• $\min \text{RSS} + \lambda \sum \beta_i^2$
 $\lambda = \infty, \beta_i \approx 0$ underfit model.



Ridge Regression

RR → inc model interpretability. (adv)
inc Complexity Computation (dis).



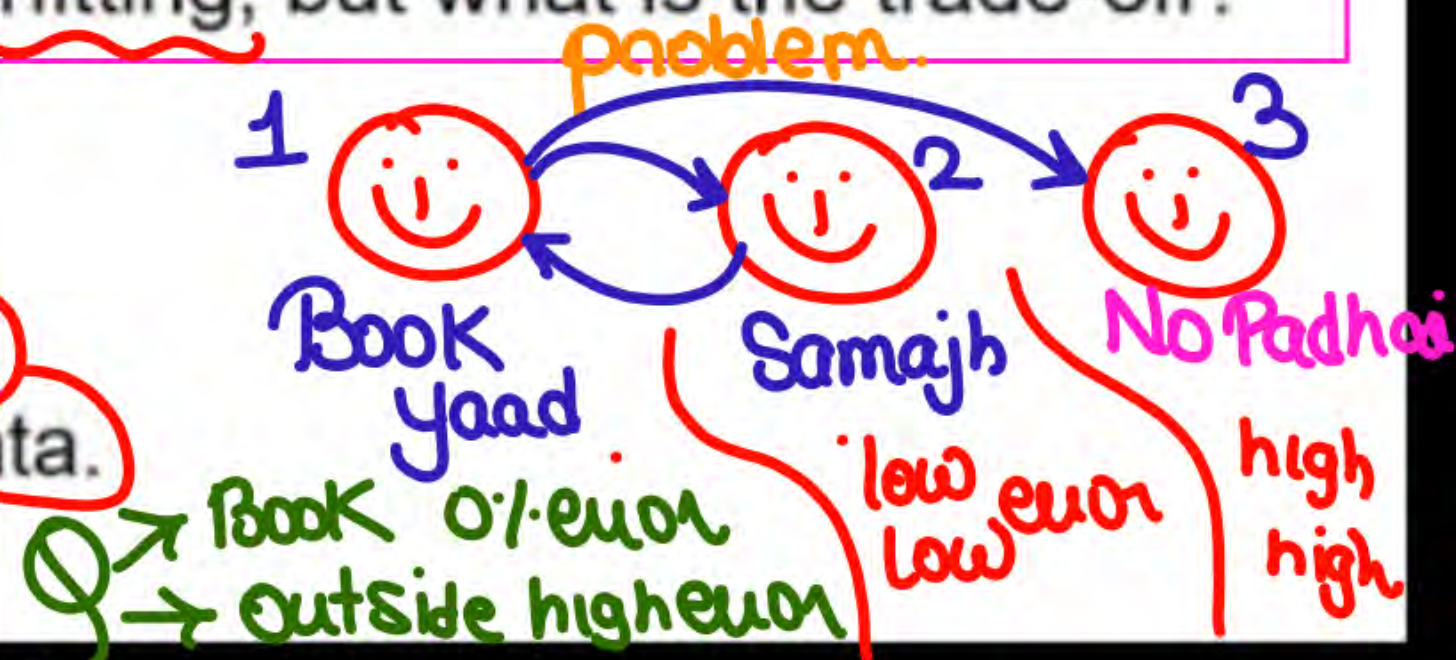
Ridge Regression – lets practise

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

Benefit

RR 1 → 2

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



overfit → Best fit

100 dimension
K_a
data.

→ LR ⇒

100 β's K_a eq.

Jis main har dimension is
given importance

→ RR ⇒

85 β's → 0 ⇒ less imp dimension

15 β's ≠ 0 ⇒ imp dimension

→

So model will give more error on
training data



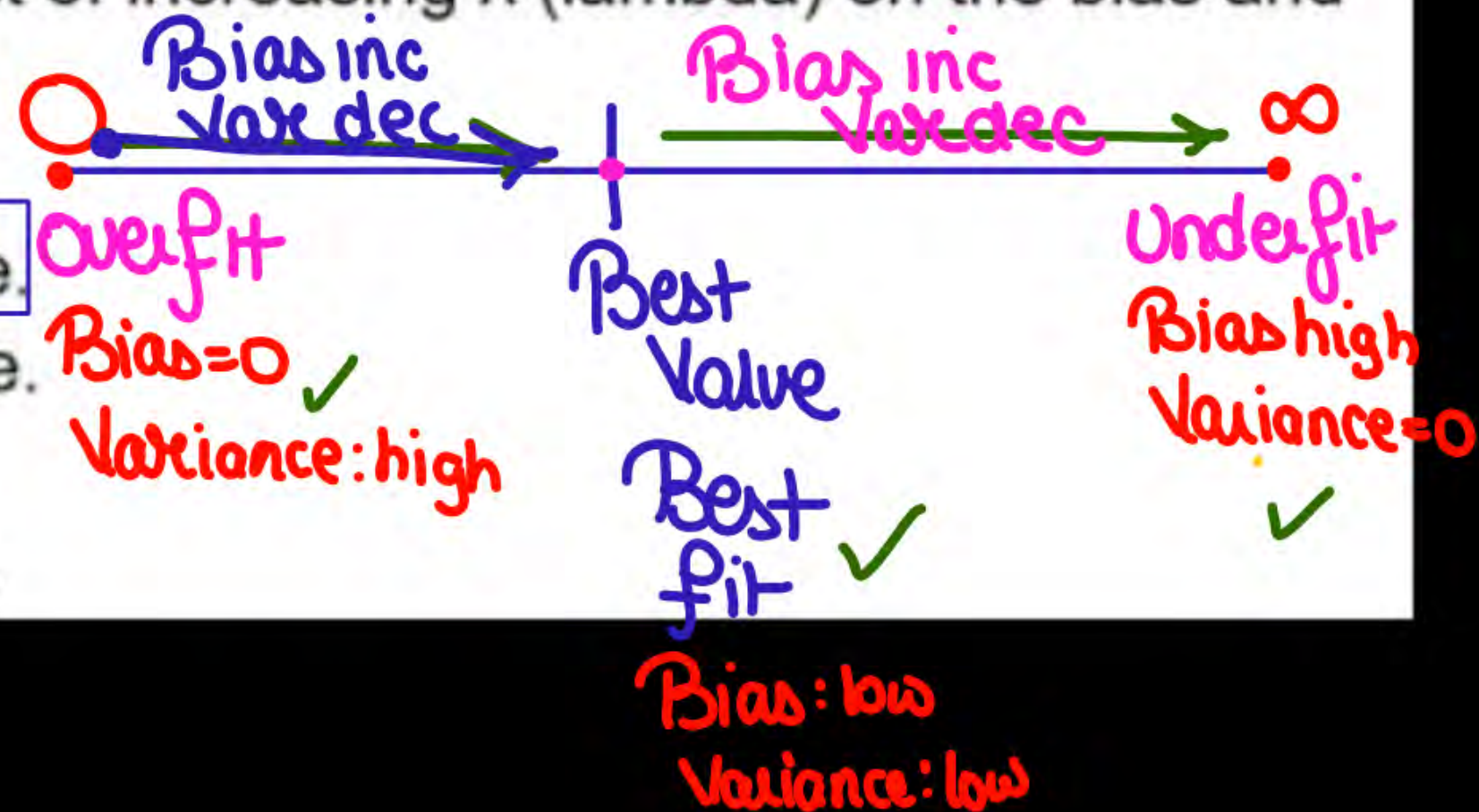
Ridge Regression



Ridge Regression – lets practise

In Ridge Regression, what is the effect of increasing λ (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 λ (lambda) is 3, what is the modified penalty term in the Ridge Regression cost function?

a) 150 ✓

b) 135

c) 123

d) 578

$$\begin{aligned}\text{Penalty term} &= \lambda \sum \beta^2 \\ &= 3 \times 50 \\ &= 150\end{aligned}$$



Ridge Regression

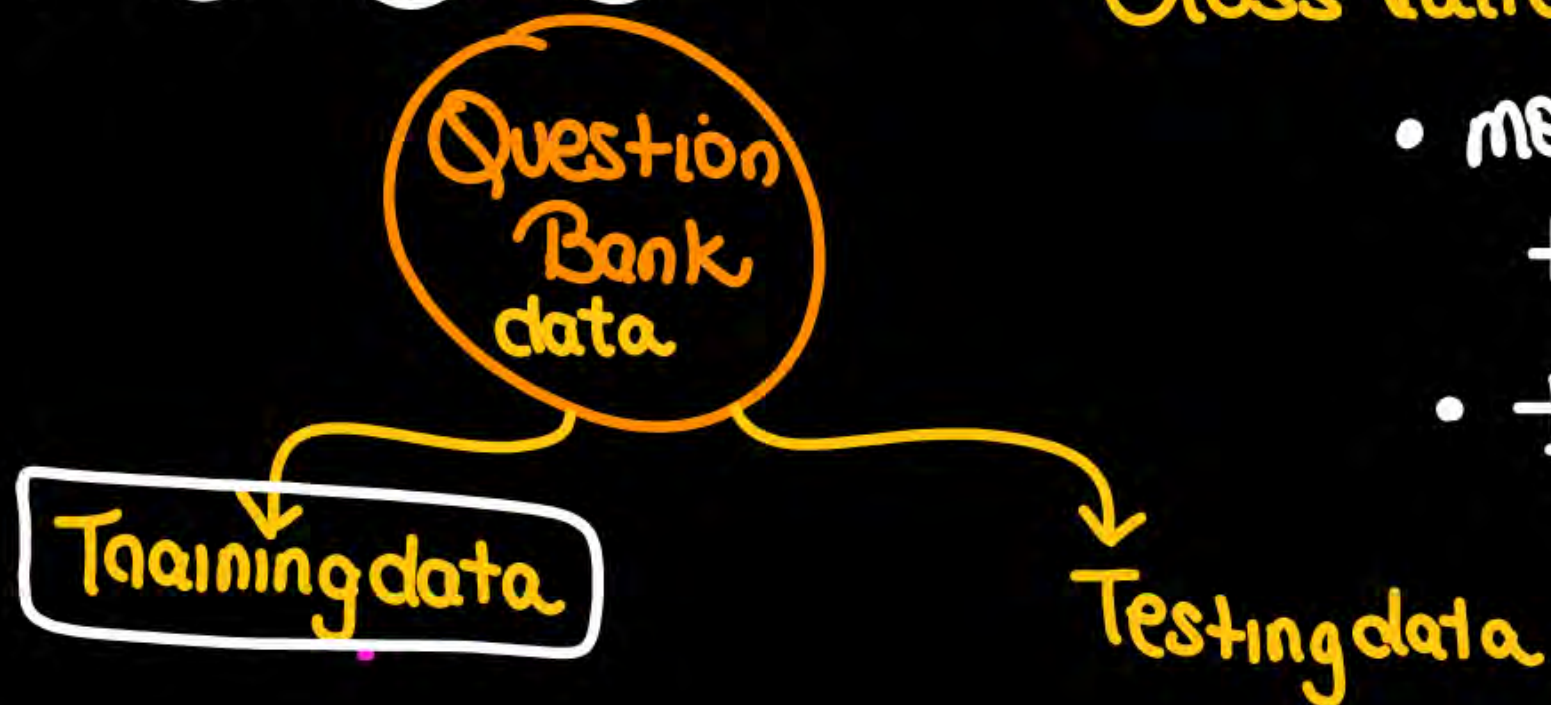


How to find the best hyperparameter

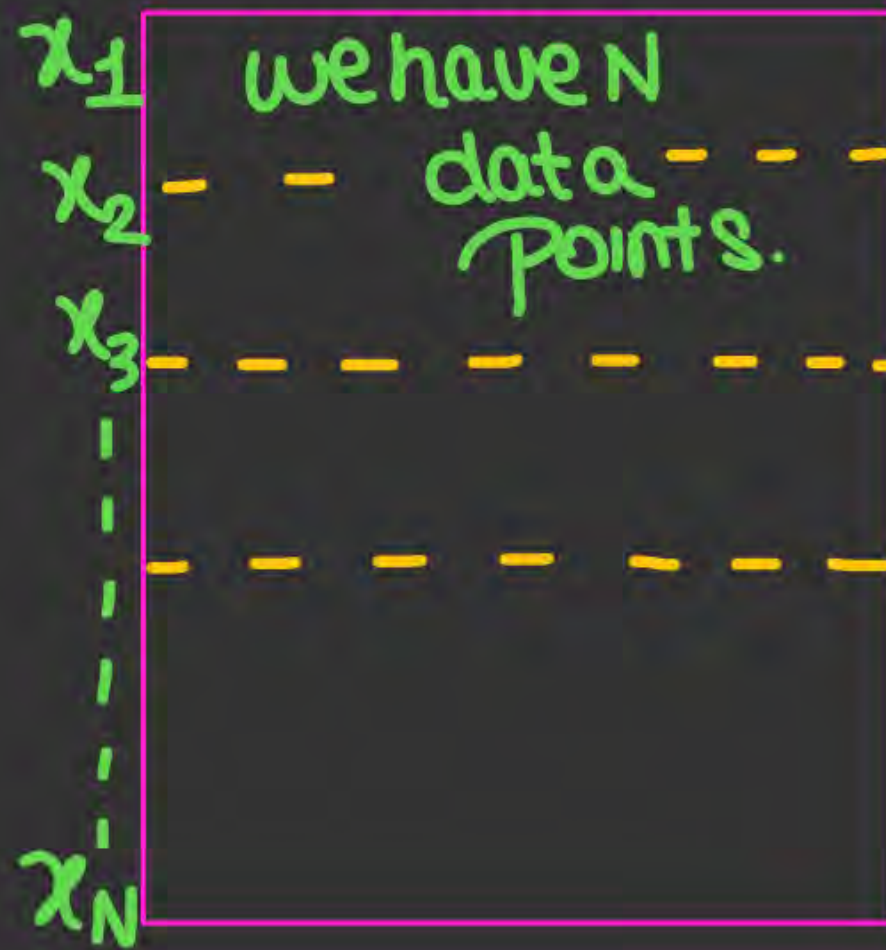
How to find best λ

The method is called
Cross validation \Rightarrow

- model is prepared using training data
- then accuracy is tested on Testing data.



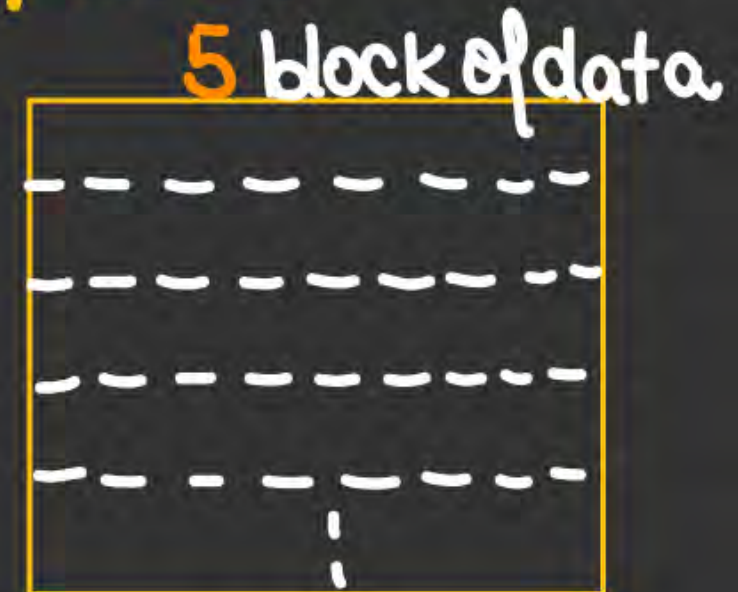
Training data.



K fold Cross validation
→ Step ① divide training data into K folds.

Each block/fold has $\frac{N}{K}$ data points.

Step 2 : let $K = 5$



K=5, 5 fold CV.

		Toain	Validate	%accuracy
we choose λ $\lambda = 0.1$	1	1,2,3,4	5	a
	2	2,3,4,5	1	b
	3	1,3,4,5	2	c
	4	1,2,4,5	3	d
	5	1,2,3,5	4	e

for $\lambda = 0.1$

accuracy =

$$\frac{a+b+c+d+e}{5}$$

So by this we check that model prepared by $\lambda = 0.1$, is good for overall data or not

K=5, 5 fold CV.

Now

		Toain	Validate	%accuracy
we choose λ $\lambda = 0.2$	1	1,2,3,4	5	a
	2	2,3,4,5	1	b
	3	1,3,4,5	2	c
	4	1,2,4,5	3	d
	5	1,2,3,5	4	e

For $\lambda = 0.1$

accuracy =

$$\frac{a+b+c+d+e}{5}$$

So by this we check that model prepared by $\lambda = 0.1$ is good for overall data or not

K=5, 5 fold CV.

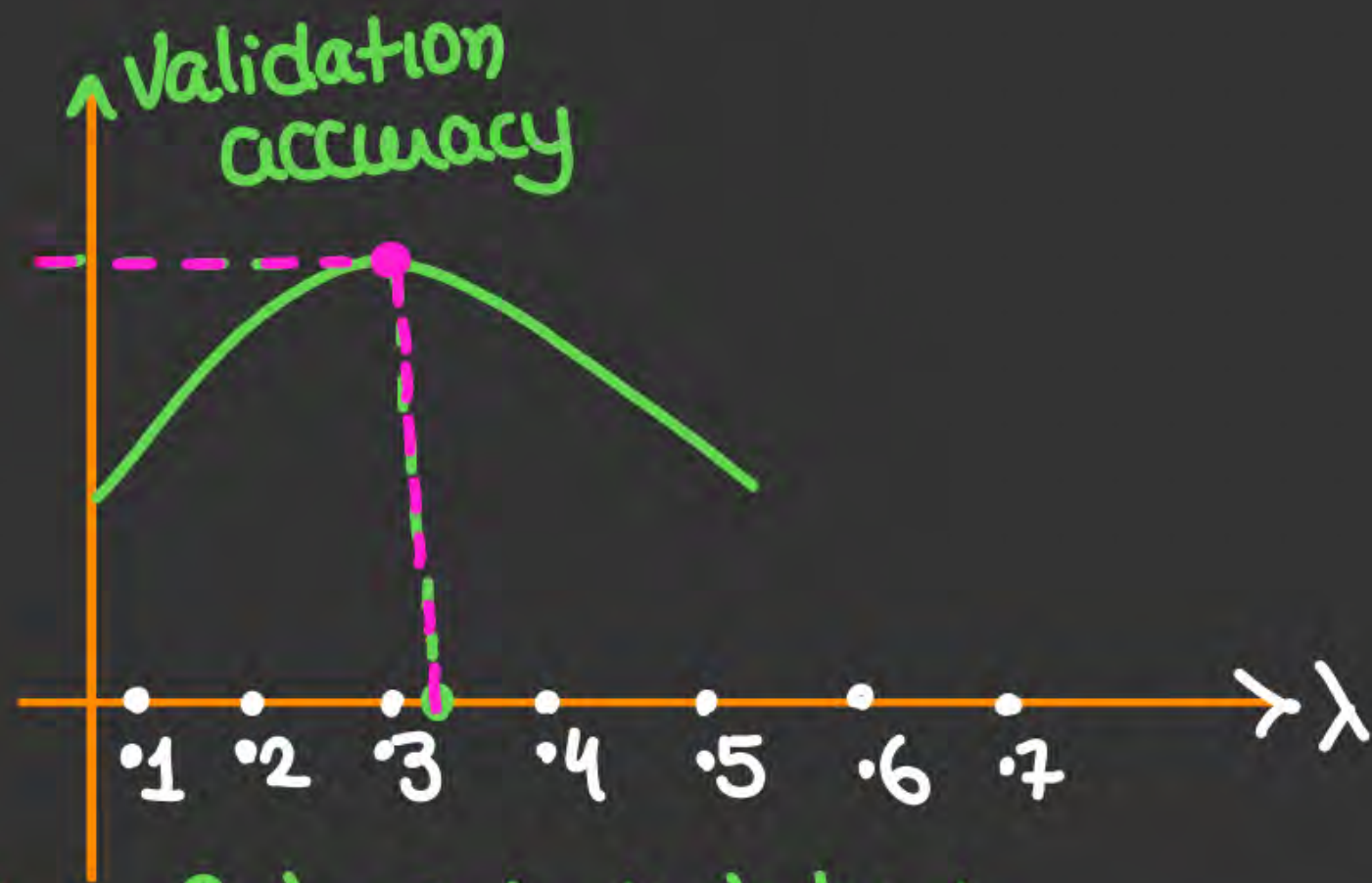
Now

		Toain	Valdate	%accuracy
we choose λ $\lambda = 0.3$	1	1,2,3,4	5	a
	2	2,3,4,5	1	b
	3	1,3,4,5	2	c
	4	1,2,4,5	3	d
	5	1,2,3,5	4	e

for $\lambda = 0.1$

$$\text{accuracy} = \frac{a+b+c+d+e}{5}$$

So by this we check that model prepared by $\lambda = 0.1$ is good for overall data or not



Finding best λ is really very lengthy.

So $\lambda = 0.3$ to 0.4 is best
Repeat process for $\lambda = 0.31$
 $\lambda = 0.32$
 $\lambda = 0.33$
!
To find best λ .

• K fold CV

• This process is used to find best \rightarrow

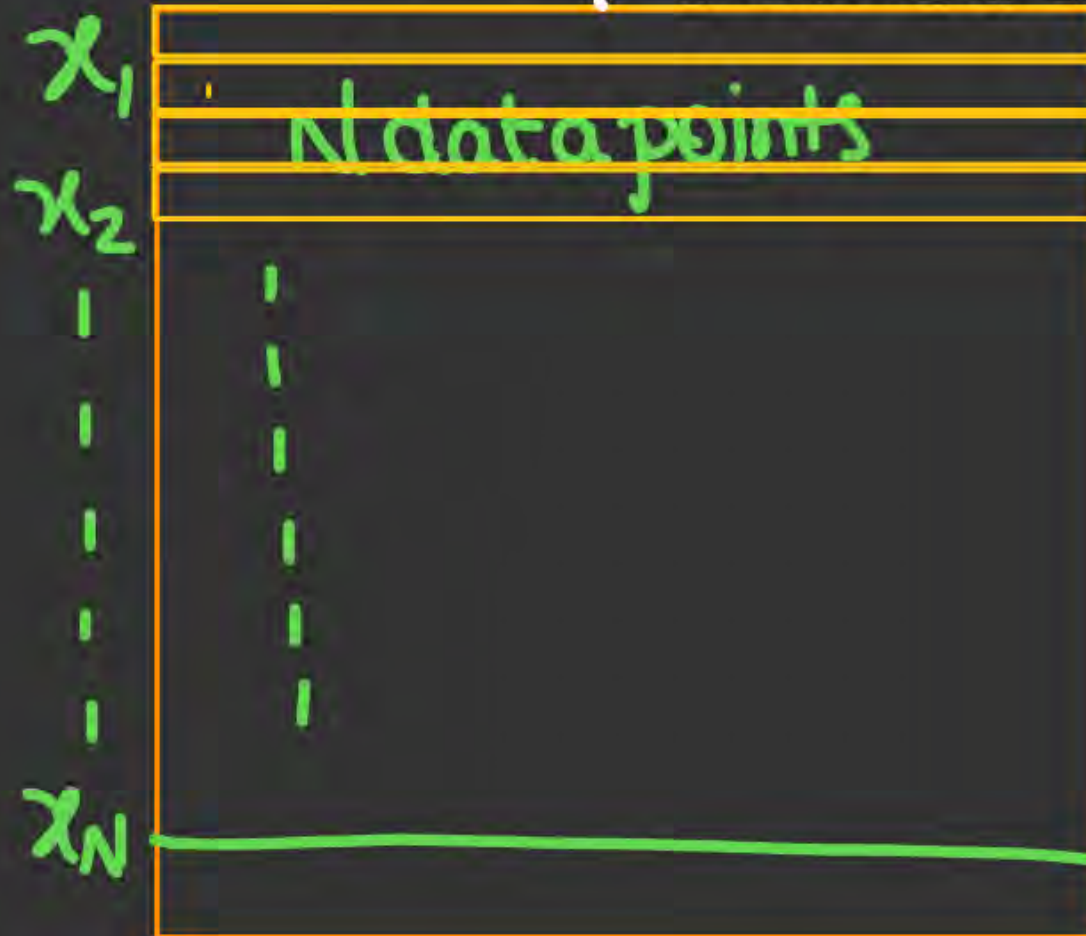
① divide training data into K fold

② 1 fold for validation, $(K-1)$ fold for training

\rightarrow Repeat process K times.

• LOOCV \Rightarrow leave one out Cross Validation.

\rightarrow we break data into N folds.



each fold has 1 data point
@ each step $(N-1)$ fold Training
1 fold Testing

\rightarrow Process Repeat $\sim N$ No of times.

Ex 1000 data points

↓
LOOCV

↓
1000 folds

↳ 1000 times process
Repeat

1 fold Validation
999 fold training.



Ridge Regression



Constraint representation of Ridge Regression

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

minimize loss function

• another way \Rightarrow

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

Such that

$$\sum_{i=1}^D \beta_i^2 \leq C$$

\hookrightarrow Constant

- C will be a hyper parameter, which will be calculated by CV



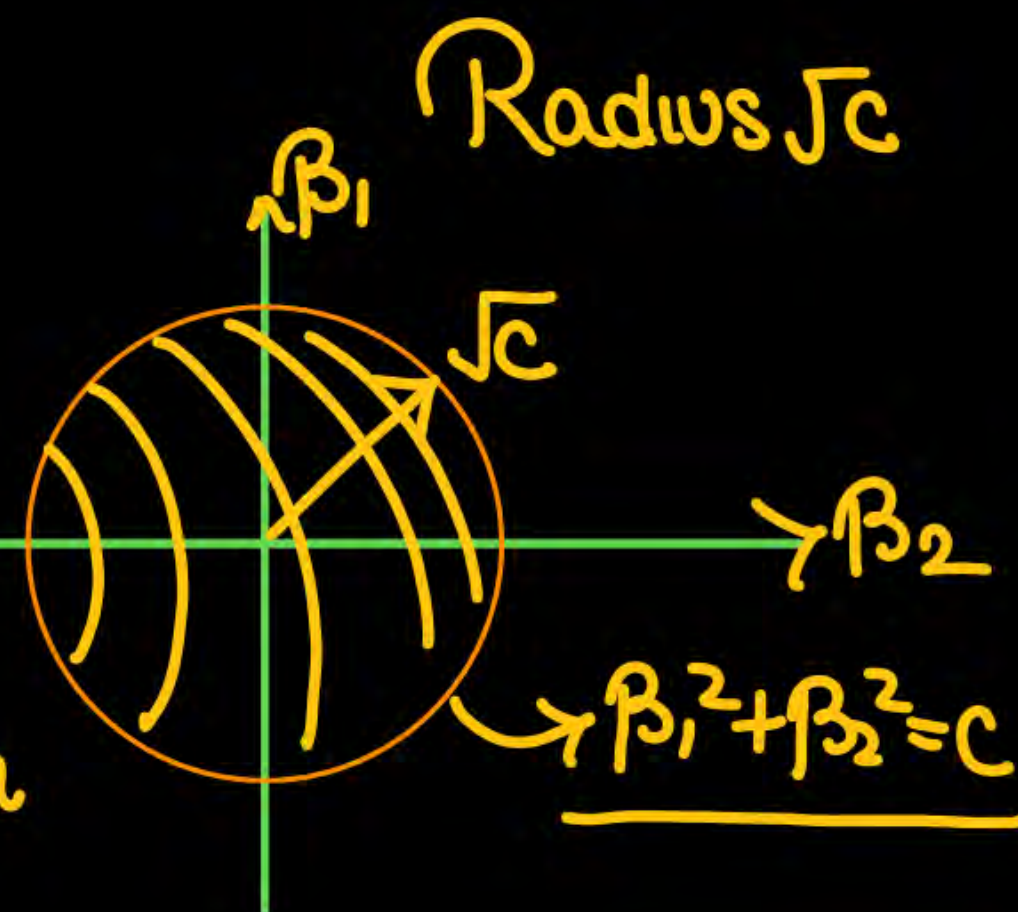
Constraint representation of Ridge Regression

Let data is of 2 dimension

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

St. $\beta_1^2 + \beta_2^2 \leq C$ \Rightarrow Plot

C: Constant
 β_1, β_2 unknown





Ridge Regression



What is Lasso Regularisation

Basic (not in Syllabus).

→ loss function = $\sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^D |\beta_i|$

The penalty term has absolute value of β .

we minimize the loss function

$f(x) = x^2$ Anjun $\rightarrow x = \sqrt{.001} = .031$
 $f(x) = x$ Karan. $\rightarrow x = .001$

- I want $\min f(x)$, $f(x)$ ki value zero nahi ho sakti
 $\underline{f(x) = .001}$ Karo.

When $x = .031$ then $x^2 = .001$
 When $x = .001$ then $x = .001$

$RR \Rightarrow \min RSS + \lambda \sum_{i=1}^D \beta_i^2$ (RR Keha hai $\beta_i = .031$ then $\beta_i^2 \Rightarrow .001$)

$lasso \Rightarrow \min RSS + \lambda \sum |\beta_i|$

$\beta_i = .001$

- lasso make β 's closer to 0 than RR.

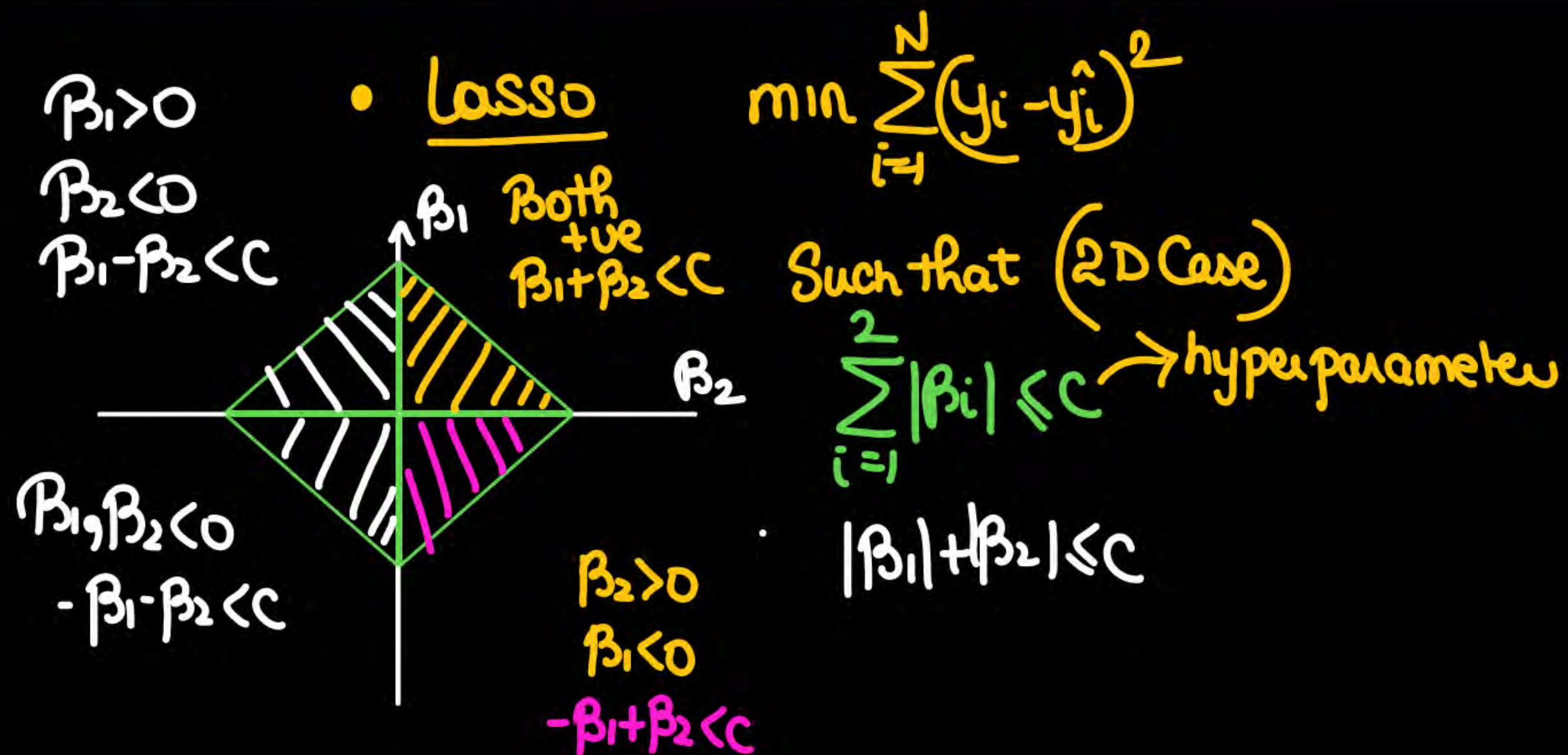
Lasso
Give more
Sparse
model.

→ isiliye lasso model mein β zyada β 's = 0 hote
hai,
Zyada dimension gayab from model.

Sparse model



Constraint representation of Lasso Regularisation

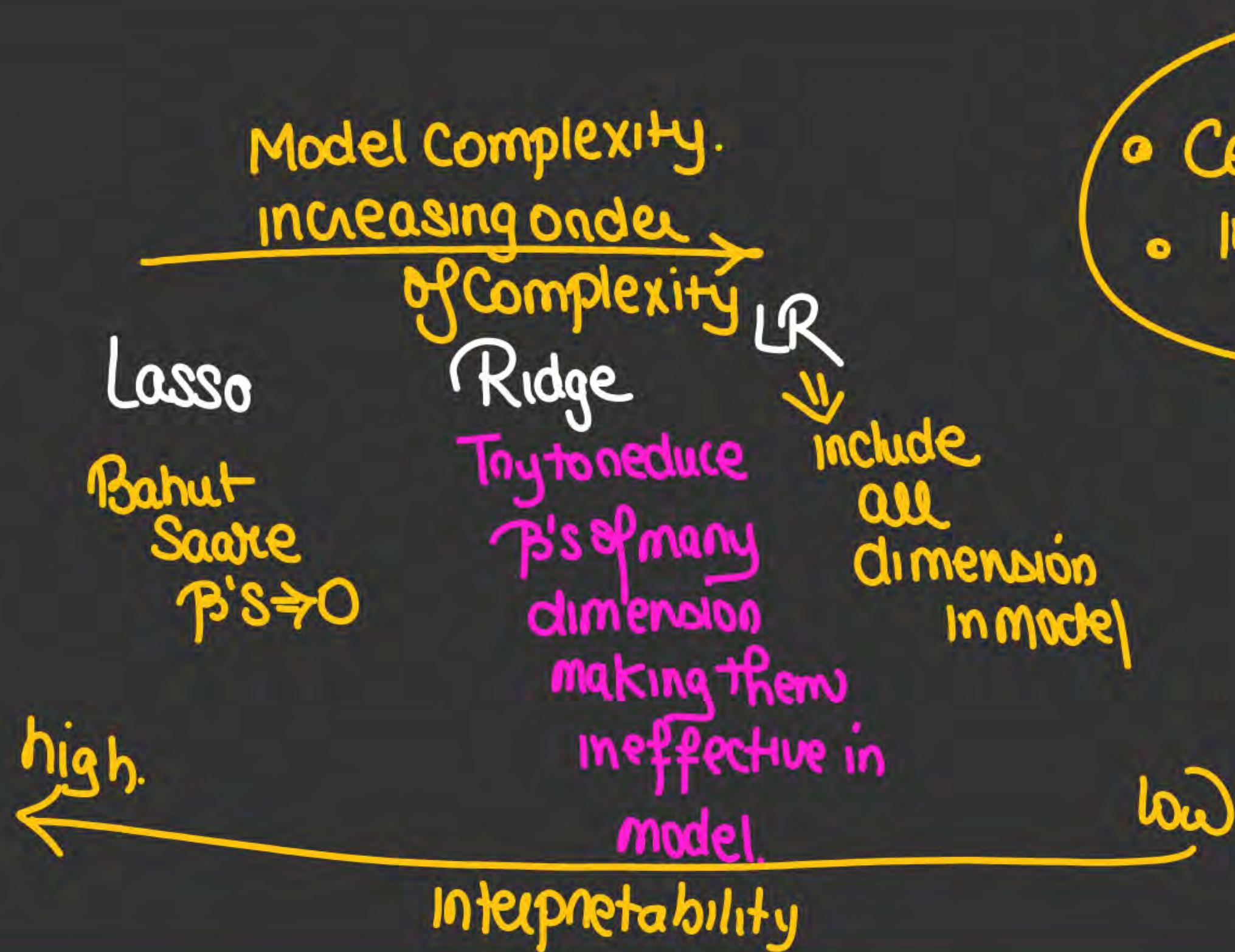


we have a relation
to find β 's directly.

• LR \rightarrow Closed form Sol. $\beta = (X^T X)^{-1} X^T Y$

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

Lasso \Rightarrow no closed form Sol.



- Complexity high
- Interpretability low

Penalty term $|\beta|$

Lasso is called L1 Regularisation

Ridge is called L2 Regularisation

Penalty term β^2

dpp
* H.W
* 80Q
* drive



Ridge Regression

Lasso Vs Ridge Regression

Read



Parameter	Ridge Regression	Lasso Regression
Regularization Type	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.
Primary Objective	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.
Feature Selection	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.
Coefficient Shrinkage	Coefficients are shrunk towards zero but not exactly to zero.	Coefficients can be shrunk to exactly zero, effectively eliminating some variables.
Suitability	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.
Bias and Variance	Introduces bias but reduces variance.	Introduces bias but reduces variance, potentially more than Ridge due to feature elimination.
Interpretability	Less interpretable in the presence of many features as none are eliminated.	More interpretable due to feature elimination, focusing on significant predictors only.
Sensitivity to λ	Gradual change in coefficients as the penalty parameter λ changes.	Sharp thresholding effect where coefficients can abruptly become zero as λ changes.
Model Complexity	Generally results in a more complex model compared to Lasso.	This leads to a simpler model, especially when irrelevant features are abundant.



- **Linear Classification**

- **Linear classification**

Classification vs Regression...



Linear Classification



Linear Regression of an Indicator Matrix

Let's consider a 2-class case

What is an Indicator Matrix



Linear Classification



Linear Regression of an Indicator Matrix

Let's understand
using figures





Linear Regression of an Indicator Matrix

So, now the analysis is as follows :



Linear Regression of an Indicator Matrix

Lets extend the case for K classes



Linear Classification



Linear Regression of an Indicator Matrix

How to find the variables for the linear regression



Linear Classification



**So linear regression can
be used for classification
also**



Linear Classification



Here we will have the error of $1/3$, hence the linear regression fails to classify even the seperable points.



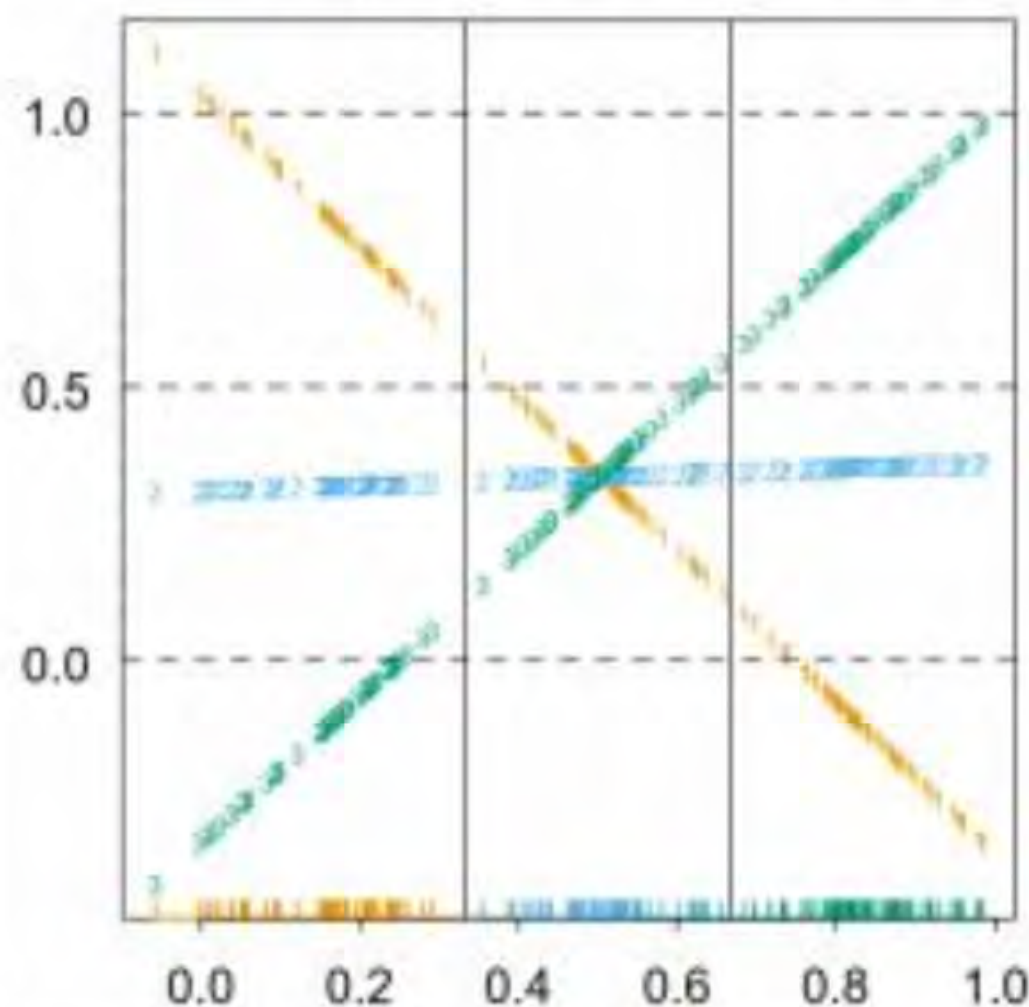


Linear Classification

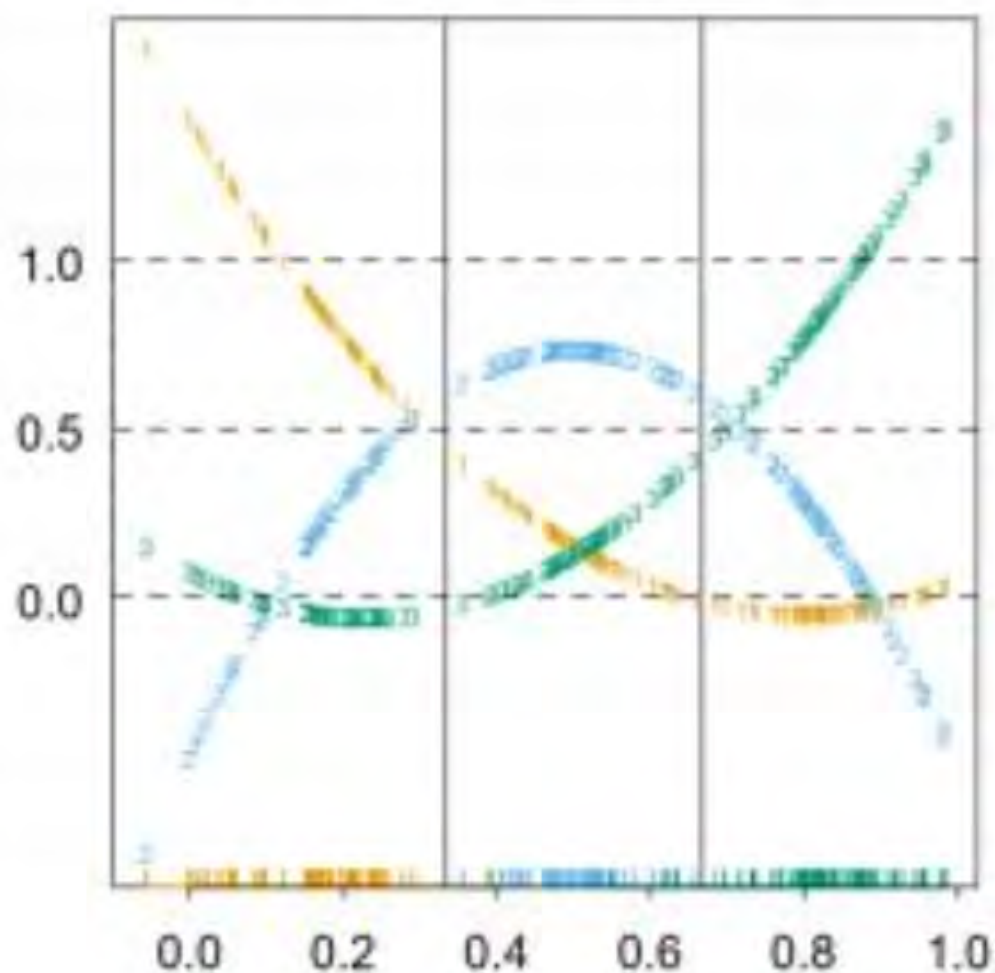


Linear Regression of an Indicator Matrix

Degree = 1; Error = 0.33



Degree = 2; Error = 0.04



The three classes are perfectly separated by linear decision boundaries, yet linear regression misses the middle class completely.

But we can classify if we use the quadratic curves.

A loose but general rule is that if $K \geq 3$ classes are lined up, polynomial terms up to degree $K - 1$ might be needed to resolve them.



Linear Regression of an Indicator Matrix

In general p -dimensional input space, one would need general polynomial terms and cross-products of total degree $K - 1$, $O(p^{K-1})$ terms in all, to resolve such worst-case scenarios.



Linear Classification



Linear Regression of an Indicator Matrix

Lets consider a 2 class problem... We can have a single classifier for a 2 class problem...



Linear Classification



Linear Regression of an Indicator Matrix

The loss function for a
2 class case...



Linear Classification



Linear Regression of an Indicator Matrix

But this loss function
has 2 problems 1.
outlier and 2. value of
predicted \hat{Y}



- **Linear Classification**

- **Linear Classification**

Problem of outliers



- **Linear Classification**

- **Linear Classification**

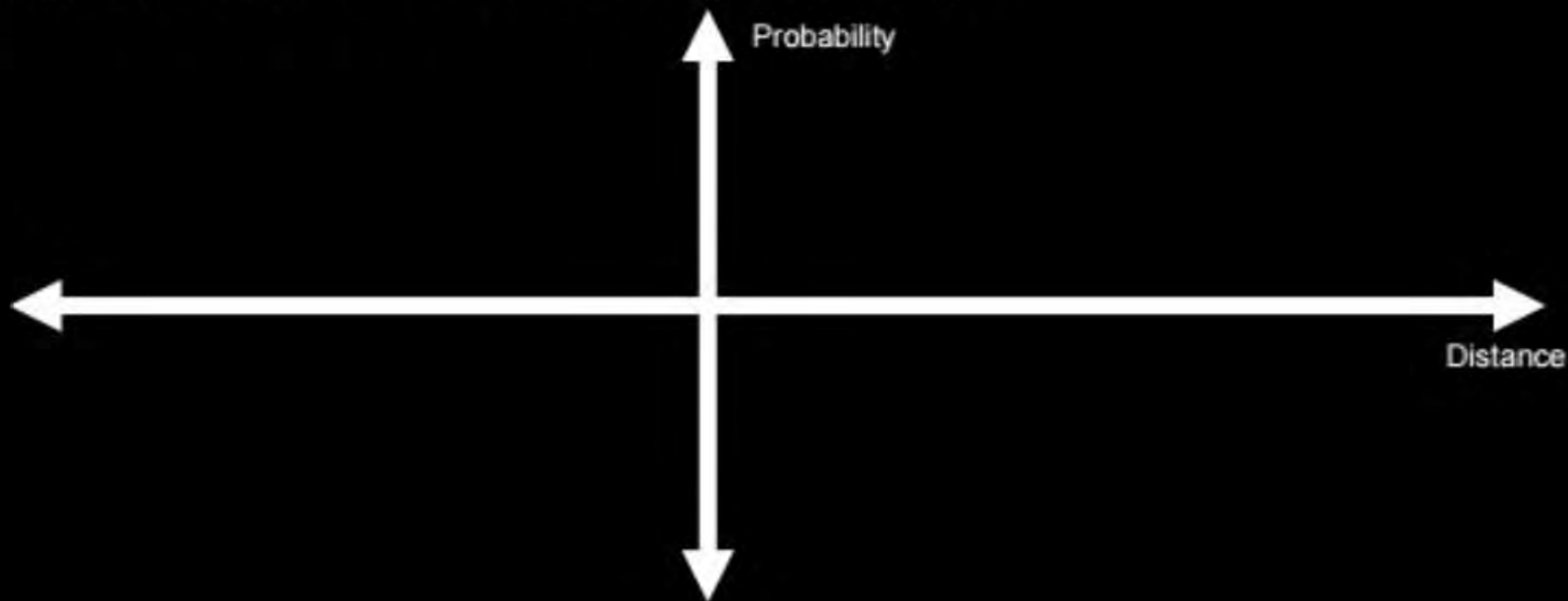
To solve the problem of outlier we will not use the distance in the analysis rather we will use the probability.



- **Linear Classification**

- **Linear Classification**

To solve the problem of outlier we will not use the distance in the analysis rather we will use the probability.



THANK - YOU