# Data Science and Artificial Intelligence

## Machine Learning

**Classification**

**Lecture No. 1**

GATE WALLAH

By- SIDDHARTH SABHARWAL SIR

# Recap of Previous Lecture

**Topic** — Ridge Reg.

**Topic** — Cross Validation → K-Fold Cross Validation

**Topic** — effect of $\lambda$

**Topic**

**Topic**

# Topics to be Covered

**Topic** — Lasso

**Topic** — Lasso vs RR

**Topic** — Classification

**Topic** — doubts

**Topic**

Parade

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

Plan and work for it

**Ridge Regression Final expression**

$$\rightarrow \quad \min \; \frac{1}{2} \sum (y_i - \hat{y}_i)^2$$

$$\text{Const} \quad \sum_{i=1}^{D} \beta_i^2 < C$$

Ridge Regression Final expression

# Linear Reg

$$\mathcal{L} = \min_{i=1}^{N} \sum (y_i - \hat{y}_i)^2$$

$$\frac{\partial L}{\partial \beta} = \begin{bmatrix} \partial y / \partial \beta_0 \\ \partial y / \partial \beta_1 \\ \vdots \\ \partial y / \partial \beta_D \end{bmatrix} \Rightarrow -2 \left[ X^T Y - (X^T X) \beta_{old} \right] \checkmark$$

# Ridge Reg

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

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

$$\left( \beta^{new} = \beta^{old} - \eta \cdot \frac{\partial L}{\partial \beta} \right)$$

## 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 – lets practise**

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

$$\lambda/2 \sum_{i=1}^{P} \beta_i^2$$

Siddharth Sir AI/ML.

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.

## Ridge Regression – lets practise

Ridge Regression is particularly useful when:

*RR is used*

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 – lets practise

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

→ Strongest Reg.

B→0

A) λ = 0
B) λ = 1
C) λ = 10
D) λ = ∞

## 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.

*Handwritten annotations:*

Regularisation best $\lambda$

Overfitting X
generalise model ✓

* This is not a problem bcoz our model is generalising better

→ The model give more error on training data

## Ridge Regression – lets practise

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

→ Testing error.

A) Increases bias, decreases variance.

B) Decreases bias, increases variance.

C) Increases both bias and variance.

D) Decreases both bias and variance.

as λ inc → model generalise better

model complexity reduce, Training error

→ training error inc / testing err dec

Bias inc / Var dec

## 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

$$\text{Penalty term} = \lambda \sum \beta_i^2 / 2$$
$$\Rightarrow 3 \times 50 / 2$$
$$\Rightarrow 150 / 2 \Rightarrow \boxed{75}$$

## 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    *done*

b)135

c)123

d)578

## Interpretability

So in L·R we get

$$y = \beta_0 + \beta_1 x^1 + \beta_2 x^2 \cdots -$$

from this equation we can imagine our model i.e the pattern of data.

⟹ dR ⟹ high Interpretability

## Ridge Regre.

generalise better, remove effect of few useless dimensions

Thus gives a better/easier model that Represent pattern of data ⟹ More Interpretable model.

## 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?

we are comparing RR with LR

* $LR \Rightarrow RSS \Rightarrow 60$

* $RR \Rightarrow RSS \Rightarrow 50$

So $R^2 \Rightarrow 1 - \dfrac{RSS \text{ of } RR}{RSS \text{ of } LR} \Rightarrow 1 - \dfrac{50}{60} \Rightarrow \left(\dfrac{1}{6}\right)$

$\left(R^2 \Rightarrow \dfrac{1}{6}\right) \Rightarrow 16.66\%$

## What is Lasso Regularisation

$$\text{(L1 Reg)}$$

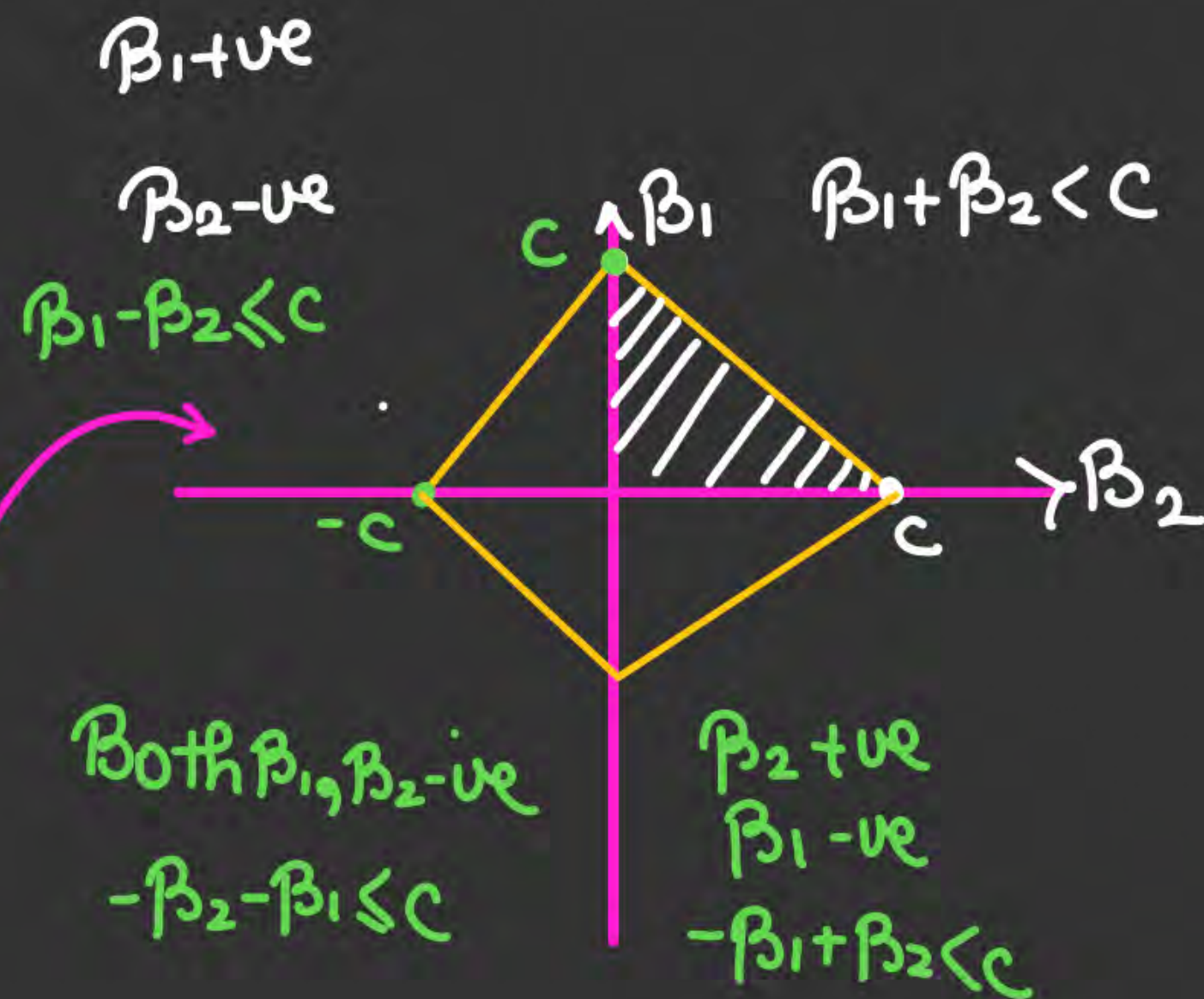$$\text{loss fxn} \Rightarrow \left( \frac{1}{2} \sum (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^{D} |\beta_i| \right)$$

$$\rightarrow OR \Rightarrow \min \frac{1}{2} \sum (y_i - \hat{y}_i)^2$$

$$\text{Such that } \sum_{i=1}^{D} |\beta_i| < C$$

$$\begin{cases} |x| \to x & \text{when } x\ +ve \\ \quad\ -x & \text{when } x\ -ve \end{cases}$$

Taking a 2D Case

$$|\beta_1| + |\beta_2| \leqslant C$$

So Constraint
Kite/Square.

$\beta_1\ +ve$

$\beta_2\ -ve$

$\beta_1 - \beta_2 \leqslant C$

$\beta_1 + \beta_2 < C$

Both $\beta_1, \beta_2\ -ve$

$-\beta_2 - \beta_1 \leqslant C$

$\beta_2\ +ve$
$\beta_1\ -ve$
$-\beta_1 + \beta_2 < C$

1. The algorithm of RR try to reduce $\beta$ values but lass Reg. has the algorithm that try to make $\beta$'s $\to 0$.

2. lasso is more sensitive to `$\lambda$' than Ridge Regression

*Interpret:*
- *linear R*
- *^*
- *Ridge R*
- *^*
- *Lasso R*

*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 $\lambda$ | Gradual change in coefficients as the penalty parameter $\lambda$ changes. | Sharp thresholding effect where coefficients can abruptly become zero as $\lambda$ 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. |

*When # of feature is v.high*

Interpretability   LR $<$ RR $<$ lasso

Complexity          LR $>$ RR $>$ lasso

generalise          LR $<$ RR $<$ lasso

**Q7** Using the data $X=[-3,5,4]$ and $Y=[-10,20,20]$, assuming a ridge penalty $\lambda = 50$, what ratio versus the Maximum Likelihood Estimate (MLE) estimate $w_{mle}$ do you think the ridge regression L2 estimate estimate $w_{ridge}$ estimate will be?

(A) 2

(B) 1

(C) 0.6

(D) 0.5

W is the $\beta$

- $\beta_{MLE}$ = the $\beta$ of LR

- $\beta_{RRv}$

we will have $(w_0, w_1)$

$\Rightarrow$ Ratio of $\dfrac{w_{MLE}}{w_{RR}} = \dfrac{w_{LR}}{w_{RR}} = \dfrac{w_{1LR}}{w_{1RR}}$

$$RR \Rightarrow X \Rightarrow \begin{bmatrix} -3 \\ 5 \\ 4 \end{bmatrix} \xleftarrow{\quad} \qquad Y = \begin{bmatrix} -10 \\ 20 \\ 20 \end{bmatrix} \Rightarrow \text{Centering}$$

$$\Rightarrow dR \Rightarrow \begin{bmatrix} 1 & -3 \\ 1 & 5 \\ 1 & 4 \end{bmatrix}$$

$$X = \begin{bmatrix} \\ \\ \end{bmatrix} \qquad Y = \begin{bmatrix} \\ \\ \end{bmatrix}$$

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

$$\Rightarrow B_1 = (X^T X + \lambda I)^{-1} (X^T Y)$$

$$\Rightarrow B_0 = \bar{y} - B_1 \bar{x}$$

Concept $\Rightarrow$

$\|\beta\|^2$

$\|Y\|^2$

$$A = \begin{bmatrix} a \\ b \\ c \\ a \end{bmatrix}$$

$\|A\|^2 \Rightarrow (\text{norm})^2$

$\Rightarrow$ Sum of Square of Elements

$\|A\|^2 \Rightarrow a^2 + b^2 + c^2 + d^2$

**Q1** Consider the linear regression model $Y = X\beta + \varepsilon$ with $\varepsilon \sim N(0n, \sigma_\varepsilon^2\, Inn)$. This model (without intercept) is fitted to data using the ridge regression estimator $\hat{\beta}(\lambda) = \arg\min_\beta \|Y - X\beta\|_2^2 + \lambda\|\beta\|_2^2$ with $\lambda > 0$.

**The data are:**

$$X^T = (-1\ 1\ 1\ -1) \text{ and } Y^T = (-1.5\ 2.9\ -3.5\ 0.7\ )$$

What is the maximum likelihood/ordinary least squares estimator of the regression parameter for $\lambda = 0$?

(A) $[-0.3, 0.05]$         (B) $[-0.5, 0.1]$

(C) $[0.1, -0.2]$         (D) $[0.05, -0.3]$

**Q2** Suppose you are training a Ridge Regression model for a particular task and notice the following training error and validation RSS

Train: 5?

You have a dataset with 30 observations. After applying linear regression, you find that the residual standard error (RSE) is 5. If the coefficient of determination (R^2) is 0.8, what is the root mean square error (RMSE) for this model?

(A) 2

(B) 3

(C) 4

(D) 5

$$RSE = RMSE$$

$$5 = RMSE$$

- Linear classification

**Classification vs Regression...**

- The X value can be anything, but the Y values are Categorical

↳ we try to create a model which predict Y based on X

↳ Value of Y can be any real number.

**Linear Regression of an Indicator Matrix**

One hot Coding

2 classes

Yes ← → No

Let's consider a 2-class case

What is an Indicator Matrix

| Fever | H·R | Sugar | $Y_1$ | $Y_2$ |
|-------|-----|-------|-------|-------|
| – | – | – | 1 | 0 |
| – | – | – | 0 | 1 |
| – | – | – | 1 | 0 |
| – | – | – | 0 | 1 |
| | | | 0 | 1 |

## Linear Regression of an Indicator Matrix

- So we have 2 classes and we Create 2 Y values $Y_1, Y_2$

- For each data point only one Y value will be `1` rest will be `0`.

Let's consider a 2-class case

What is an Indicator Matrix

So Now we use regression for classification ⟫

→ So Now taking $Y_1$ as Y we do L.R and we will get a line that will give '1' for points of 'Yes' Class '0' for points of 'No Class

→ Similarly taking $Y_2$ as Y we do L.R and we will get a line that give '1' for 'No' points and '0' for Yes points.
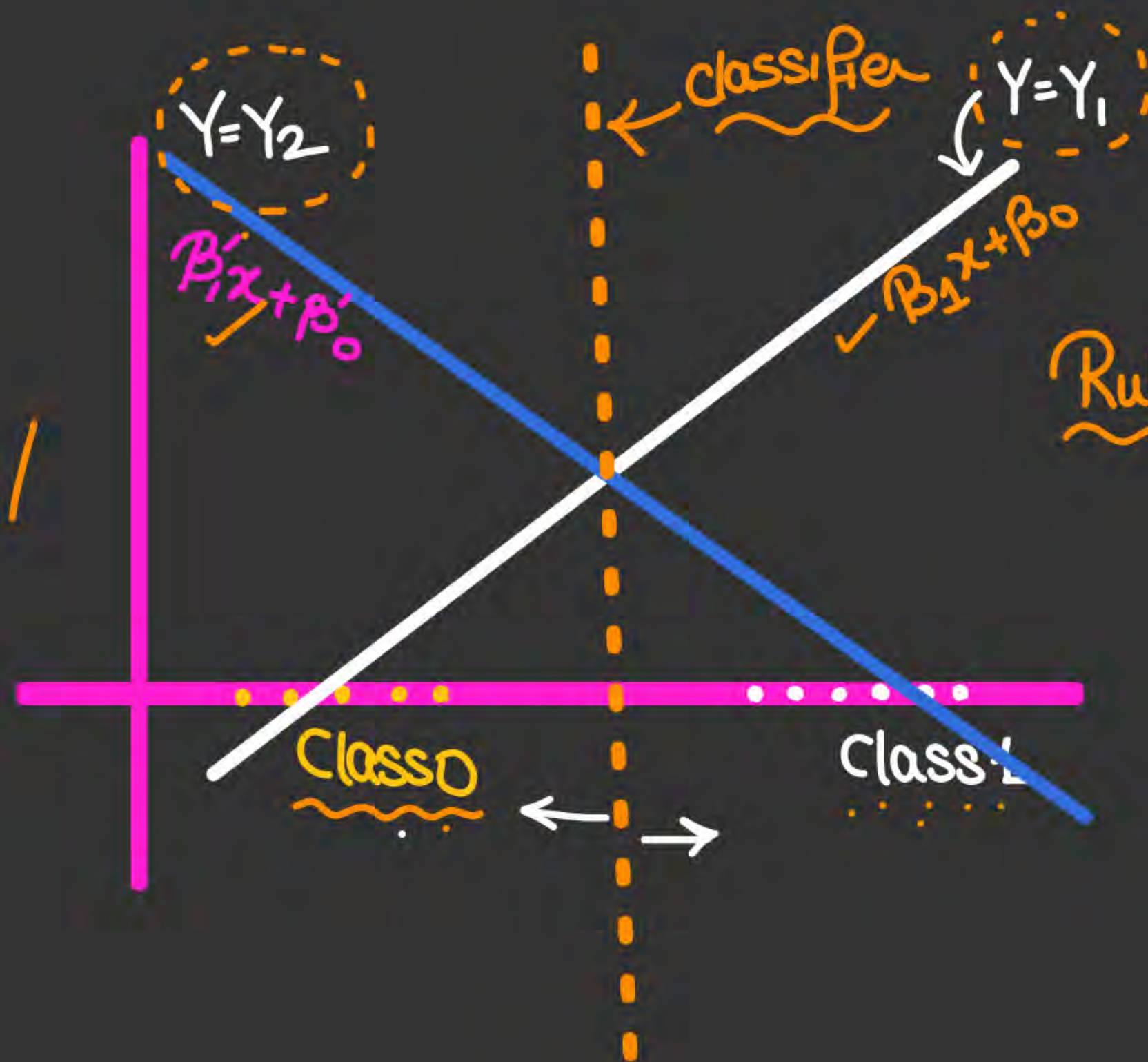
## Linear Regression of an Indicator Matrix

doubt: We can work with a Single
line such that if value of
Y from line at any point is close to 1 ⇒ Yes
Y from " " " " " " " 0 ⇒ No

Let's understand using figures



Purely Reg. Activity not a Seperator.

No

Yes

- But actually in classification
  our motive is not to find y value,
  from line
  Rather we need a line that
  separate the points belonging
  to diff classes

good classifier

good for
Regne.

Class 0

Class 1

## Linear Regression of an Indicator Matrix

**So, now the analysis is as follows :**

So if we have data with 2 Classes

↳ → One hot Coding ⇒ Indicator matrix

→ two times L.R ⇒ 2 lines

→ Condition $\beta_1' x + \beta_0' > \beta_1 x + \beta_0 \Rightarrow$ class0

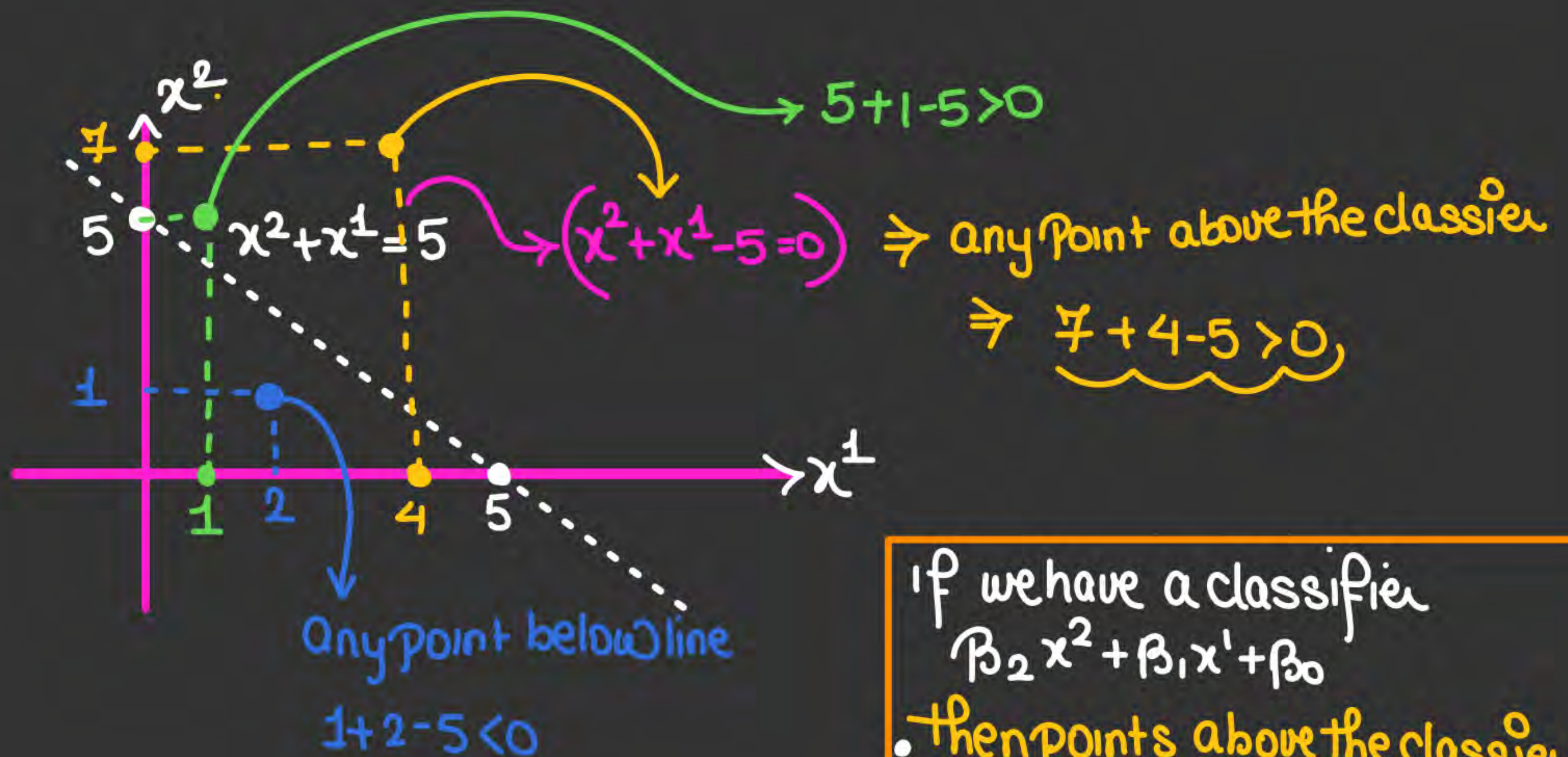$\beta_1 x + \beta_0 > \beta_1' x + \beta_0' \Rightarrow$ class1

**Linear Regression of an Indicator Matrix**

**Lets extend the case for K classes**

from these 2 conditions we actually get a single classifier
line b/w 2 clases.

Can we find this classifier directly ?? ⇒ Yes by linear
Classifier

$x^2$

$7$

$5$

$x^2 + x^1 = 5$

$5 + 1 - 5 > 0$

$(x^2 + x^1 - 5 = 0)$ ⇒ any point above the classifier

⇒ $7 + 4 - 5 > 0$

$1$

$1$  $2$  $4$  $5$  $x^1$

Any point below line

$1 + 2 - 5 < 0$

'If we have a classifier
$\beta_2 x^2 + \beta_1 x^1 + \beta_0$

• then points above the classifier
$\beta_2 x^2 + \beta_1 x^1 + \beta_0 > 0$

• Similarly points below classifier
$\beta_2 x^2 + \beta_1 x^1 + \beta_0 < 0$

## Linear Regression of an Indicator Matrix

Lets extend the case for K classes

11 am

## Linear Regression of an Indicator Matrix

How to find the variables for the linear regression

So linear regression can be used for classification also

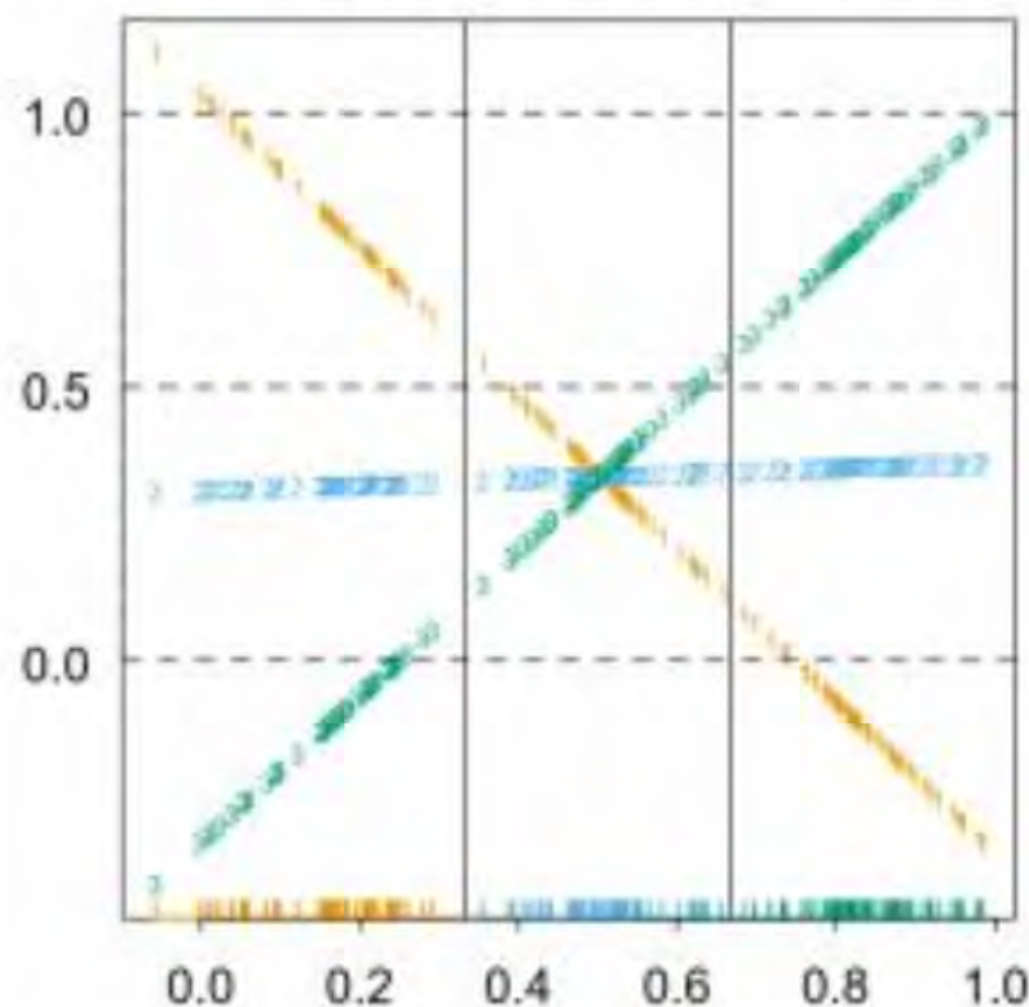Here we will have the error of 1/3, hence the linear regression fails to classify even the seperable points.
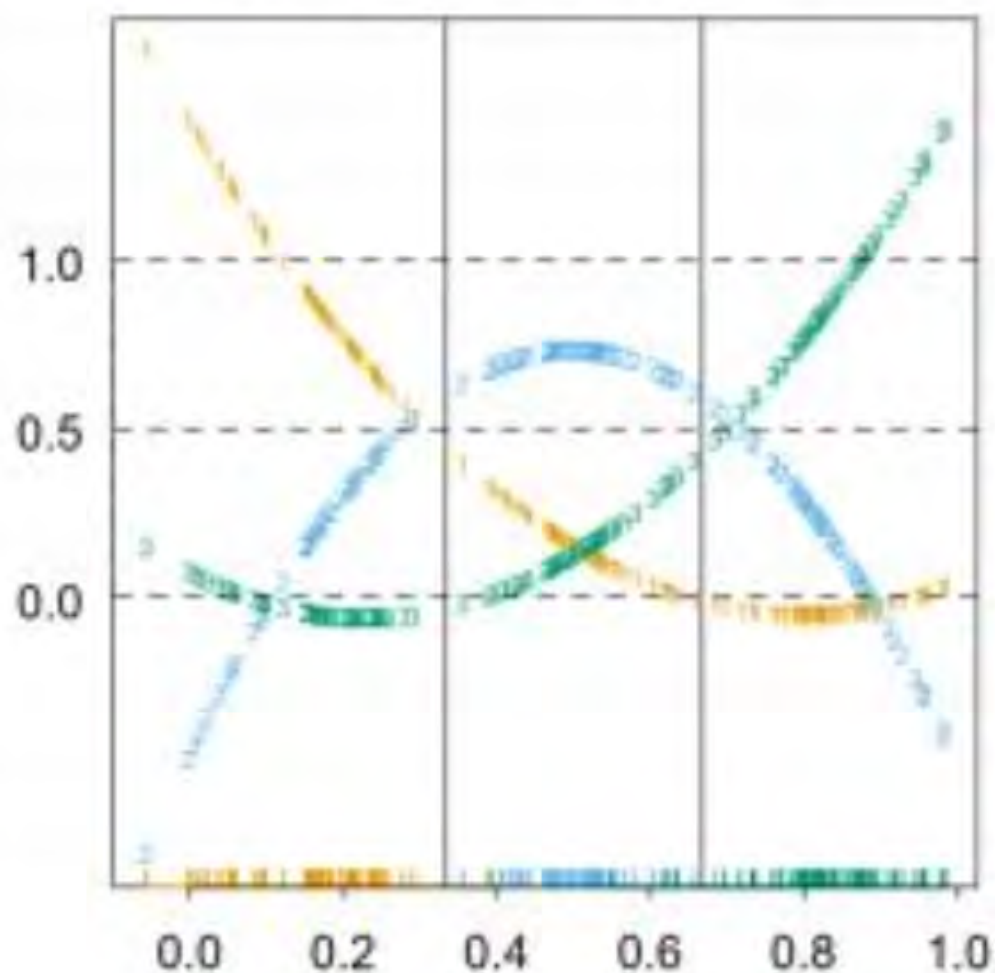
## Linear Regression of an Indicator Matrix



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 ≥ 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 Regression of an Indicator Matrix

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

## Linear Regression of an Indicator Matrix

The loss function for a
2 class case...

## Linear Regression of an Indicator Matrix

But this loss function has 2 problems 1. outlier and 2. value of predicted Y

- **Linear Classification**

**Problem of outliers**

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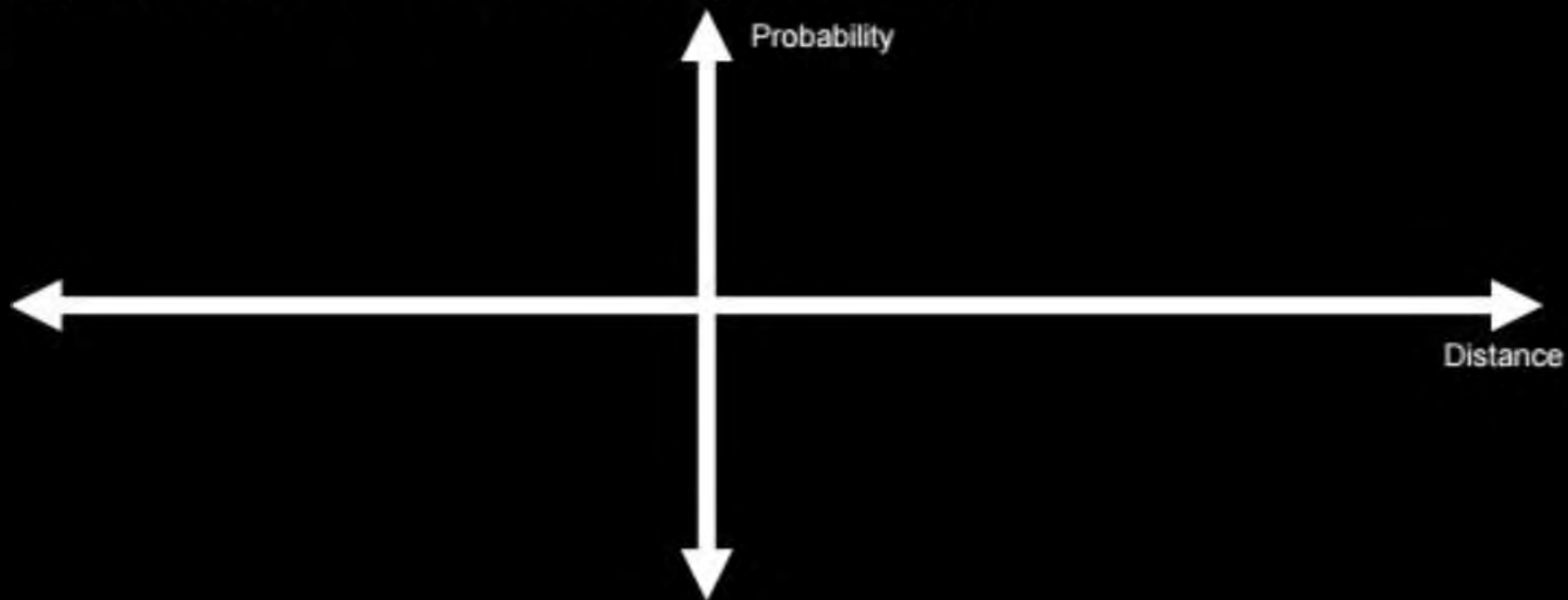
To solve the problem of outlier we will not use the distance in the analysis rather we will use the probability.

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

Probability

Distance

# 2 mins Summary

Topic

Topic

Topic

Topic

Topic

THANK - YOU