

(h=4)

## Classification

The linear regression we discussed in last chapter assumes response variable ( $Y$ ) is quantitative.

but what if that variable is qualitative.  
Predicting a qualitative response for an observation is referred to as classification.

### An Overview of Classification:

Classification problems occur often, even more than regression problems:

1. Credit Card Fraud Detection: An online banking system uses it to detect fraud.
2. Classifying disease causing cells in a DNA Sequence: On the basis of DNA sequence, a biologist would like to figure out DNA mutations.
3. Classifying an Email is Spam or Not: On the basis of text vectors a software developer would like to classify if the mail is spam or not.

Given a feature vector  $X$  & a qualitative response  $Y$  taking values in set  $G$ , the task is to build a function  $G(X)$  that takes as input the feature vector  $X$  and predicts its value for  $Y_i$  i.e.  $G(X) \in G$

### Why not linear regression?

We have stated that linear regression is not appropriate in the case of qualitative response. Why not?

Let us suppose we can, we'll look at it with an example:

Let us say, we are trying to predict of 3 possible diagnoses:-

$$Y = \begin{cases} 1; & \text{stroke} \\ 2; & \text{drug overdose} \\ 3; & \text{epileptic seizure} \end{cases}$$

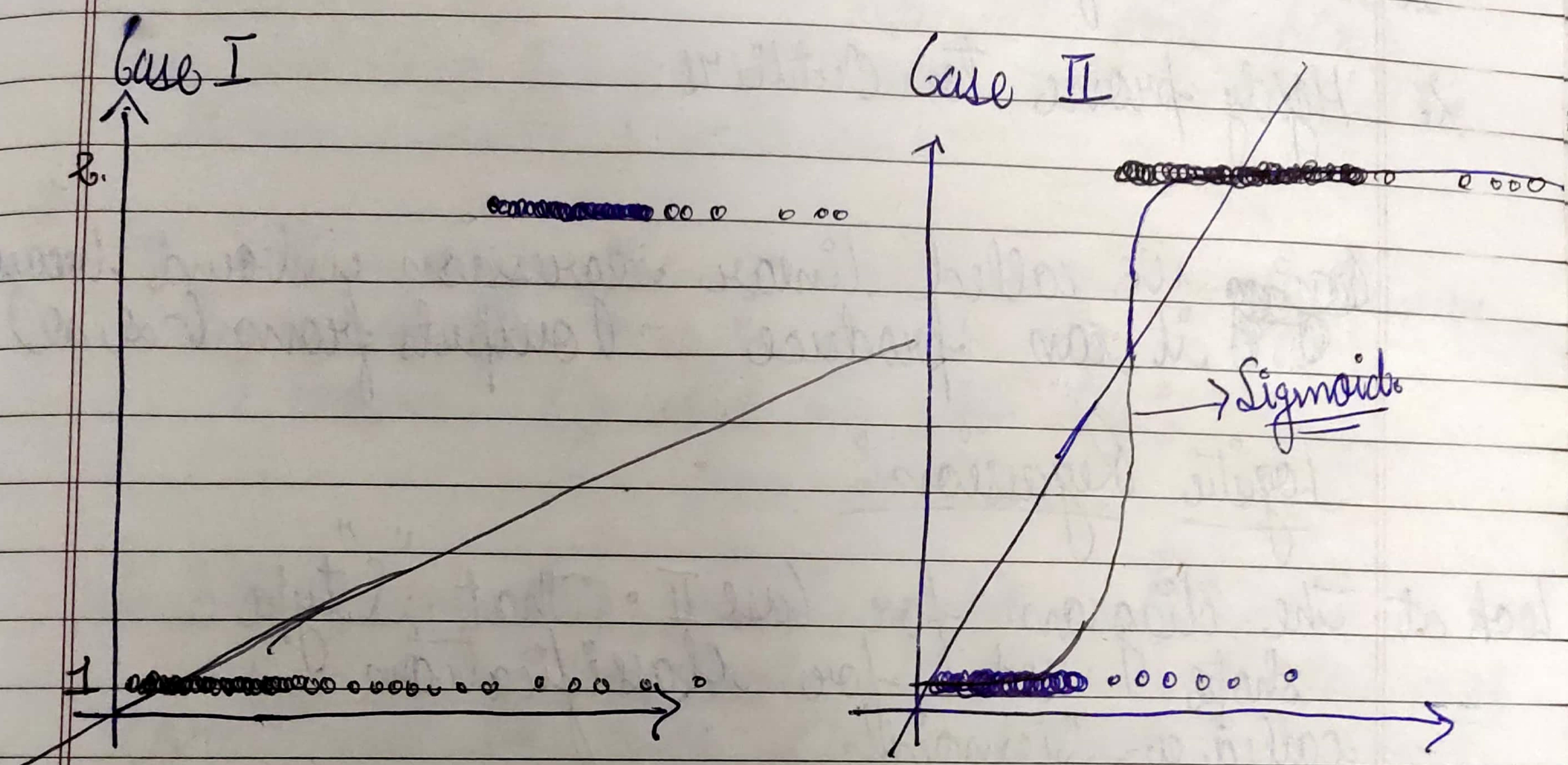
Applying L.R. This impression of "Y" implies that the difference between a stroke & drug overdose is equal to epileptic seizure & drug overdose.

But in real life that's not the case.

For one instance we could encode this as:

$$Y = \begin{cases} 1; & \text{epileptic seizure} \\ 2; & \text{stroke, drug overdose} \\ 3; & \text{drug overdose, stroke} \end{cases}$$

Using linear regression & fitting the model & in Case I & II would be totally different.



o Stroke

o Drug Overdose

The graphs above implies 2 things :-

1. If the encoding / priority changes the graph as well best fit line changes widely. This implies high variation.
2. There is high variance in line (best fit) so the  $r^2$  would be low.

2 major reasons why we can't use linear regression for classification:

1. Linear Regression is Unbound
2. Highly prone to outliers

We called linear regression unbound because it can produce outputs from  $(-\infty, \infty)$ .

Logistic Regression:

Look at the diagram for Case II. That "S" type shape used for classification is called a Sigmoid.

Let's write  $p(x) = P(Y=1|x)$

Logistic Regression uses the form:

$$p(x) = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$$

where,  $C \approx 3.718$

& whatever values  $B_0, B_1$  or  $x$  takes,  $p(x)$  will have values between 0 & 1.

A piece of rearrangement gives:

$$\log \left( \frac{p(x)}{1-p(x)} \right) = B_0 + B_1 x$$

This is called log odds or a logit of  $P(x)$

Because of the output of log odds, it is called Logistic Regression

where,  $\frac{P(t)}{P(1-t)}$  is called odds of  $t$ .

As shown in Case II graph, logistic regression ensures that our estimate for  $P(x)$  lies between 0 & 1 & fits our data much better than linear regression.

\* It solves both the problems given by linear regression.

where, Sigmoid is given by:  $\frac{1}{1+e^{-x}}$

& Because the relationship between  $x$  &  $P(x)$  is not linear,  $\beta_1$  does not correspond to one.  $\Delta$  change in  $P(x)$  with one unit increase in  $x$  as there was the case in linear regression.

For a binary classification:

$$\begin{array}{ll} P(x) \geq 0.5 & ; \text{ Class} = 1 \\ P(x) < 0.5 & ; \text{ Class} = 0 \end{array}$$

### Estimation of Parameters:

The coefficients  $\beta_0$  &  $\beta_1$  are unknown and has to be estimated based on the available training data.

We use MAXIMUM LIKELIHOOD ESTIMATION to estimate the parameters of logistic regression.

## MAXIMUM LIKELIHOOD ESTIMATION:

MLE is a method of estimating the parameters of a probability distribution by maximizing the likelihood function so that under the assumed model the observed data is most probable, in our case the model is logistic & parameters are  $B_0$  &  $B_1$ .

QUESTION

HOW DOES A LIKELIHOOD FUNCTION USED / CALCULATED:

First, of all lets try to understand the difference between likelihood & probability.

~~Probability @ ~~discrete~~ probability functions~~

Probability of  $Y$  given Parameters is probability  
at a point "X" where distribution is known.

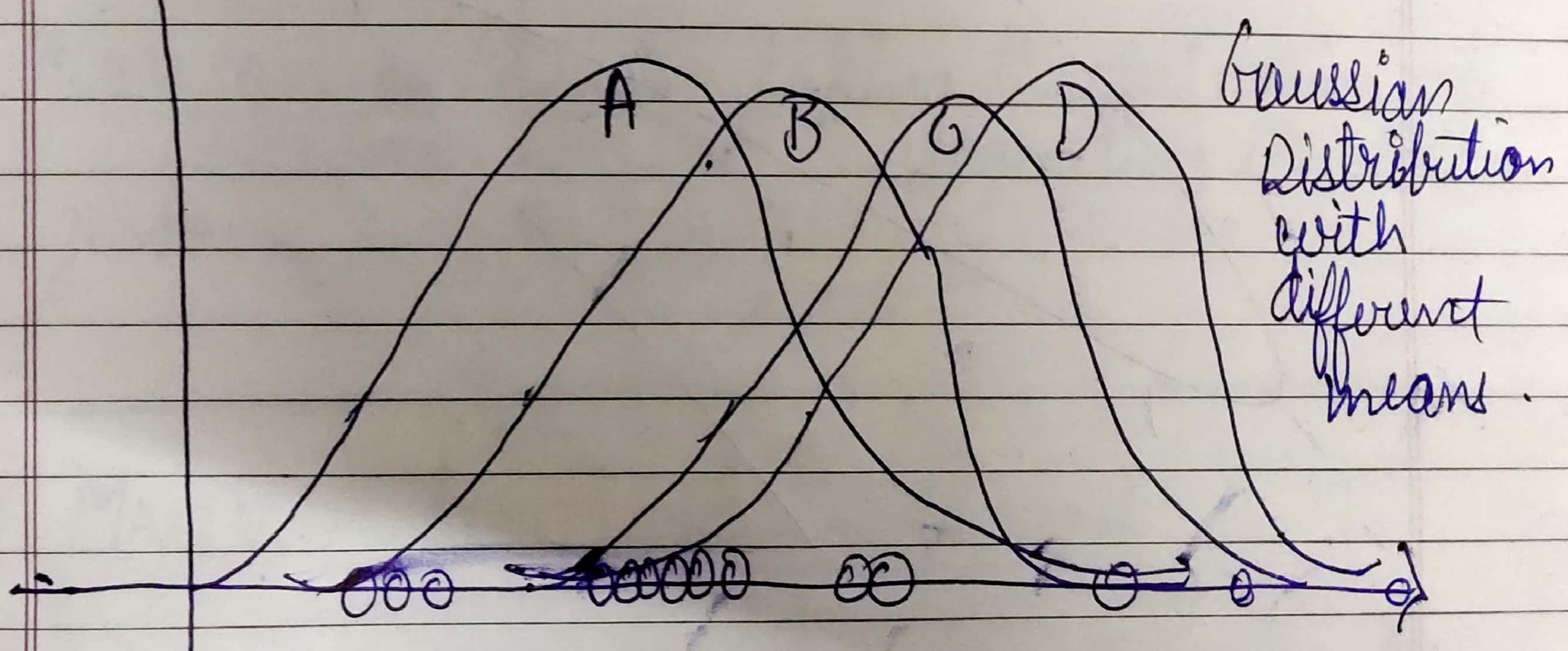
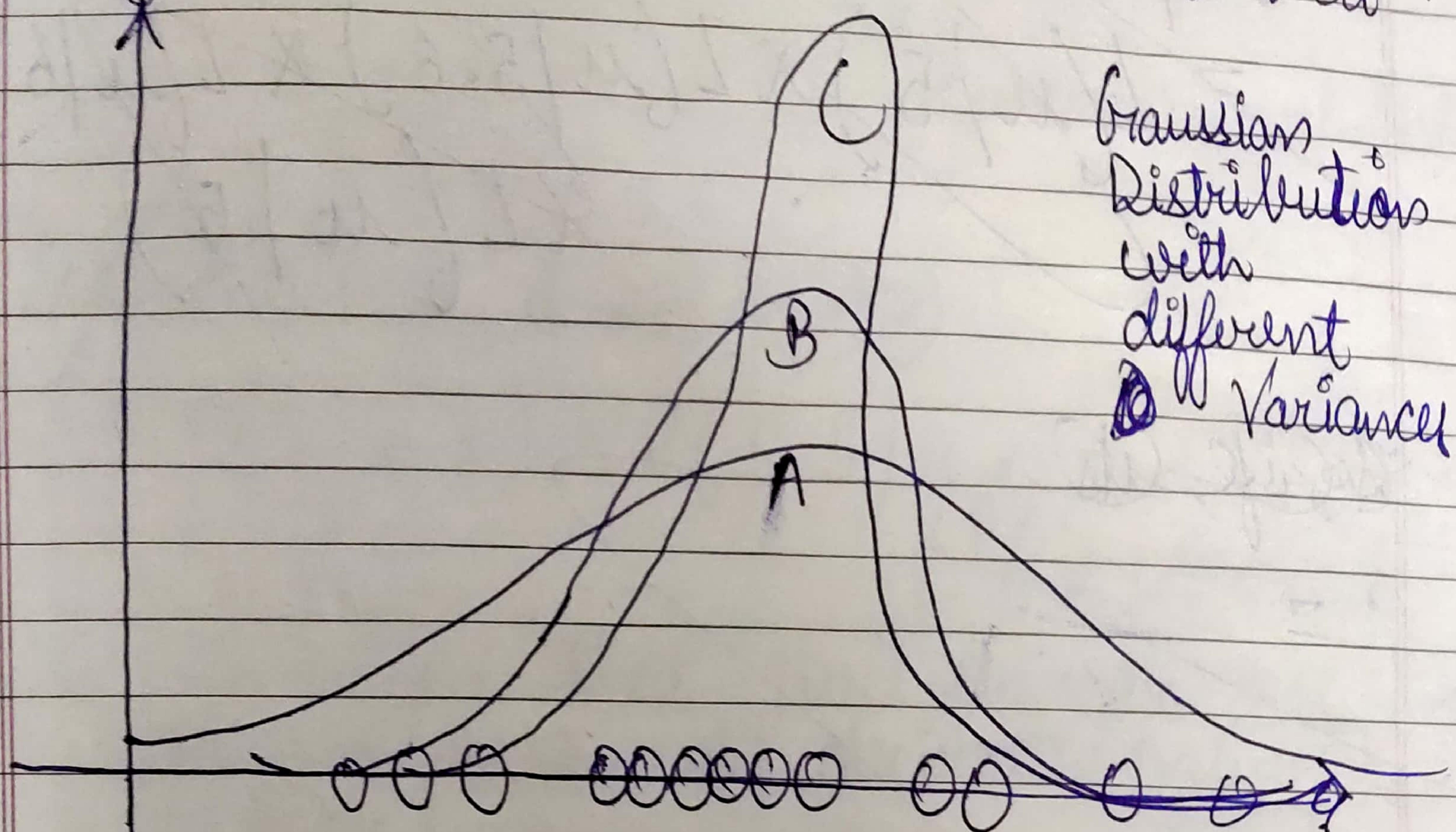
While

Likelihood is probability of data belonging to a certain distribution.

It sounds like a hickey, but you'll understand through an example.

Let us assume we are given a data & we are trying to find a best distribution that fits that specific data.

We decided that we want to fit gaussian Normal distribution to it.



So, Maximum Likelihood estimation lets you get the best fit distribution from infinite distributions with different means & variances.

Let's say 4 data points are given & we have been given a task to decide mean & variance for distribution using MLE.

$$P(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\left(\frac{(x-\mu)^2}{2\sigma^2}\right)}$$

(Gaussian Normal Distribution)

~~$y(\mu, \sigma | x=5, 5.5, 6, 5)$~~

$$= P(5|\mu, \sigma) \times P(5.5|\mu, \sigma) \times P(6|\mu, \sigma) \times P(5|\mu, \sigma)$$

The distribution that gives the max value for likelihood here will decide the values for  $\mu$  &  $\sigma$ .

\* We can do this by taking the derivative of likelihood function & equating it to zero.

$$X - X -$$

## MAXIMUM LIKELIHOOD ESTIMATION:

We use Maximum Likelihood Estimation to estimate the parameters.

$$L(\beta_0, \beta_1) = \prod_{y_i=1} P(x_i) \prod_{y_i=0} (1 - P(x_i))$$

## MARINER PREDICTIONS

Once the coefficients are estimated we put

those estimates in  $p(x)$

$$\hat{P}(x) = \frac{e^{\hat{B}_0 + \hat{B}_1 x}}{1 + e^{\hat{B}_0 + \hat{B}_1 x}}$$

if this value is  $> 0.5$  then class-1 is predicted

else class-2 will be predicted  
\*\* we can always manipulate this threshold.

## MULTI-VARIABLE LOGISTIC REGRESSION:

This changes nothing, the only thing that changes is the number of predictors.

$$\hat{P}(x) = \frac{e^{\hat{B}_0 + \hat{B}_1 x_1 + \hat{B}_2 x_2 + \hat{B}_3 x_3}}{1 + e^{\hat{B}_0 + \hat{B}_1 x_1 + \hat{B}_2 x_2 + \hat{B}_3 x_3}}$$

## MULTI-CLASS LOGISTIC REGRESSION:

$$P(X=k | x) = \frac{e^{\hat{B}_{0k} + \hat{B}_{1k} x_1 + \hat{B}_{2k} x_2 + \dots + \hat{B}_{pk} x_p}}{\sum_{l=1}^K e^{\hat{B}_{0l} + \hat{B}_{1l} x_1 + \dots + \hat{B}_{pl} x_p}}$$

Multiclass logistic regression is also called Multinomial Regression

## \* Linear Discriminant Analysis:

Here approach is to model the distribution of  $X$  in each of the classes separately & then use Bayes theorem to flip things around and obtain  $\text{Pr}(Y|X)$ .

### BAYES THEOREM FOR CLASSIFICATION:

Thomas Bayes was a famous mathematician whose name represents a big subfield of stats & probability.

Here we'll focus on BAYES THEOREM:

$$\text{Pr}(Y=k | X=x_0) = \frac{\text{Pr}(X=x_0 | Y=k) \cdot \text{Pr}(Y=k)}{\text{Pr}(X=x_0)}$$

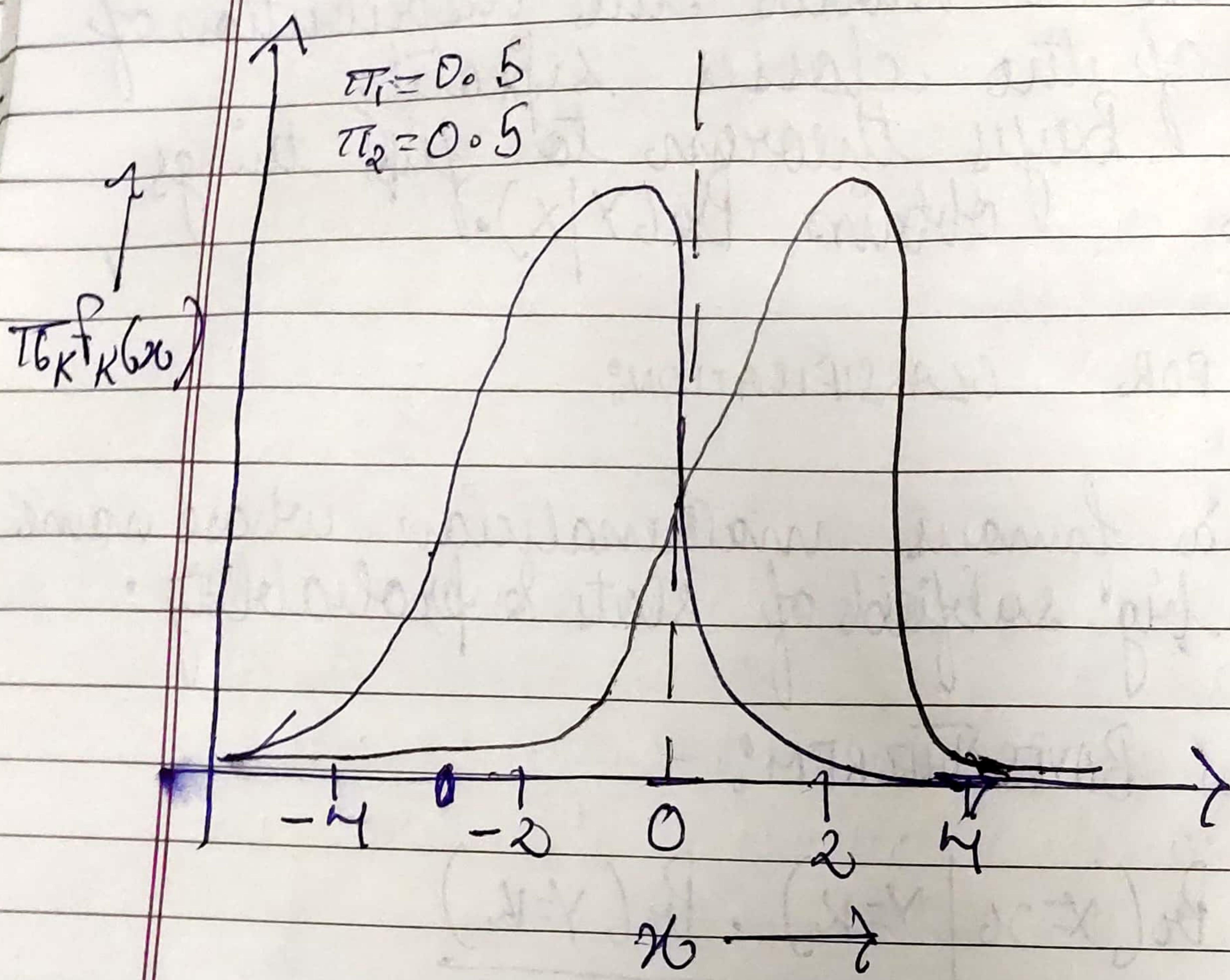
In discriminant analysis it is written quite differently:

$$\text{Pr}(Y=k | X=x_0) = \frac{\pi_k f_{k|x_0}}{\sum_{l=1}^K \pi_l f_{l|x_0}} \quad \text{where}$$

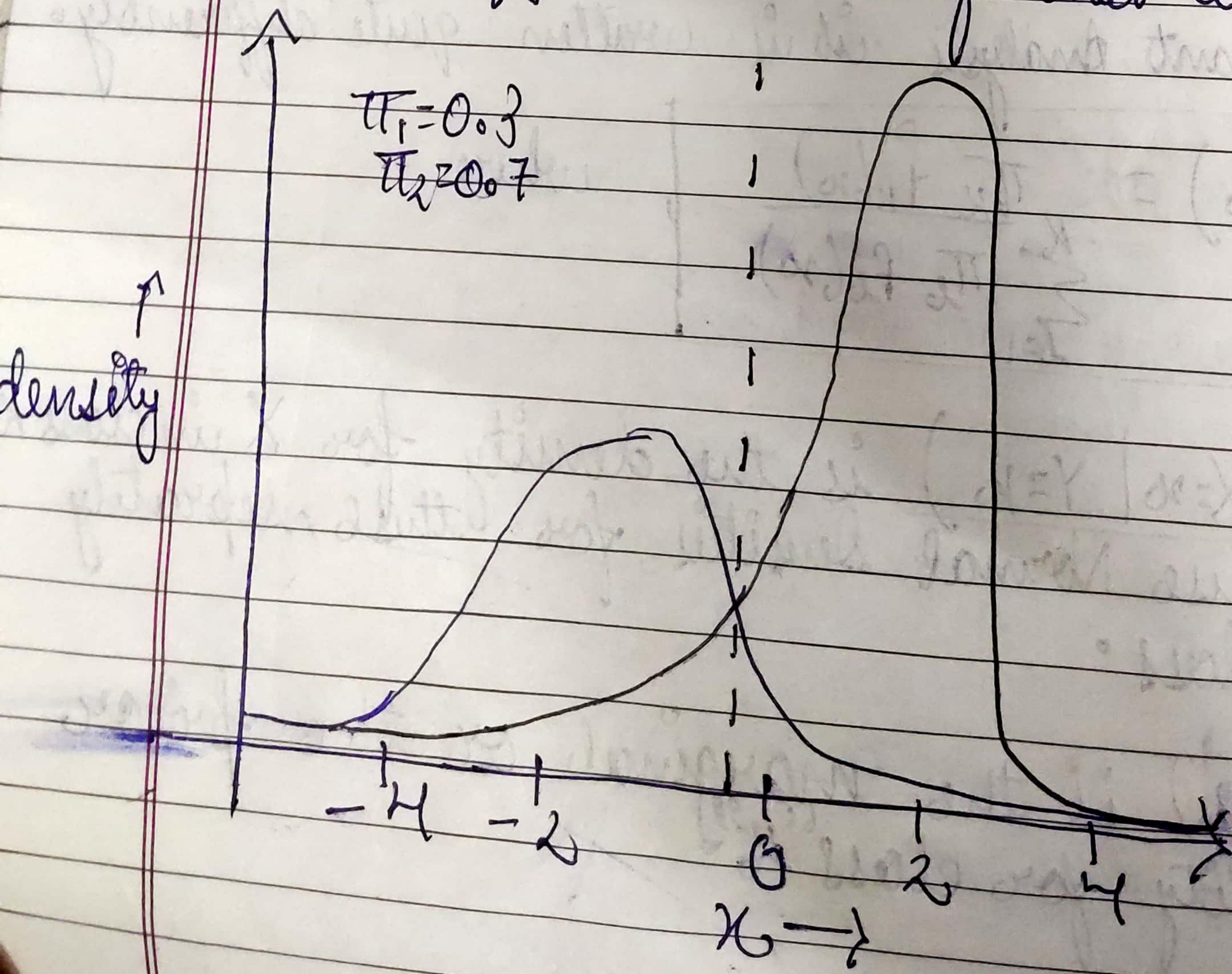
$f_{k|x_0} = \text{Pr}(X=x_0 | Y=k)$  is the density for  $X$  in class  $k$ .  
 Here, we use Normal densities for this, separately in each class.

$\pi_k = P(Y=k)$  is the marginal or the prior probability for class  $k$ .

\* How do we classify on the basis of density function



We classify on the basis of what density is best.



When the priors ( $\pi_k$ ) were different we had to take them into account.

In the second graph we favor the black class. - the decision boundary was shifted to the left.

## Why discriminant analysis?

- When the classes are well separated estimates of logistic regressions are surprisingly unstable. LDA doesn't suffer from this problem. (Because predictors will overshoot to 0)
- If  $n$  is small & distribution of predictors are approximately normal, then obviously LDA will perform better.
- LDA performs better in multiclass-classification, because it also provides low dimensional view of the data.

## Linear Discriminant Analysis when $f=1$ .

When there is only one predictor The Gaussian Density has the form:

$$f_K(x) = \frac{1}{\sqrt{2\pi} \sigma_K} e^{-\frac{(x-\mu_K)^2}{2\sigma_K^2}}$$

LDA takes the assumptions that  $\sigma_K^2 = 0$

Plugging the density function into the Bayes formula gives:

$$P_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu_k}{\sigma})^2}}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu_l}{\sigma})^2}}$$

The formula above looks pretty ugly.

What we can do is take the log of probability function & discard all the terms which are independent of  $\mu$ .

Why can we discard those terms because we focus on what probability is higher & constant terms don't make any change.

This is called a discriminant score and is given by:

$$\delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

\* Note this, the function above is linear in Nature.

We can simplify the application even further when  
 $R=2$  &  $\pi_1 = \pi_2 = 0.5$ .

$f(x)$  will be equal at decision boundary.

$$\chi \frac{m_1}{\sigma^2} - \frac{m_1^2}{2\sigma^2} + \log(\pi_1) = \chi \frac{m_2}{\sigma^2} - \frac{m_2^2}{2\sigma^2} + \log(\pi_2)$$

$$\chi(m_1 - m_2) = \frac{m_1^2 - m_2^2}{2}$$

$$\chi = \frac{m_1 + m_2}{2} \quad \text{at the decision boundary}$$

When  $\pi_1 = \pi_2$  &  $k=2$  classes.

## Estimating the parameters:

$$\hat{\pi}_k = \frac{n_k}{n}$$

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i$$

$$\hat{\sigma}_k^2 = \frac{1}{n-k} \sum_{k=1}^K \sum_{y_i=k} (x_i - \hat{\mu}_k)^2$$

$$= \sum_{k=1}^K \frac{n_k - 1}{n - k} \cdot \hat{\sigma}_k^2$$

where,  $\hat{\sigma}_k^2 = \frac{1}{n_k-1} \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2$  is usual formula for the estimated variance in  $k^{\text{th}}$  class.

# Linear Discriminant Analysis when p > 1

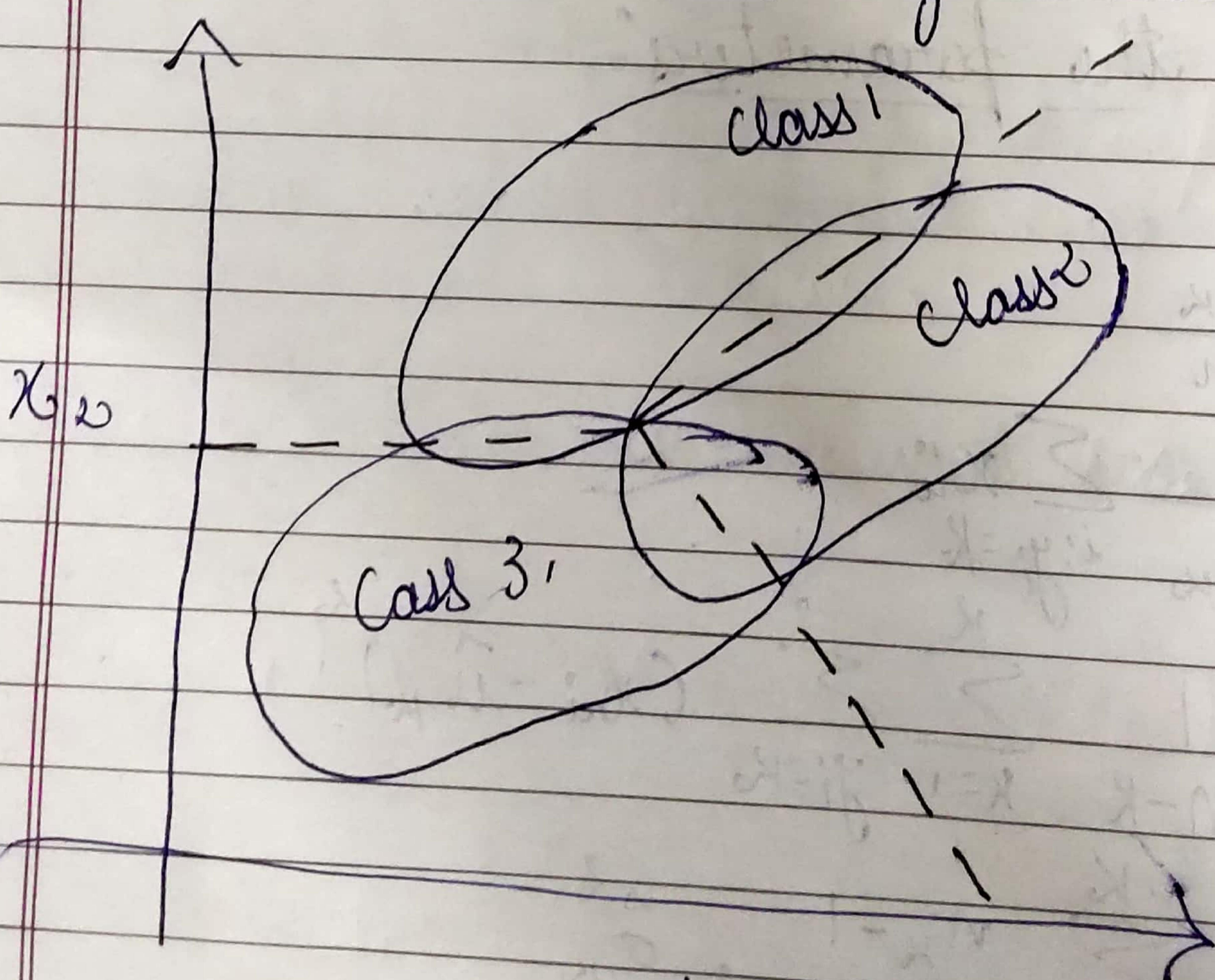
$$\text{Density: } f(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

The density function changes when we move to multivariate data.

where,  $\Sigma$  = Covariance Matrix  
&  $p$  = # predictors.

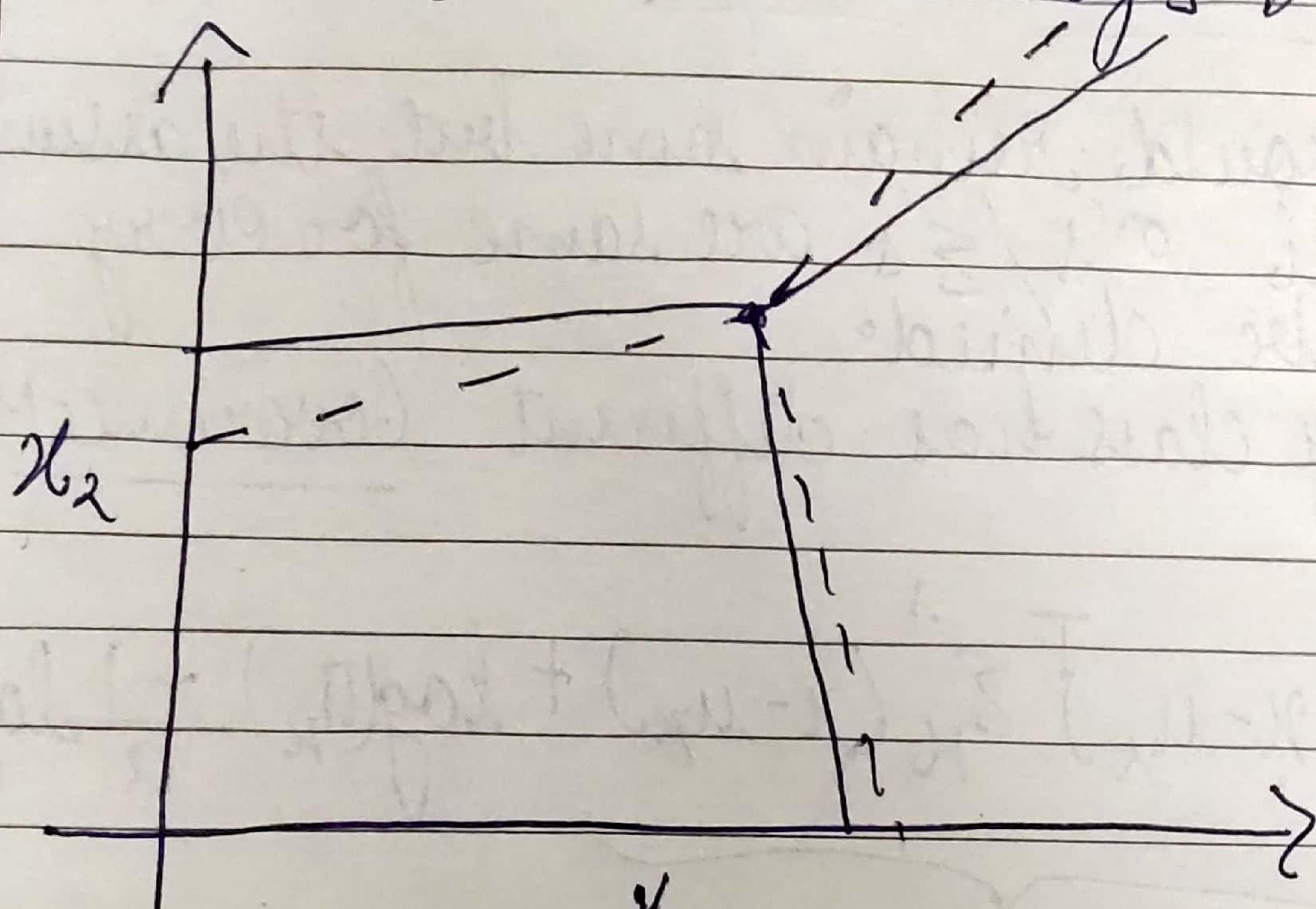
$$\text{Discriminant Function: } x^T \Sigma^{-1} \mu_K - \frac{1}{2} \mu_K^T \Sigma^{-1} \mu_K + b_0$$

This is still linear as you can see.



In the plot above we have shown 3 classes as contours of data rather than a scatterplot

The dashed lines are Bayes Decision Boundary.



In the plot above dashed lines are Bayes decision boundary  
 & solid lines are LDA boundaries.

From  $\hat{\delta}_k(x)$  to probabilities:

Once we have estimates  $\hat{\delta}_k(x)$ , we can turn it into estimates of probability.

$$P(Y=k \mid X=x) = \frac{e^{\hat{\delta}_k(x)}}{\sum_{l=1}^L e^{\hat{\delta}_l(x)}}$$

So, classifying largest  $\hat{\delta}_k(x)$  leads to classifying largest  $P(Y=k \mid X=x)$

## Quadratic Discriminant Analysis

Everything would remain same but the assumption that all  $\sigma^2$ 's/ $\Sigma$ 's are same for every class will be denied.  
Now, every class has different Covariance Matrix.

$$\delta_k(x) = -\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log(\pi_k) - \frac{1}{2} \log |\Sigma_k|$$

This Quadratic is Non-convex.  
Because,  $\Sigma_k$  is different.

## Naive Bayes Classifier

This assumes features are ~~not~~ independent.  
This is useful when  $p$  is large, here methods like LDA & QDA Breakdown.

→ Gaussian Naive Bayes assumes  $\Sigma_k$  is diagonal.

$$\delta_k(x) \propto \log[\pi_k \prod_{j=1}^p f_{kj}(x_j)]$$

$$= -\frac{1}{2} \sum_{j=1}^p \left[ \frac{(x_j - \mu_{kj})^2}{\sigma_{kj}^2} + \log \sigma_{kj}^2 \right] + \log \pi_k$$

when the data is discrete,  $P_{X,Y}(x,y)$  is replaced by probability mass function (histogram)

Logistic Regression  $\checkmark$  vs LDA.

For a two class problem, we can show that

$$\log \left( \frac{p_1(x)}{1-p_1(x)} \right) = \log \left( \frac{p_1(x)}{p_2(x)} \right) \text{ is of the form} \\ = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

so it has same form as logistic regression.

The difference is how the parameters are estimated:

→ Logistic regression uses conditional likelihood based on  $P(Y|X)$  known as discriminative learning.

→ LDA uses full likelihood based on  $P(X, Y)$  known as generative learning.

Logistic regression uses only the distribution of  $Y$   
 LDA uses distribution of both  $X$  &  $Y$

## SUMMARY

- Logistic Regression works well when  $K=2$
- LDA works well when "n" is small &  
    classes are well separated; Also when  $K>$
- Naive Bayes is useful when  $p$  is very large.