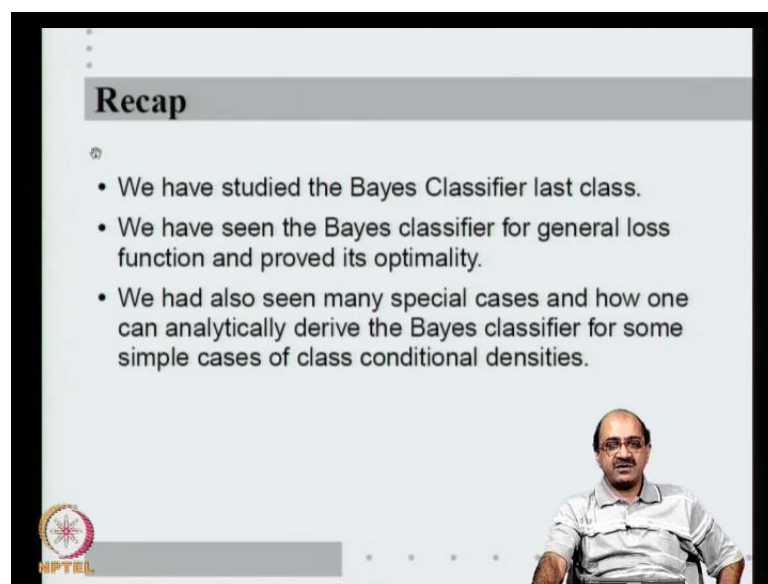


**Pattern Recognition**  
**Prof. P. S. Sastry**  
**Department of Electronics and Communication Engineering**  
**Indian Institute of Science, Bangalore**

**Lecture - 04**  
**Estimating Bayes Error Minimax and Neymann - Pearson Classifiers**

Hello, let us continue with the next class just briefly recapitulate last class, we have studied the Bayes classifier in detail.

(Refer Slide Time: 00:29)



We have derived the Bayes classifier for the general loss function and we have also proved the optimality of the general Bayes classifiers, optimality in the sense that no other classifier can achieve lower risk than Bayes classifiers. So for risk minimization, we showed that the Bayes classifier is the optimal classifier and we have also seen several special cases of what the classifier looks like, for example for 0,1 loss function how it looks like and also for special classes of loss conditional density such as normal and so on. We have analytically derived the Bayes classifiers, so we seen some examples of how one derives Bayes classifier. So, this class we will start with another simple example of deriving Bayes classifier in a slightly different setting. Then we will look at Bayes error, how to calculate the error of Bayes classifier and then, we will move on to a few other criteria for classification.

(Refer Slide Time: 01:26)

**An Example**

Consider another example of deriving Bayes classifier.

- Suppose we have  $K$  classes. The classifier is allowed the option to 'reject' a pattern and this is done by the classifier assigning class  $K + 1$  to the pattern. Define the loss function by

$$L(i, j) = \begin{cases} 0 & \text{if } i = j \text{ and } i, j = 1, \dots, K \\ \rho_m & \text{if } i = 1, \dots, K, \text{ and } i \neq j \\ \rho_r & \text{if } i = K + 1 \end{cases}$$

Now we want to derive the Bayes classifier in terms of the posterior probabilities.

NPTEL

So let us start with an example, let us say we have  $K$  classes and as I said when, we did the Bayes classifier last class that the actions of the classifier need not necessarily have to be class labels they can be more than the class labels. So, let us take an example where, the classifier is allowed the option to reject a pattern that is you look at the pattern and say no I can not classify and let us assume that this is done by the classifier assigning a class  $K$  plus 1 to the pattern.

So, this  $K$  classes, so class label takes values 1 to  $K$  where, the classifier actions take values from 1 to  $K$  plus 1 and will interpret the action  $K$  plus 1 of the classifier as the classifier rejecting the pattern. It is rejecting the pattern because, may be it does not have enough confidence to classify, so now for this case let us put some loss function. So, as you know the loss function has 2 arguments, the first argument is what I call,  $i$  here is the actions of the classifier, second argument of the class labels. So,  $i$  can take values 1 to  $k$  plus 1 where as,  $j$  takes values only 1 to  $K$ . So, is this is something like a 0, 1 loss function. So, if  $i$  did correct, so if both  $i$  and  $j$  belong 1 to  $K$ , that means, the class the classifier has classified the pattern to one of the classes and if  $i$  is equal to  $j$  then loss is 0.

Similarly, if  $i$  is 1 to  $K$  that is the classifier has decided to call a particular class. But,  $i$  is not equal to  $j$  then, we have misclassified let us call that cost  $\rho_m$ ,  $m$  for misclassification and there had irrespective of  $j$ , if  $i$  is  $K$  plus 1, that is the classifier has decided to reject the pattern the loss is  $\rho_r$ . So, this is the loss function, so if  $i$

misclassify a pattern  $x$  suffer a loss of  $r_m$ , if I reject a pattern  $x$  suffer a loss of  $r_r$ , ofcourse, correct classification is 0 loss given this can we now derive the Bayes classifier in terms of the posterior probabilities as usual.

(Refer Slide Time: 03:24)

**Example Contd.**

- Recall that the Bayes classifier is
 
$$h_B(\mathbf{X}) = \alpha_i \text{ if } R(\alpha_i | \mathbf{X}) \leq R(\alpha_j | \mathbf{X}), \forall j.$$

where

$$R(\alpha_i | \mathbf{X}) = \sum_{j=0}^K L(\alpha_i, C_j) q_j(\mathbf{X})$$

- So, we now need to calculate  $R(\alpha_i | \mathbf{X})$  for different actions,  $\alpha_i$  available to the classifier.

The slide also features a video feed of a man in a white shirt and a small NPTEL logo in the bottom left corner.

Recall that the general Bayes classifiers is given an  $X$ , the Bayes classifier will say  $\alpha_i$  or  $i$ , if the risk of that action is less than the risk of all other actions, where risk of that action given  $X$  is defined as,  $L(\alpha_i, C_j) q_j(X)$ , where  $j$  the or  $C_j$  are the class labels. So,  $j$  goes from not 0 to 1, I am sorry 1 to  $K$ , it is not 0 to  $K$ , but 1 to  $K$ , I am sorry about that. So, that  $K$  classes now, we have we we will calculate this risk for various actions and say, when is classification pattern and when is reject pattern.

(Refer Slide Time: 04:03)

- For  $\alpha_i = 1, \dots, K$ , we have  $L(\alpha_i, C_j) = \rho_m$  if  $\alpha_i \neq C_j$  and it is zero otherwise.
- Hence,  $R(i | \mathbf{X}) = \sum_{j \neq i} \rho_m q_j(\mathbf{X}) = \rho_m(1 - q_i(\mathbf{X}))$ .
- Also,  $R(K + 1 | \mathbf{X}) = \sum_j \rho_r q_j(\mathbf{X}) = \rho_r$
- Hence,  $h_B(\mathbf{X}) = i, 1 \leq i \leq K$ , if
 
$$\rho_m(1 - q_i(\mathbf{X})) \leq \rho_m(1 - q_j(\mathbf{X})), \forall j$$
 and
 
$$\rho_m(1 - q_i(\mathbf{X})) \leq \rho_r$$

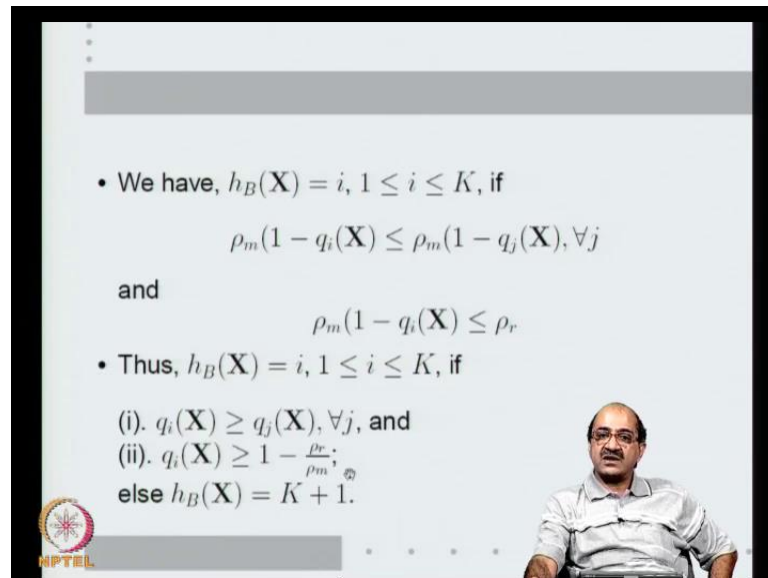
NIPTEL

When the classifier wants to take actions 1 to K that means, classify into one of the classes then the loss is rho m, if i misclassify, otherwise it is 0. So for i i between 1 to K the risk of i given x is rho m, q j z not equal to I, because when i correctly classify there is nonono loss, otherwise all other misclassifications have the same loss rho m.

So, now rho m comes for this examination and it becomes rho m into 1 minus q i X on the other hand, if classifier takes the action K plus 1, then for all class labels j the loss is rho r, so the risk will be summation over j rho or q j, rho r comes with a summation and summed over q j is equal to 1, so this becomes rho r. Here because, this summation over j is not equal to i, summation q j will become 1 minus q i, here this is for all j, so summation q j becomes 1. So, this becomes rho r.

So, what does it mean, when can I call a particular class or i given x, should be better than r j given x, for a all other class labels j and it is also should be better than r K plus 1 given x right. So, my my new my best classifier can say i for one of the class labels, if the risk associated with i, that is rho m into 1 minus q i should be less than or equal to risk associated with any other class label rho m into 1 minus q j, also risk associated with action i should be less than the risk associated with the action K plus 1, which is rho r. So, both these have to be satisfied for me to call class i.

(Refer Slide Time: 05:51)



• We have,  $h_B(\mathbf{X}) = i, 1 \leq i \leq K$ , if

$$\rho_m(1 - q_i(\mathbf{X})) \leq \rho_m(1 - q_j(\mathbf{X})), \forall j$$

and

$$\rho_m(1 - q_i(\mathbf{X})) \leq \rho_r$$

• Thus,  $h_B(\mathbf{X}) = i, 1 \leq i \leq K$ , if

(i).  $q_i(\mathbf{X}) \geq q_j(\mathbf{X}), \forall j$ , and

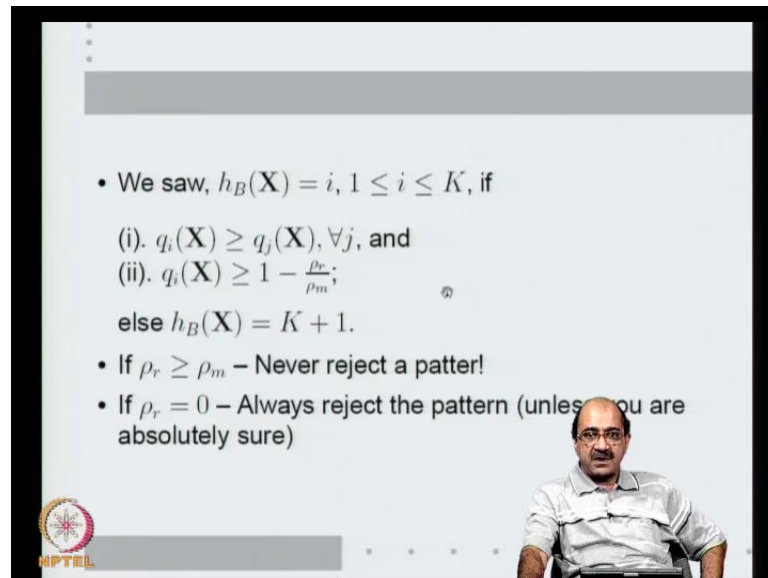
(ii).  $q_i(\mathbf{X}) \geq 1 - \frac{\rho_r}{\rho_m}$ ;

else  $h_B(\mathbf{X}) = K + 1$ .

So, let us simplify this, so the Bayes classifier will, now say a class label  $i$ , if  $\rho_m(1 - q_i(\mathbf{X}))$  is less than  $\rho_m(1 - q_j(\mathbf{X}))$  for all  $j$  as well as  $\rho_m(1 - q_i(\mathbf{X}))$  is less than  $\rho_r$ . The first inequality, I can cancel  $\rho_m$  from both sides and first inequality becomes same as  $q_i \geq q_j$ . So, let us simplify the second also, so I can call  $i$ , the first inequality says  $q_i$  is greater than or equal to  $q_j$  and the second inequality says I will bring  $\rho_m$  this side this is  $1 - q_i$  less than  $\rho_r$  by  $\rho_m$  or if you bring  $q_i$  the other way  $q_i$  is greater than  $1 - \rho_r$  by  $\rho_m$ .

So, if both these conditions are satisfied then I can call a class  $i$ , obviously there will be some  $i$  for, which  $q_i$  is greater than or equal to  $q_j$ , which ever is the highest posterior probability class. But, earlier when I did not have reject the Bayes classifier simply puts it in the class corresponding with the highest posterior probability, but no that is not enough for me to call a class because, I am allowed reject not only, I should be the highest posterior probability class. But, the probability of the highest posterior probability class itself should be greater than some threshold, if this is not true then it is better to call reject. Otherwise,  $h_B(\mathbf{X})$  will be  $K + 1$ , because this inequality are not satisfied, which means the least risk will be for the action  $K + 1$ .

(Refer Slide Time: 07:11)



• We saw,  $h_B(\mathbf{X}) = i, 1 \leq i \leq K$ , if


(i).  $q_i(\mathbf{X}) \geq q_j(\mathbf{X}), \forall j$ , and

(ii).  $q_i(\mathbf{X}) \geq 1 - \frac{\rho_r}{\rho_m}$ ;

else  $h_B(\mathbf{X}) = K + 1$ .

• If  $\rho_r \geq \rho_m$  – Never reject a pattern!

• If  $\rho_r = 0$  – Always reject the pattern (unless you are absolutely sure)



So, let us try and understand it again. So, in the reject case my new Bayes classifier will say a class label  $i$ , if this is true where, the highest posterior probability is greater than 1 minus  $\rho_r$  by  $\rho_m$  and then, I will call the highest posterior probability class. Otherwise I will call  $K + 1$  to understand this let us look at few special cases. Suppose, cost of rejection is greater than cost of misclassification, what does that mean, if I misclassify a pattern, I suffer less loss than, if I reject a pattern then, it should not be there should be no condition under, which reject is good. Now that is what this will tell me, if  $\rho_r$  is greater than  $\rho_m$  then  $\rho_r$  by  $\rho_m$  will be greater than 1. So, 1 minus  $\rho_r$  by  $\rho_m$ , will be negative and hence  $q_i$  being a probability will always be greater than this.

So, then it boils down to whole thing for some  $i$  is of that  $q_i(\mathbf{X})$  is greater than  $q_j$ , I will, call that  $i$ , I will never ever call  $K + 1$  right, so I never reject a pattern. Now consider the other extreme case, suppose cost of rejection is 0, cost of rejection is same as cost of correct classification. So, what should I do, I might just reject everything what does my derivation say if  $\rho_r$  is 0, I will call one of the class labels only if  $q_i(\mathbf{X})$  is greater than or equal to 1.

So, what does that mean, if cost of rejection is 0, I always reject a pattern unless of course, I am absolutely sure if  $q_i(\mathbf{X})$  is equals to 1, then I can call  $i$ , unless I am absolutely sure it is better to reject a pattern because, rejection cost me nothing. So, these

are just extreme cases for us to understand that the senice, so this is another example of how I may derive a Bayes classifier, where classifiers actions may be different from class labels. Here, we have one extra action namely the reject option, so with this example, we will stop discussing examples of Bayes classifiers.

(Refer Slide Time: 09:12)

**Finding Bayes Error**

- Given class conditional densities, the Bayes classifier is easily computed.
- We may also want to compute the Bayes error.
- Gives us the expected performance. Also lets us decide whether we need better features.
- For the case of 0-1 loss function, we need to evaluate

$$\int_{\mathbb{R}^n} \min(p_0 f_0(\mathbf{X}), p_1 f_1(\mathbf{X})) d\mathbf{X}$$

- In general, a difficult integral to evaluate

Let us move to one other important issue with Bayes classifiers, how do I find the Bayes error from all the examples, we have considered, so far given class conditional densities the Bayes classifier is easily computed right. For various loss functions, we know how to compute it, now we may also want to compute the Bayes error, because that tells me, what is the expected error of the classifier, what is the expected risk of the classifier.

Now for example, if the expected risk is not within acceptable limits accept expected error rate is not within acceptable limits, then I may have to rethink my whole problem in the sense, I may want to get better features because, with these features this is the best performance I can get. So, estimating or finding Bayes error is useful for me to know whether my classifier will meet the specification requirements. So for the 0, 1 loss function, we have already seen when, we derived the optimality of Bayes classifier, that this is the error. So, we have to integrate minimum of  $p_0 f_0, p_1 f_1$  over the entire feature space to find the error rate of the classifier, this is the probability of misclassification by the Bayes classifier. In general it is a very difficult integral to evaluate, because it is a mean inside the integrand.

(Refer Slide Time: 10:27)


• Let us consider the simplest case:  
2-class problem,  $X \in \mathcal{R}$ , normal class conditional densities and 0-1 loss function.

• Assume equal priors. Let  $\sigma_0 = \sigma_1 = \sigma$  and  $\mu_0 < \mu_1$ .

• Then  $h_B(X) = 0$  if  $X < (\mu_0 + \mu_1)/2$ .

• Then, Bayes error is

$$P(\text{error}) = 0.5 \int_{-\infty}^{\frac{\mu_0 + \mu_1}{2}} f_1(X) dX + 0.5 \int_{\frac{\mu_0 + \mu_1}{2}}^{\infty} f_0(X) dX$$

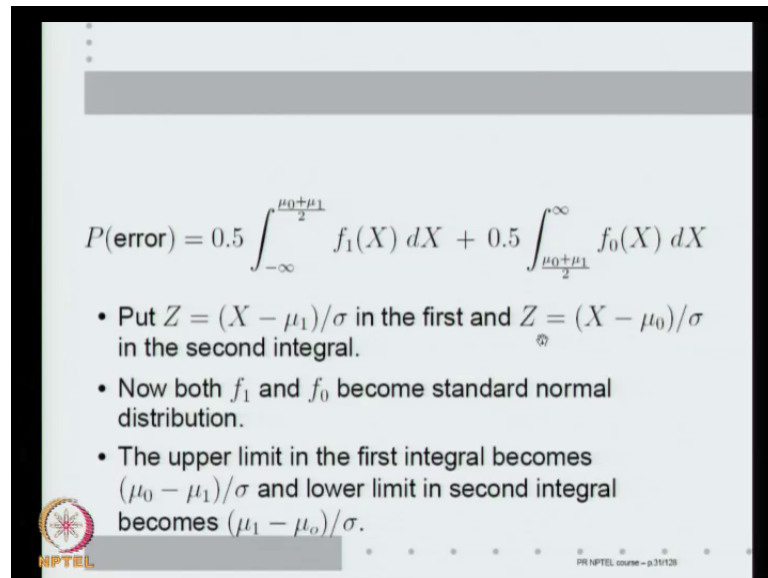
 PR NPTEL course - p27128

First let us look at a very simple case, let us consider one dimensional feature space, 2 class problem 0, 1 loss function, assume equal priors assume normal class conditional densities with equal variance. This is about the simplest special case, you can consider and for notational convenience, let us assume the mean of class 0 is less than mean of class 1. In this case, we have already known, because the both variances are same the threshold for the Bayes classifier is midway between the 2 means. So, if  $X$  is less than  $\mu_0 + \mu_1$  by 2, I will call class 0, if it is greater than that, I will call class 1, this is the Bayes classifier. So, what will be the error of the Bayes classifier.

So, for  $X$  less than  $\mu_0 + \mu_1$  by 2, I will always call 0. So, the probability of class 1 patterns coming with  $X$  less than this is 1 half of the error. So, probability that  $X$  belongs to class 1 and  $X$  is less than  $\mu_0 + \mu_1$  by 2 is this probability this integral, integral of the  $f_1$  the class 1 class conditional density over this range. Similarly, the other error occurs, when a class 0 pattern comes with  $X$  greater than  $\mu_0 + \mu_1$  by 2. So, this is the Bayes error in for normal density, it is easy to evaluate, so let us evaluate it.




(Refer Slide Time: 11:49)



$$P(\text{error}) = 0.5 \int_{-\infty}^{\frac{\mu_0 + \mu_1}{2}} f_1(X) dX + 0.5 \int_{\frac{\mu_0 + \mu_1}{2}}^{\infty} f_0(X) dX$$

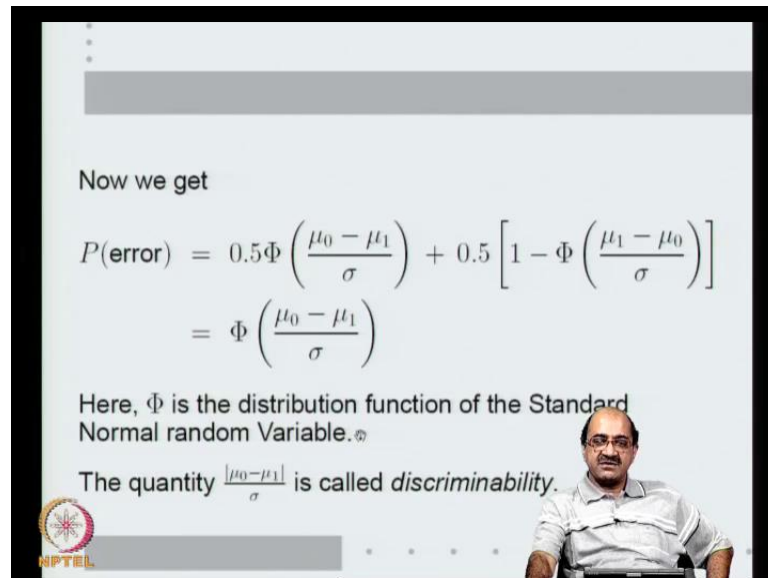
- Put  $Z = (X - \mu_1)/\sigma$  in the first and  $Z = (X - \mu_0)/\sigma$  in the second integral.
- Now both  $f_1$  and  $f_0$  become standard normal distribution.
- The upper limit in the first integral becomes  $(\mu_0 - \mu_1)/\sigma$  and lower limit in second integral becomes  $(\mu_1 - \mu_0)/\sigma$ .

 NPTEL course - p31128

So, this is the error integral  $f_1$  is normal with mean  $\mu_1$  and variance  $\sigma^2$  and  $f_0$  is normal with mean  $\mu_0$  on variance  $\sigma^2$ . So, using the standard substitution by change the variable  $X$  in this integral to  $Z$ , where  $Z$  is  $X$  minus  $\mu_1$  by  $\sigma$ , then this density assume the form of the standard normal density. Similarly, for this integral because,  $f_0$  is normal with mean  $\mu_0$  on variance  $\sigma^2$ , if I use the substitution  $Z$  is equal to  $X$  minus  $\mu_0$  by  $\sigma$ .

This becomes a standard normal density integral, now what happens to the limits, when a put  $Z$  this, when  $X$  goes up to  $\mu_0 + \mu_1$  by  $2$ ,  $Z$  goes to  $\mu_0 - \mu_1$  by  $\sigma$  and similarly, here. So, this becomes an integral of the standard normal density over minus infinity to  $\mu_0 + \mu_1$ ,  $\mu_0 - \mu_1$  by  $\sigma$  and this becomes from  $\mu_1 - \mu_0$  by  $\sigma$  to infinity of another standard normal.

(Refer Slide Time: 12:52)



Now we get

$$P(\text{error}) = 0.5\Phi\left(\frac{\mu_0 - \mu_1}{\sigma}\right) + 0.5\left[1 - \Phi\left(\frac{\mu_1 - \mu_0}{\sigma}\right)\right]$$
$$= \Phi\left(\frac{\mu_0 - \mu_1}{\sigma}\right)$$

Here,  $\Phi$  is the distribution function of the Standard Normal random Variable.

The quantity  $\frac{|\mu_0 - \mu_1|}{\sigma}$  is called *discriminability*.

The error density is this, where phi is the distribution of the standard normal density. So, the first integral is from minus infinity to mu 0 minus mu 1 by sigma. So, that is phi of mu 0 minus mu 1 by sigma, second integral is from mu 1 minus mu 0 by sigma to infinity. So, this is 1 minus phi of mu 1 minus mu 0 by sigma, Because, the distributional function standard normal is symmetric phi of 1 minus phi x is equal to phi of minus x right, what is in this big brackets here is same as this. So, the 0.5 goes away, so the error becomes this. So, for equal variants both class conditional densities being normal, this is the Bayes error.

So, essentially a Bayes error depends on mu 0 minus mu 1 by sigma, as you would expect, when sigma is same, I am just putting the point midway. So, how much error I make depends on how much the means are separated, if means are separated by a large amount then, I will make less error, if means are separated by a small amount, I will make more error and small and large amount is relative to the variants of the distribution right.

So, that is what this expression denotes, the quantity mu 0 minus mu 1 by sigma, where mu 0 minus absolute value of mu 0 minus mu 1 by sigma is called the discriminability this, when this quantity is large right, note that, we are assuming mu 0 less than mu 1. So, when discriminability is large this should be phi of some large negative quantity, so close to 0.

So, if  $\mu_0$  and  $\mu_1$  are separated by a large amount relative to the variance, then I make very small error, if they are separated by a small amount, then I will make a large error the Bayes error. So, that is why, this quantity is called the discriminability, now this is for a very special one dimensional case, you know normal densities with equal variance.

(Refer Slide Time: 14:49)

- In the general case, we need to evaluate  

$$P(\text{error}) = \int_{\mathbb{R}^n} \min(p_0 f_0(\mathbf{X}), p_1 f_1(\mathbf{X})) d\mathbf{X}$$
- A useful inequality here is  

$$\min(a, b) \leq a^\beta b^{1-\beta}, \forall a, b \geq 0, 0 \leq \beta \leq 1.$$
- Easy to prove. Suppose  $a < b$   

$$a^\beta b^{1-\beta} = a^{-1+\beta} b^{1-\beta} a = \left(\frac{b}{a}\right)^{1-\beta} a \geq a = \min(a, b)$$
- Hence we have (for 0-1 loss function)

$$P(\text{error}) \leq p_0^\beta p_1^{1-\beta} \int_{\mathbb{R}^n} f_0^\beta(\mathbf{X}) f_1^{1-\beta}(\mathbf{X}) d\mathbf{X}$$

NPTEL logo on the left and 'PR NPTEL course - p.39/28' on the bottom right.

What about the general case, in general, we have to evaluate this integral as the already said mean inside the integrand is often difficult to evaluate. But, I can use, I very standard and useful inequality, for any 2 real numbers a and b mean of a comma b can always be bounded above by a to the power beta into b to the power 1 minus beta, for any beta between 0 and 1. I here, we are assuming that both a and b are positive numbers, this bond is not difficult to prove.

So, let us prove this, let us suppose a is less than b, now what is a power beta into B power 1 minus beta, I can write it as a power minus 1 plus beta, that is divided by a and then multiply by a. So, this becomes b by a, whole to the power 1 minus beta into a, now I am assuming b greater than a. So, b by a is greater than 1, so b by a to the power 1 minus beta some quantity greater than 1. So, when you multiply that with a, I get some quantity greater than 1 and a is ofcourse mean of a comma b.

So, this shows that mean of a comma b is bounded above by a power beta into b power 1 minus beta. So, I can reduce this integral to, I can write this mean bound this mean by

this to the power beta and this to power 1 minus beta, that becomes p 0 to the power beta p 1 to the power 1 minus beta integral of f 0 to the power beta and f 1 to the power 1 minus beta ah. This is slightly easier integral to evaluate than, this because, I do not have mean inside, it is the normal density raised to some fractional power.

(Refer Slide Time: 16:21)

Suppose  $f_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  and  $f_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ . Then we can show

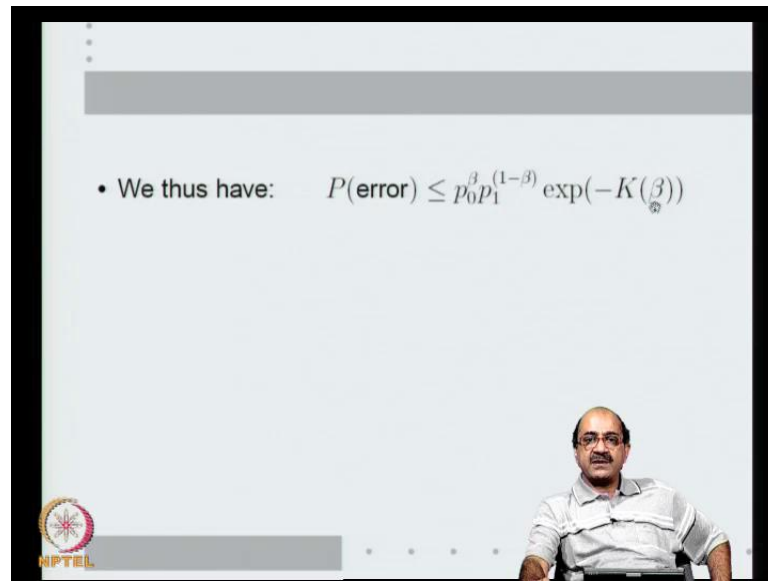
$$\int f_0^\beta(\mathbf{X}) f_1^{1-\beta}(\mathbf{X}) d\mathbf{X} = \exp(-K(\beta))$$

where

$$K(\beta) = \frac{\beta(1-\beta)}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^t (\beta \boldsymbol{\Sigma}_0 + (1-\beta) \boldsymbol{\Sigma}_1)^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \frac{1}{2} \ln \left( \frac{|\beta \boldsymbol{\Sigma}_0 + (1-\beta) \boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|^\beta |\boldsymbol{\Sigma}_1|^{1-\beta}} \right)$$

Even then it is a difficult integral to evaluate, but suppose, if I assume that f 0 and f 1 are some multidimensional Gaussian densities f 0 has mean mu 0 and covariance matrix sigma 0. F 1 has mean mu 1 and covariance matrix sigma 1, then one can show that this integral can be bounded above by e power minus K beta, where that K beta is some involved expression like this essentially some kind of a quadratic form, involving mu us and sigma 0. But, anyway this can be shown it is say just say algebraically difficult, but otherwise, this derivation is straight forward. So, I can show this.

(Refer Slide Time: 17:03)

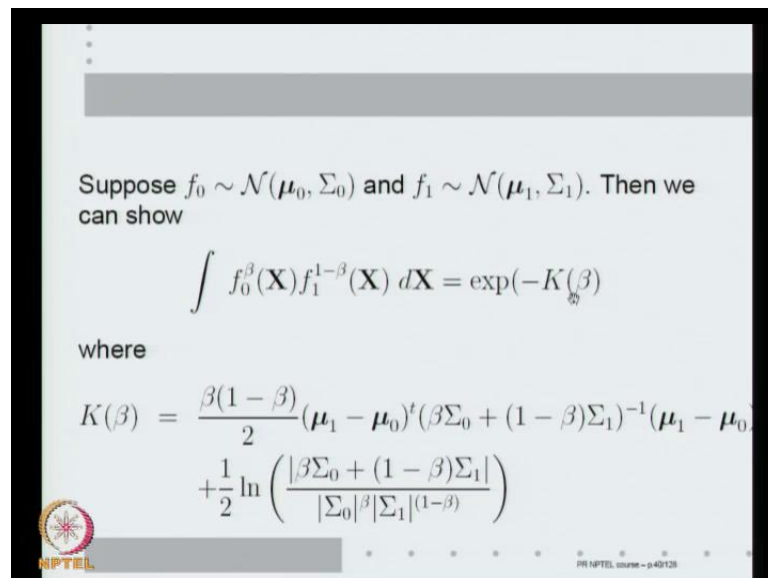


• We thus have:  $P(\text{error}) \leq p_0^\beta p_1^{(1-\beta)} \exp(-K(\beta))$

The slide features a presenter in the bottom right corner and an NPTEL logo in the bottom left corner.

So, what does this mean, that P error is less than or equal to p 0 power of beta p 1 to the power 1 minus beta exponential minus K beta.

(Refer Slide Time: 17:09)



Suppose  $f_0 \sim \mathcal{N}(\mu_0, \Sigma_0)$  and  $f_1 \sim \mathcal{N}(\mu_1, \Sigma_1)$ . Then we can show

$$\int f_0^\beta(\mathbf{X}) f_1^{1-\beta}(\mathbf{X}) d\mathbf{X} = \exp(-K(\beta))$$

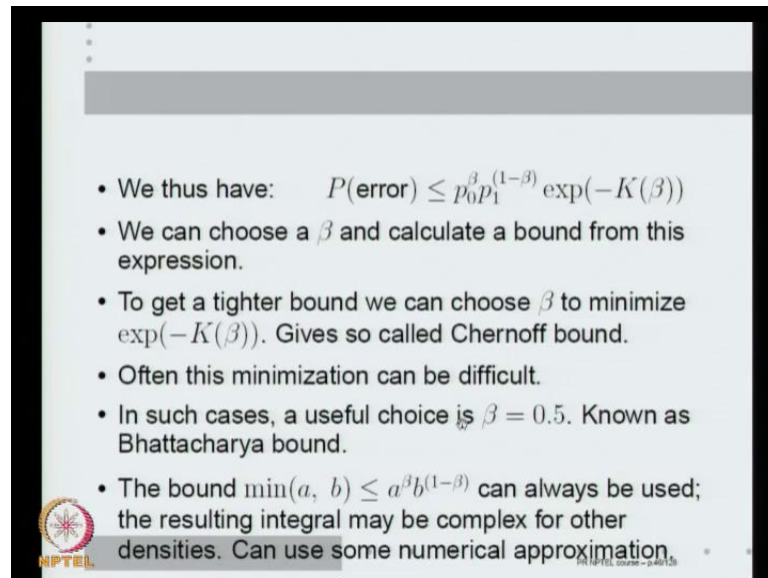
where

$$K(\beta) = \frac{\beta(1-\beta)}{2} (\mu_1 - \mu_0)^t (\beta \Sigma_0 + (1-\beta) \Sigma_1)^{-1} (\mu_1 - \mu_0) + \frac{1}{2} \ln \left( \frac{|\beta \Sigma_0 + (1-\beta) \Sigma_1|}{|\Sigma_0|^\beta |\Sigma_1|^{(1-\beta)}} \right)$$

The slide includes an NPTEL logo in the bottom left and a course ID 'PR NPTEL course - p 49128' in the bottom right.

Where, the K beta term is given by this, which I can calculate, if I know mu 1 mu 0, sigma 0 sigma 1.

(Refer Slide Time: 17:17)

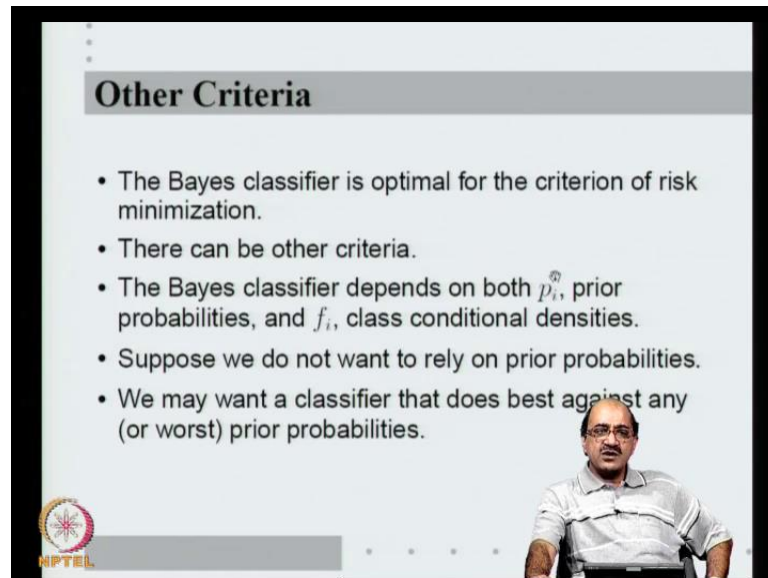


- We thus have:  $P(\text{error}) \leq p_0^\beta p_1^{(1-\beta)} \exp(-K(\beta))$
- We can choose a  $\beta$  and calculate a bound from this expression.
- To get a tighter bound we can choose  $\beta$  to minimize  $\exp(-K(\beta))$ . Gives so called Chernoff bound.
- Often this minimization can be difficult.
- In such cases, a useful choice is  $\beta = 0.5$ . Known as Bhattacharya bound.
- The bound  $\min(a, b) \leq a^\beta b^{(1-\beta)}$  can always be used; the resulting integral may be complex for other densities. Can use some numerical approximation.

Now, how can I use this bound, this bound is true for all beta between 0 and 1. So, for example, I can ask, which beta will give me the tightest bond. So, we can choose a beta and calculate a bound from this expression. A bond calculated like that is often called Chernoff bound, we can get a tighter bound by choosing beta, that minimizes this expression.

Such a bound is called a Chernoff bound, if you do not want to do all that work, in practice a beta that often works is beta is equal to 0.5 and the bound and the Bayes error obtained through this expression, where I choose beta to be 0.5 is known as the Bhattacharya bound. There is another bound on the Bayes error. Of course, in general this bound can always be used though, for general class conditional densities, I would not have this exponential minus K beta, I will have that actual integral and we need to know how to evaluate it. But, even for normal class conditional density as you can see, we can only bound the Bayes error, it is not easy, to actually compute the Bayes error. But, this this is a this is one way in, which I can estimate the Bayes error, either using Chernoff bounds or the Bhattacharya bound all right.

(Refer Slide Time: 18:40)



The slide is titled "Other Criteria" and contains the following text:

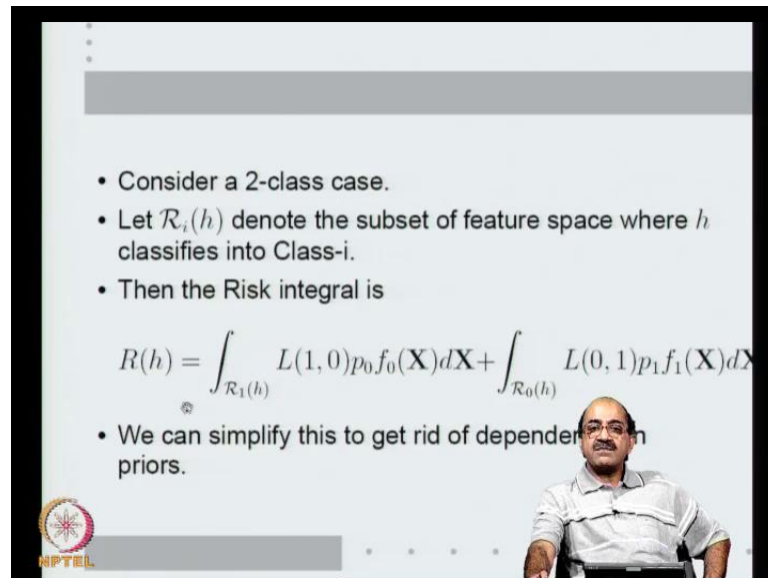
- The Bayes classifier is optimal for the criterion of risk minimization.
- There can be other criteria.
- The Bayes classifier depends on both  $p_i$ , prior probabilities, and  $f_i$ , class conditional densities.
- Suppose we do not want to rely on prior probabilities.
- We may want a classifier that does best against any (or worst) prior probabilities.

In the bottom right corner of the slide, there is a small inset image of a man with glasses and a light-colored shirt, sitting at a desk. In the bottom left corner, there is a logo for "NIPTEL" with a circular emblem.

Let us move on from Bayes classifier, Bayes classifier is optimal for the criteria of risk minimization. But, risk minimization is only 1 criteria right, we can have many other criteria, what does risk there is one thing about Bayes classifiers, it depends both on the prior probabilities and class conditional densities. Now very often, I may not know prior probabilities, out there in the field, which patterns will come I may not know though, I may be able to estimate the class conditional densities, were sometimes, we may want a classifier that does well against any worst kind of prior probabilities. So, without knowing what is the prior probabilities is I do not want my Bayes classifiers to depend on priors.

Because, one day I might have to work with predominantly 1 class patterns, another day I may have to do work with predominantly another class patterns. So, I can ask minimizing risk is not what I want, I want a classifier that has the best risk against the worst possible prior probabilities. Now, this ofcourse, would not be the Bayes classifiers, because I do not know the priors, let us just intuitively see, what this will mean.

(Refer Slide Time: 19:56)



- Consider a 2-class case.
- Let  $\mathcal{R}_i(h)$  denote the subset of feature space where  $h$  classifies into Class- $i$ .
- Then the Risk integral is

$$R(h) = \int_{\mathcal{R}_1(h)} L(1,0)p_0f_0(\mathbf{X})d\mathbf{X} + \int_{\mathcal{R}_0(h)} L(0,1)p_1f_1(\mathbf{X})d\mathbf{X}$$

- We can simplify this to get rid of dependence on priors.

So, let us say, we will take a 2 class case as earlier let us say  $\mathcal{R}_i(h)$  denotes the subset of feature space, where the classifier  $h$  will put things in class  $i$ , that is  $\mathcal{R}_0$  is the region of class 0,  $\mathcal{R}_1$  is the region of class 1 and by region, I mean not the actual region of class 1 or class 0 feature vectors. But, that subset of the feature space, where the classifier  $h$  will put the patterns in that particular class.

Then, if from what we derived earlier the risk integral is the probability that a class 0 pattern comes into a region where,  $h$  will put in class 1 and the probability that a class 1 pattern will come into a region that  $h$  will put in class 0, this is the same integral of that we got earlier. Now, what we want to do is we want to manipulate this expression. So, that it becomes independent of the prior probabilities.



(Refer Slide Time: 20:53)

• Using  $p_0 = 1 - p_1$ , we get

$$R = \int_{\mathcal{R}_1} L(1,0)p_0 f_0(\mathbf{X}) d\mathbf{X} + \int_{\mathcal{R}_0} L(0,1)p_1 f_1(\mathbf{X}) d\mathbf{X}$$
$$= L(1,0)p_0 \int_{\mathcal{R}_1} f_0(\mathbf{X}) d\mathbf{X} +$$
$$L(0,1)(1 - p_0) \int_{\mathcal{R}_0} f_1(\mathbf{X}) d\mathbf{X}$$

The slide includes the NPTEL logo in the bottom left corner and a small inset image of a man in the bottom right corner.

We can do that because, we know that  $p_0$  is  $1 - p_1$  or  $p_1$  is  $1 - p_0$ , so we can eliminate one of them. So, this is my risk integral over  $\mathcal{R}_1$ , it is  $L(1,0), p_0, f_0$ , over  $\mathcal{R}_0$ , it is  $L(0,1), p_1, f_1$ , now I can eliminate one of  $p_0$  and  $p_1$  let us say I will substitute  $p_1$  is equal  $1 - p_0$ , so that is my risk integral. Now, this integral has one term, which is constant  $L(0,1)$  into integral of  $f_1$  over  $\mathcal{R}_0$  and another term that depends on  $p_0$ ,  $p_0$  into this the this first integral minus the second integral.

(Refer Slide Time: 21:34)

• Thus we get

$$R = \int_{\mathcal{R}_1} L(1,0)p_0 f_0(\mathbf{X}) d\mathbf{X} + \int_{\mathcal{R}_0} L(0,1)p_1 f_1(\mathbf{X}) d\mathbf{X}$$
$$= L(0,1) \int_{\mathcal{R}_0} f_1(\mathbf{X}) d\mathbf{X} +$$
$$p_0 \left[ L(1,0) \int_{\mathcal{R}_1} f_0(\mathbf{X}) d\mathbf{X} - L(0,1) \int_{\mathcal{R}_0} f_1(\mathbf{X}) d\mathbf{X} \right]$$

The slide includes the NPTEL logo in the bottom left corner and the text "PR NPTEL course - p.5/128" in the bottom right corner.

So, I can write the risk as 1 constant term plus p 0 times difference of 2 integrals. Now, for a classifier, we choose the regions  $R_0$  and  $R_1$  in such a way that, this expression in these big brackets goes to 0, that classifier's risk is independent of priors right. The way the risk is written, if there is a classifier, which chooses class 1 and class 0 decision regions  $R_1$  and  $R_0$  in such a way that this expression becomes 0 right. The second term in this expression has become 0, for that classifier the risk will be independent of the priors right.

(Refer Slide Time: 22:23)

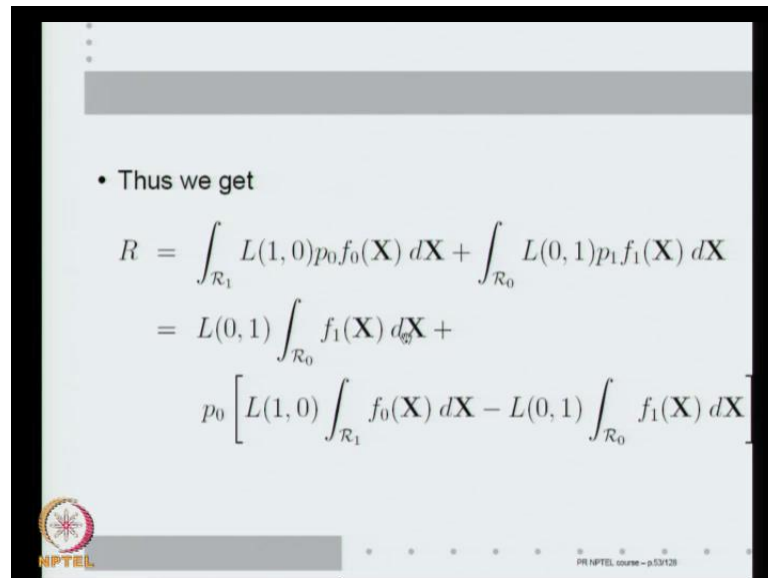
**Minmax Classifier**

- Consider a classifier such that

$$L(1, 0) \int_{R_1} f_0(\mathbf{X}) d\mathbf{X} = L(0, 1) \int_{R_0} f_1(\mathbf{X}) d\mathbf{X}$$

So, consider a classifier, which chooses any classifier, what do you mean by design of a classifier, once we design a classifier, I have designed a function from the feature space to the set 0 1, because we are considering a 2 class problems. So, which means each classifier, simply assigns some subset of the feature space where, it will if a feature vector falls in that subset, I will call class 0, similarly, the remaining substrate, I will call it class 1. So, if design of every classifier is simply choosing a region  $R_1$  where, I will call class 1 and choosing a region  $R_0$  where, I will call class 0. So, a classifier designed in such a way that the regions  $R_0$  and  $R_1$  are so chosen. So, that this equation is satisfied.

(Refer Slide Time: 23:07)



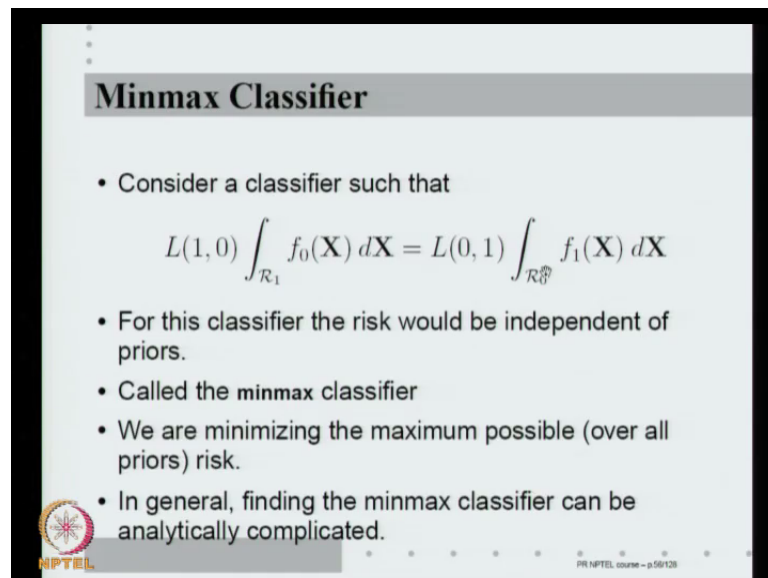
• Thus we get

$$R = \int_{\mathcal{R}_1} L(1,0)p_0f_0(\mathbf{X}) d\mathbf{X} + \int_{\mathcal{R}_0} L(0,1)p_1f_1(\mathbf{X}) d\mathbf{X}$$
$$= L(0,1) \int_{\mathcal{R}_0} f_1(\mathbf{X}) d\mathbf{X} +$$
$$p_0 \left[ L(1,0) \int_{\mathcal{R}_1} f_0(\mathbf{X}) d\mathbf{X} - L(0,1) \int_{\mathcal{R}_0} f_1(\mathbf{X}) d\mathbf{X} \right]$$

NPTEL PR NPTEL course - p.53128

Where did I get this equation from this is nothing but, the term here right. I wanted to make this term 0.

(Refer Slide Time: 23:13)



### Minmax Classifier

- Consider a classifier such that

$$L(1,0) \int_{\mathcal{R}_1} f_0(\mathbf{X}) d\mathbf{X} = L(0,1) \int_{\mathcal{R}_0} f_1(\mathbf{X}) d\mathbf{X}$$

- For this classifier the risk would be independent of priors.
- Called the minmax classifier
- We are minimizing the maximum possible (over all priors) risk.
- In general, finding the minmax classifier can be analytically complicated.

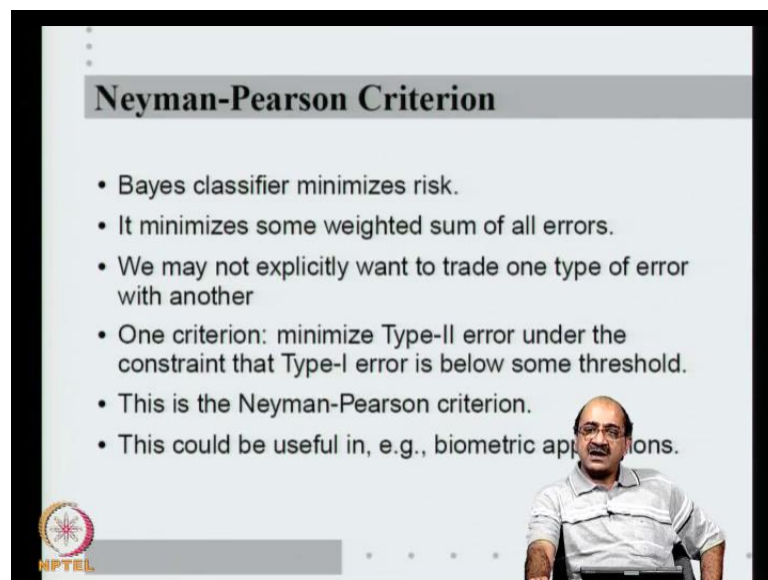
NPTEL PR NPTEL course - p.59128

So, that term will become 0, if this equation is satisfied for such a classifier the risk would be independent of priors, this classifier is known as minmax classifier, because it can be shown it is enough to see that, we are minimizing the maximum possible risk or all possible priors, because we are canceling out the prior dependence and risk. So, we

are budgeting for the maximum possible risk where, maximum is over all possible priors for the same class conditional densities.

Of course, finding analytically a classifier that satisfies such expression is not easy, in general finding minmax classifiers is an analytically complicated issue. But, the purpose of mentioning minmax classifier here is just to say that risk minimization is not necessarily the only criterion, we can have when, we are looking for classifiers, here is another example of a classifier. Which is different from Bayes, but it has its own optimality criterion the minmax classifier, which minimizes the maximum possible risk, where maximum is over all possible priors.

(Refer Slide Time: 24:21)



**Neyman-Pearson Criterion**

- Bayes classifier minimizes risk.
- It minimizes some weighted sum of all errors.
- We may not explicitly want to trade one type of error with another
- One criterion: minimize Type-II error under the constraint that Type-I error is below some threshold.
- This is the Neyman-Pearson criterion.
- This could be useful in, e.g., biometric applications.

The slide includes a small inset image of a man in a white shirt and glasses, and a logo for NPTEL in the bottom left corner.

Let us look at 1 more criterion, this is also a very famous criterion called Neyman Pearson criterion. To understand this criterion, let us go back to Bayes classifiers again Bayes classifier minimizes risk, what is risk, risk is expectation of loss. So, each loss is what, I pay for an error and because, in expectation it is some weighted sum of the errors right, weighted sum of probability of errors weighted sum of losses. So, if I classify a class 0 pattern as class 1 pattern, there is some cause associated with it class 1 pattern as, class 0 pattern, there is some other cause associated with it I find weighted sum of all such cause and asking, which classifier minimizes it and that happens to be the Bayes classifiers.

But, the 2 kinds of errors may not always be tradable, this assumes that, if cost of 1 kind of error is 3 times cost of another kind of error, we are saying it is better to make 2 errors of 1 kind than, 1 error of this kind right. That is the trade of we are doing in minimizing risk. But, there would be situations where, we may not want to trade, 1 type of error with another type of error.

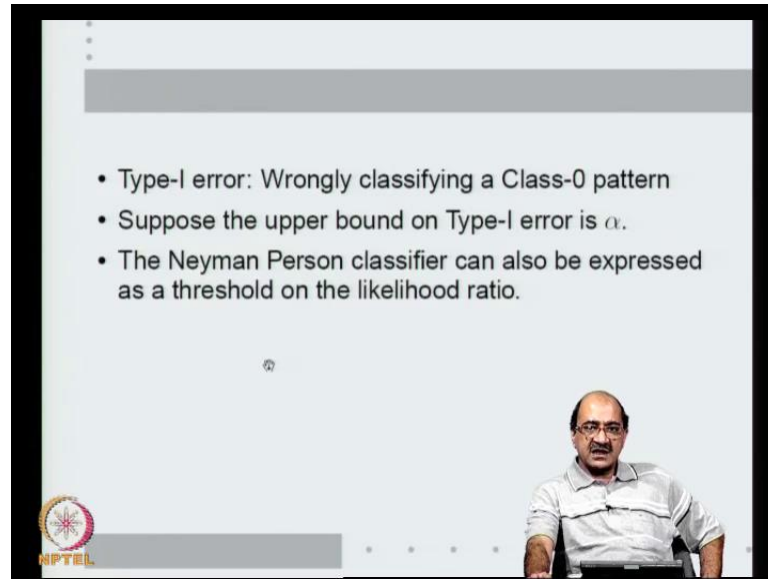
So for example, instead of trading errors, we can say for a fixed constraint and type 1 error minimized type 2 error, recall that type 1 error is wrongly classifying a class 0 pattern, why would this may be possible, let us suppose you are in a biometric application, suppose you are authenticating identity of somebody. So, there are 2 kinds of errors, when somebody is an imposter allowing him access is 1 kind of error, a authorized person not being allowed access is another kind of error.

Now, these 2 kinds of errors are qualitatively different and I may have do not want to find optimal by saying, so many time this error plus, so many time that error should be minimized. On the other hand I may say that I do not want more than 1 percent more than 0.1 percent of the time, an unauthorized person gaining access, while maintaining that can you give me, the best possible error rate for the other kind right.

So, I will put a particular threshold for type one error, that is I do not want more than 1 in 1000 more than once in 1000 times an unauthorized person should gain access. And among all classifiers, that meet this specification, I want a classifier, which minimizes the error of throwing away an authorized person, because, throwing away an authorized person is only an irritation.

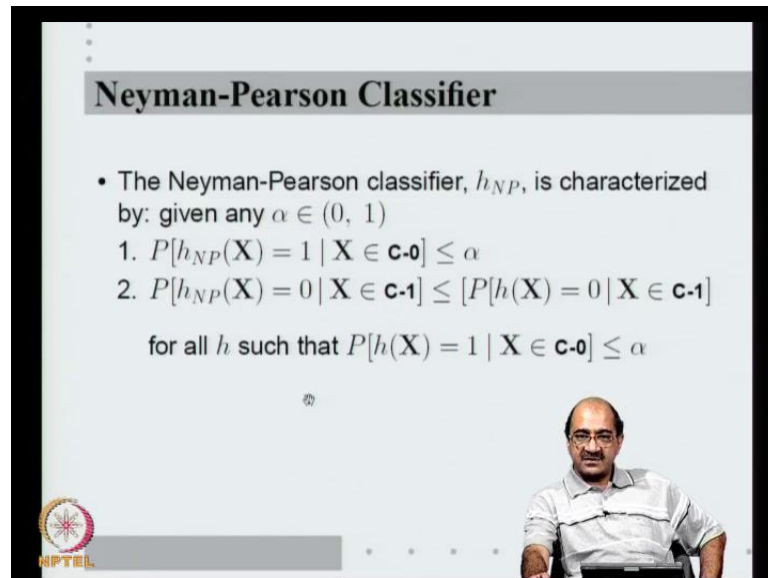
So, there are applications where, you do not want to trade errors of kind with another, I may want to put a threshold on the error of 1 kind, that is I want error of one kind should not exceed a probability of error of 1 kind should not exceed something and given that I want to then then minimize the error kind of error. So, it is as I said it is generally useful in biometric applications where, as I said I may want to put an absolute bound on how often, I may allow an unauthorized person to gain access and while, I am satisfying this specification. I want to minimize the number of times, I will throw away an authorized person.

(Refer Slide Time: 27:40)



So, type 1 error is let us say wrongly classifying class 0 pattern and let us say upper bound and probability of type 1 error is alpha. So, the Neyman Pearson classifier is a classifier, that achieves a bound of alpha and type 1 error, that is the probability of type 1 error by Neyman Pearson classifier is less than or equal to alpha. And in addition, it minimizes the other kind of error. A matter of fact Neyman Pearson classifier can also be expressed as threshold on likelihood ratio as we have seen in Bayes case, it is simply a ratio on the posterior Probabilities are class conditional densities. We have put a threshold on this ratio Neyman Pearson classifier can also be expressed and we will see how.

(Refer Slide Time: 28:27)



**Neyman-Pearson Classifier**

- The Neyman-Pearson classifier,  $h_{NP}$ , is characterized by: given any  $\alpha \in (0, 1)$ 
  1.  $P[h_{NP}(X) = 1 | X \in \mathbf{c-0}] \leq \alpha$
  2.  $P[h_{NP}(X) = 0 | X \in \mathbf{c-1}] \leq [P[h(X) = 0 | X \in \mathbf{c-1}]$for all  $h$  such that  $P[h(X) = 1 | X \in \mathbf{c-0}] \leq \alpha$

NIPTEL

©

So, let us first define the Neyman Pearson classifier, the Neyman Pearson classifiers let us call it  $h_{NP}$ , we were calling Bayes classifier as  $h_b$ , so we will call a Neyman Pearson classifier  $h_{NP}$ . So, given any  $\alpha$  in  $(0, 1)$ , this is the upper bound on the type 1 error what does  $h_{NP}$  have to satisfy. Firstly, probability  $h_{NP}(X) = 1$  given  $X$  belongs to  $\mathbf{c-0}$ , that is wrongly classifying a class 0 pattern is bounded above by  $\alpha$  right.

This is 1 thing that  $h_{NP}$  has to satisfy, then what is it have to satisfy, wrongly classifying a class 1 pattern, that is the other kind of error  $h_{NP}(X) = 0$  given  $X$  belongs to  $\mathbf{c-1}$ . That should be, less than the probability of wrongly classifying a class 1 pattern by any other classifier  $h$ . But, this is not for all  $h$ , but only those  $h$ , which also meet the bound on the type1 error because,  $NP$  minimizes the second kind of error while satisfying the bound on the first kind of error.

So, among all classifiers  $h$ , that satisfy the bound on type 1 error. So, if  $h$  is such that probability  $h(X) = 1$  given  $X$  belongs to  $\mathbf{c-0}$  is less than or equal to  $\alpha$ , that means, this classifier  $h$  also satisfies the type 1 error bound. Then the probability of the type 2 error by  $h_{NP}$  is less than or equal to probability of type 2 error by  $h$  right, I hope it is clear. So, the Neyman Pearson classifier is characterized by firstly, it is type 1 error is bounded by  $\alpha$  and it is type 2 error, that is wrongly classifying class 1 pattern is less than the type 2 error of any other classifier  $h$ , if  $h$  also satisfies the bound on the type 1 error.

(Refer Slide Time: 30:18)

**Neyman-Person Classifier**

- Let the bound on Type-I error be  $\alpha$ . Then

$$h_{NP}(\mathbf{X}) = \begin{cases} 1 & \text{if } \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} > K \\ 0 & \text{Otherwise} \end{cases}$$

where  $K$  is such that

$$P \left[ \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} \leq K \mid \mathbf{X} \in \mathbf{c}_0 \right] = 1 - \alpha$$

(We assume  $P\{\mathbf{X} : f_1(\mathbf{X}) = K f_0(\mathbf{X})\} = 0$ , for simplicity)

NPTEL course - p15128

Now this is what the Neyman Pearson classifier is let us suppose the bound is alpha as, we are saying then the h N P classifier is the Neyman Pearson classifier is defined by it will assign class 1 to X. If  $f_1(\mathbf{X}) / f_0(\mathbf{X})$  is greater than K, otherwise assign class 0 where,  $f_1$  is the class conditional density for 1 and  $f_0$  is class conditional density for class 0 where, the K itself is obtained by under the distribution of class 0, that is the another distribution  $f_0$ . The probability that, this ratio is less or equal to K is bound above by is equal to 1 minus alpha.

So, I choose a K see Bayes is some probability involving the random variable X under the distribution  $f_0$ , you think of this  $f_1$  and  $f_0$ , simply as some functions. So,  $f_1$  by  $f_0$  is some other function g of X. So, this is probability g of X less than or equal to K. Under the condition that X is distributed as  $f_0$ , X belongs to  $f_0$  means X is distributed as  $f_0$  right. What does this ensure, this ensures that my type 0 error is equal to alpha, when will I wrongly classify a 0 pattern. I will call class 1, if this ratio is greater than K, when X belongs to  $\mathbf{c}_0$ . This ratio will be greater than K the probability of the ratio is greater than K is equal to alpha because, K is chosen to satisfy this equation, the ratio is less than or equal to K is 1 minus alpha.

So, the probability ratio is greater than K is equal to alpha right. So, by construction the Neyman Pearson classifier satisfies the bound on type 1 error ok. We will we will see the see it once more, just for completeness is sake the way, we stated this, we are assuming



that  $f_1$  and  $f_0$  are true density functions. That is probability  $X$  belonging to any lower dimensional subspace, for example, a sub space the characters be  $f_1(X)$  equals to  $K f_0(X)$  this is 0. So, because, we are assuming that the ratio is either greater than  $K$  or less than equal to  $K$ .

So we we we will we will not allow any any kind of derived delta part in the  $f_1$  and  $f_0$ , this is only a technical condition, those of you do not understand this, do not worry about it simply assume that  $f_1$  and  $f_0$  are nice smooth density functions then this will be all right. So, now let us prove that the classifiers that, we put here right. This is this is the specification this is the this is how Neyman pearson classifier will classify a new pattern  $X$ , this classifier is actually the Neyman pearson classifier that is it satisfies, the 2 condition that, we shut down for then a Neyman pearson classifier.

(Refer Slide Time: 33:17)

• We now prove that this satisfies the NP Criterion. By construction, we have

$$P[h_{NP}(X) = 1 | X \in \mathbf{c-0}] = P\left[\frac{f_1(X)}{f_0(X)} > K | X \in \mathbf{c-0}\right] = \alpha$$

What are the 2 conditions, first is type 1 error should be less than or equal to alpha, for the type 1 error  $h_{NP}(X)$  equal to 1 given  $X$  belongs to  $\mathbf{c-0}$ .

(Refer Slide Time: 33:30)

**Neyman-Person Classifier**

- Let the bound on Type-I error be  $\alpha$ . Then

$$h_{NP}(\mathbf{X}) = \begin{cases} 1 & \text{if } \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} > K \\ 0 & \text{Otherwise} \end{cases}$$

where  $K$  is such that

$$P \left[ \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} \leq K \mid \mathbf{X} \in \mathbf{c}_0 \right] = 1 - \alpha$$

(We assume  $P\{\mathbf{X} : f_1(\mathbf{X}) = Kf_0(\mathbf{X})\} = 0$ , for simplicity)

NPTEL course - p.67128

When will  $h_{NP}(\mathbf{X})$  be 1 by definition  $h_{NP}(\mathbf{X})$  is 1, if  $f_1$  by  $f_0$  is greater than  $1/K$ . So, probability  $h_{NP}(\mathbf{X})$  equal to 1 is probability  $f_1$  by  $f_0$  is greater than  $K$  and when  $\mathbf{X}$  belongs to  $\mathbf{c}_0$ , this where the  $K$  in the Neyman Pearson classifier is obtained by this equation. So, this equation ensures that probability  $f_1$  by  $f_0$  greater than  $K$  is equal to  $\alpha$  right.

(Refer Slide Time: 33:55)

- We now prove that this satisfies the NP Criterion. By construction, we have

$$P[h_{NP}(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c}_0] = P \left[ \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} > K \mid \mathbf{X} \in \mathbf{c}_0 \right] = \alpha$$

- So, we need to show that its Type-II error is less than that for any other classifier satisfying the constraint on Type-I error.

NPTEL course - p.67128

So, the probability of type 0, type 1 error by Neyman Pearson classifier is alpha. So, it satisfies the first criterion. So, now we have to show there is type 2 error is less than that for any other classifier, which also satisfies the constraint on the type 1 error.

(Refer Slide Time: 34:16)

• Let  $h$  be any classifier such that

$$P[h(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c}\text{-}0] \leq \alpha$$

• To complete the proof we have to show that

$$P[h_{NP}(\mathbf{X}) = 0 \mid \mathbf{X} \in \mathbf{c}\text{-}1] \leq P[h(\mathbf{X}) = 0 \mid \mathbf{X} \in \mathbf{c}\text{-}1]$$

Or, equivalently

$$P[h_{NP}(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c}\text{-}1] \geq [P[h(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c}\text{-}1]$$

NPTEL logo and course information are visible at the bottom of the slide.

So let us prove this, so let  $h$  be any other classifier, which also satisfies the constraint on the type 1 error that is probability  $h(\mathbf{X}) = 1$ , when  $\mathbf{X}$  belongs to  $\mathbf{c}\text{-}0$  is less than or equal to alpha. Then to complete the proof you have to show that, the type 2 error of  $h_{NP}$ , that is probability  $h_{NP}(\mathbf{X}) = 0$ , when  $\mathbf{X}$  belongs to  $\mathbf{c}\text{-}1$  is less than or equal to probability of  $h(\mathbf{X}) = 0$ , for when  $\mathbf{x}$  belongs to  $\mathbf{c}\text{-}1$ .

So, this is what we will next show, we actually won't show it in this form the way, we will show it is, we will show the complement of this event. So, we will show that the probability  $h_{NP}(\mathbf{X}) = 1$ , when  $\mathbf{X}$  belongs to  $\mathbf{c}\text{-}1$  is greater than or equal to probability of  $h(\mathbf{X}) = 1$ , where  $\mathbf{X}$  belongs to  $\mathbf{c}\text{-}1$ . So, instead of showing probability  $h_{NP}(\mathbf{X}) = 0$  is less than probability  $h(\mathbf{X}) = 0$ , when  $\mathbf{X}$  belongs to  $\mathbf{c}\text{-}1$  instead of showing this is less than this, we are showing that the complement event  $h_{NP}(\mathbf{X}) = 1$  is greater than probability  $h(\mathbf{X}) = 1$ .

(Refer Slide Time: 35:18)

• Consider the Integral

$$I = \int_{\mathbb{R}^n} (h_{NP}(\mathbf{x}) - h(\mathbf{x})) (f_1(\mathbf{x}) - K f_0(\mathbf{x})) d\mathbf{x}$$

$$= \int_{f_1 > K f_0} (h_{NP}(\mathbf{x}) - h(\mathbf{x})) (f_1(\mathbf{x}) - K f_0(\mathbf{x})) d\mathbf{x} +$$

$$\int_{f_1 \leq K f_0} (h_{NP}(\mathbf{x}) - h(\mathbf{x})) (f_1(\mathbf{x}) - K f_0(\mathbf{x})) d\mathbf{x}$$

• We first show that this integral is always non-negative.

To show this let us consider the integral  $I$ , which is the integral over the entire feature space of the product of 2 terms, the first term is  $h_{NP}(\mathbf{x}) - h(\mathbf{x})$ , second term is  $f_1(\mathbf{x}) - K f_0(\mathbf{x})$ , where  $K$  is the threshold used in the  $n$ -P classifier. Recall that  $h_{NP}$  and  $h$  in this case as binary valued functions  $h_{NP}$  takes value 0 or 1,  $h$  also takes value 0 or 1. So,  $h_{NP} - h$  is some real numbers as, a matter of fact it can be only either minus 1 or plus 1, what we are going to show first either this integral will always be positive and then, we show that, that completes the proof of the  $h_{NP}$  classifier that, we gave is the Neyman Pearson classifier. Let us first know that  $I$  can split this integral into 2 parts, integral over all  $\mathbf{x}$ . So, that  $f_1(\mathbf{x}) > K f_0(\mathbf{x})$  and integral over all  $\mathbf{x}$ . So, that  $f_1(\mathbf{x}) \leq K f_0(\mathbf{x})$  this splits  $\mathbb{R}^n$  into 2 parts and we are going to show that, for each half the integral is positive. Positive means greater than or equal to 0, but this is non negative.

(Refer Slide Time: 36:24)

• When  $f_1(\mathbf{x}) > K f_0(\mathbf{x})$ , we have  
$$h_{NP}(\mathbf{x}) - h(\mathbf{x}) = 1 - h(\mathbf{x}) \geq 0 \quad \text{which implies}$$
$$(h_{NP}(\mathbf{x}) - h(\mathbf{x}))(f_1(\mathbf{x}) - K f_0(\mathbf{x})) \geq 0$$

• Similarly, when  $f_1(\mathbf{x}) < K f_0(\mathbf{x})$ , we have  
$$h_{NP}(\mathbf{x}) - h(\mathbf{x}) = 0 - h(\mathbf{x}) \leq 0 \quad \text{which implies}$$
$$(h_{NP}(\mathbf{x}) - h(\mathbf{x}))(f_1(\mathbf{x}) - K f_0(\mathbf{x})) \geq 0$$

• This shows that  $I \geq 0$ .

So, first let us consider all  $X$ . So, that  $f_1(X)$  is greater than  $K f_0(X)$ , recall that  $h_{NP}(X)$  says  $h_{NP}(X)$  equal to 1, if  $f_1(X)$  is greater than  $K f_0(X)$ . So, for all  $X$  as that  $f_1(X)$  greater than  $K f_0(X)$ , we have  $h_{NP}(X)$  is 1, which means  $h_{NP}(X) - h(X)$  will be  $1 - h(X)$  is always greater than or equal to 0, because  $h(X)$  can be either 1 or 0 right, no matter what classifier  $h$  is  $h(X)$  is either 1 or 0.

So, for all  $X$  either  $f_1(X) - K f_0(X)$ ,  $f_1(X)$  greater than  $K f_0(X)$ ,  $h_{NP}(X) - h(X)$  will always be greater than or equal to 0, which means  $h_{NP}(X) - h(X)$  into  $f_1(X) - K f_0(X)$  is positive, because both terms here are positive. Now, let us look at the other way around let us look at all  $X$  has the  $f_1(X)$  is less than  $K f_0(X)$ , now  $h_{NP}(X)$  will say 0 right, the those  $x$  are put in class 0 by Neyman Pearson classifier. So,  $h_{NP}(X)$  is 0. So,  $h_{NP}(X) - h(X)$  will be  $0 - h(X)$ , which for any classifier  $h(X)$  is less than or equal to 0, because  $h(X)$  can be either 0 or 1. So, once again the product  $h_{NP}(X) - h(X)$  into  $f_1(X) - K f_0(X)$  is positive, because both factors here are negative, which means the integral, we started with is always positive.

(Refer Slide Time: 37:43)

• Thus, we have

$$\int_{\mathbb{R}^n} (h_{NP}(\mathbf{x}) - h(\mathbf{x}))(f_1(\mathbf{x}) - K f_0(\mathbf{x})) d\mathbf{x} \geq 0$$

• This implies

$$\int h_{NP}(\mathbf{x})f_1(\mathbf{x}) d\mathbf{x} - \int h(\mathbf{x})f_1(\mathbf{x}) d\mathbf{x} \geq K \left[ \int h_{NP}(\mathbf{x})f_0(\mathbf{x}) d\mathbf{x} - \int h(\mathbf{x})f_0(\mathbf{x}) d\mathbf{x} \right]$$

So what we have shown, so far is that this integral is positive. So, now let us expand this integral by multiplying these 2 terms. So, multiply with with  $f_1$  first. So, I get  $h_{NP}$  into  $f_1$  minus  $h$  into  $f_1$  right, the those are the first 2 integrals, take the other terms on the other side. So, this is greater than or equal to  $K$  times,  $h_{NP}$  into  $f_0$  minus  $h$  into  $f_0$ , because these, because I have taken them on the other side, now  $h_{NP}$  term will become positive right. I just multiply this term and put 2 integrals on this side 2, integrals on this side.

(Refer Slide Time: 38:29)

Since  $h_{NP}$  and  $h$  take values in  $\{0, 1\}$ ,

$$\int_{\mathbb{R}^n} h_{NP}(\mathbf{x})f_1(\mathbf{X})d\mathbf{X} = P[h_{NP}(\mathbf{X}) = 1 | \mathbf{X} \in \mathbf{c-1}]$$

and

$$\int_{\mathbb{R}^n} h(\mathbf{x})f_1(\mathbf{X})d\mathbf{X} = P[h(\mathbf{X}) = 1 | \mathbf{X} \in \mathbf{c-1}]$$

Similarly for the integrals involving  $f_0$ .

Since  $h_{NP}$  and  $h$  are binary valued functions for any given  $X$   $h_{NP}(X)$  is either 0 or 1 similarly,  $h$  is a 0 or 1. So, if I take an integral over entire feature space  $\mathbb{R}^n$   $h_{NP}$  into  $f_1$  over  $X$ , it is simply integral of  $f_1$  over the region, over the set of all  $X$ , so that  $h_{NP}(X)$  is 1 right. Because when integrating over 1, this integral is nothing but integrating with  $f_1$ , here the integral is nothing but, conditioned on  $X$  belongs to  $C_1$  probability that  $h_{NP}$  is 1 right.

Because  $h_{NP}$  is a 0, 1 valued function integral of  $h_{NP}$  into  $f_1$  over  $\mathbb{R}^n$  is same as integral  $f_1$  over  $x$  at that  $h_{NP}(X)$  is 1, which is same as because, I am integrating  $f_1$ , which is same as probability  $h_{NP}(X)$  is equal to 1, conditioned on  $X$  belongs to  $C_1$ . Similarly, for  $h$  because  $h$  is also 0, 1 function integral  $h$  into  $f_1$  is probability that  $h(X)$  is 1 conditioned  $X$  belongs to  $C_1$ . similarly, the integral with respect to  $f_0$ .

(Refer Slide Time: 39:37)

• Thus, we have

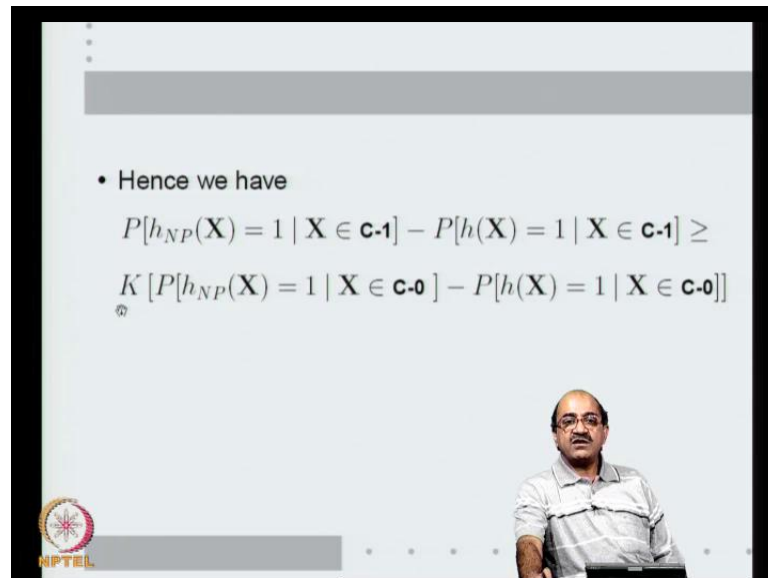
$$\int_{\mathbb{R}^n} (h_{NP}(\mathbf{x}) - h(\mathbf{x}))(f_1(\mathbf{x}) - K f_0(\mathbf{x})) d\mathbf{x} \geq 0$$

• This implies

$$\int h_{NP}(\mathbf{x}) f_1(\mathbf{x}) d\mathbf{x} - \int h(\mathbf{x}) f_1(\mathbf{x}) d\mathbf{x} \geq K \left[ \int h_{NP}(\mathbf{x}) f_0(\mathbf{x}) d\mathbf{x} - \int h(\mathbf{x}) f_0(\mathbf{x}) d\mathbf{x} \right]$$

So, which means each of these integrals this can be retained as, probability  $h_{NP}(X)$  is equal to 1 conditioned on  $X$  belongs to  $C_1$ . This is  $h(X)$  is equal to 1 conditioned on  $X$  belongs to  $C_1$ . Similarly, this is  $h_{NP}(X)$  is 1 conditioned on  $X$  belongs to  $C_0$ ,  $h(X)$  is 1 conditioned on  $X$  belongs to  $C_0$  right.

(Refer Slide Time: 40:20)



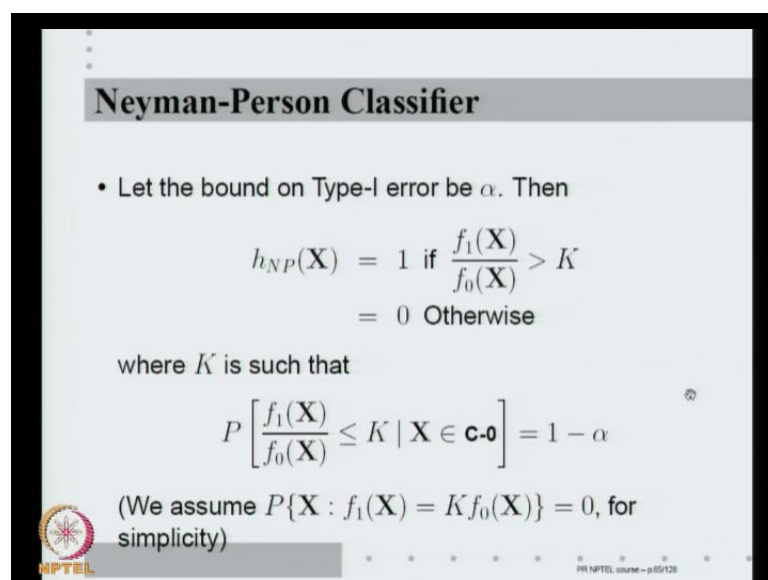
• Hence we have

$$P[h_{NP}(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c-1}] - P[h(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c-1}] \geq K [P[h_{NP}(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c-0}] - P[h(\mathbf{X}) = 1 \mid \mathbf{X} \in \mathbf{c-0}]]$$

The slide features a presenter in the bottom right corner and an NPTEL logo in the bottom left corner.

So this inequality can now be written as probability  $h_{NP}(\mathbf{X})$  is equal to 1 given  $\mathbf{X}$  belongs to  $\mathbf{c-1}$ , minus probability  $h(\mathbf{X})$  is equal to 1 given  $\mathbf{X}$  belongs to  $\mathbf{c-1}$  is greater than or equal to  $K$  times. Probability  $h_{NP}(\mathbf{X})$  is equal to 1 conditioned on  $\mathbf{x}$  belongs to  $\mathbf{c-0}$  minus probability  $h(\mathbf{X})$  is equal to 1, conditioned on  $\mathbf{x}$  belongs to  $\mathbf{c-0}$ , let us also remember that this factor  $K$  will always be positive right.

(Refer Slide Time: 40:52)



### Neyman-Person Classifier

• Let the bound on Type-I error be  $\alpha$ . Then

$$h_{NP}(\mathbf{X}) = \begin{cases} 1 & \text{if } \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} > K \\ 0 & \text{Otherwise} \end{cases}$$

where  $K$  is such that

$$P \left[ \frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} \leq K \mid \mathbf{X} \in \mathbf{c-0} \right] = 1 - \alpha$$

(We assume  $P\{\mathbf{X} : f_1(\mathbf{X}) = K f_0(\mathbf{X})\} = 0$ , for simplicity)

The slide includes an NPTEL logo in the bottom left corner and a small footer text '© NPTEL course - p 15/128' in the bottom right corner.

Because, two ways of looking at it,  $K$  is defined by, the reference of Neyman Pearson classifier. This is the reference of Neyman Pearson classifiers, both  $f_1$  and  $f_0$  are



density functions, they are always positive. So, this ratio is always positive, so, if say  $K$  is negative is forever satisfied, because so, that is the 1 way looking at it, any case because, this is some positive function of and, we want it less than or equal to  $K$  has to have some positive probability  $K$  has to be a positive number right.

(Refer Slide Time: 41:30)

• Hence we have

$$P[h_{NP}(X) = 1 | X \in \mathbf{c-1}] - P[h(X) = 1 | X \in \mathbf{c-1}] \geq K [P[h_{NP}(X) = 1 | X \in \mathbf{c-0}] - P[h(X) = 1 | X \in \mathbf{c-0}]]$$

• But for all  $h$  under consideration, the RHS above is non-negative. Hence

$$P[h_{NP}(X) = 1 | X \in \mathbf{c-1}] - P[h(X) = 1 | X \in \mathbf{c-1}] \geq 0$$

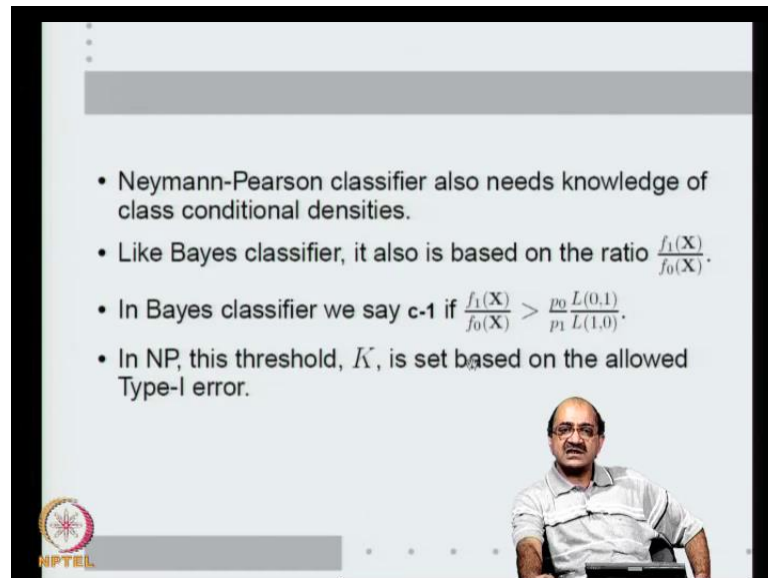
• This completes the proof.

NPTEL course - 94428

Given that case a positive number, let us look at what is there in the in the brackets here, this is  $h_{NP}(X) = 1$ , conditioned  $X$  belongs to  $\mathbf{c-0}$  minus  $h(X) = 1$  conditioned as  $X$  belongs to  $\mathbf{c-0}$ . Now this is the type 1 error of  $h_{NP}(X)$  of the Neyman Pearson classifier, which by construction is  $\alpha$ , this is the type 1 error of the classifier  $h$  and because,  $h$  is something that satisfies the conditions on type 1 error this is less than or equal to  $\alpha$ . So, this factor is greater than or equal to 0 right.

So, for all  $h$  under consideration all  $h$ , that satisfy the constraints on type 1 error right, the term on the R H S is always non negative, which means this is positive  $h_{NP}(X) = 1$  conditioned  $X$  belongs to  $\mathbf{c-1}$  minus  $h(X) = 1$  conditioned on  $X$  belongs to  $\mathbf{c-1}$  is greater than or equal to 0. This shows that the Neyman Pearson classifier has the smallest type 2 error compared to smaller type 2 error compared to any classifier, that also satisfies the type 1 error bound right. This shows that the classifiers that, we have actually put down as Neyman Pearson classifiers, satisfies the criteria for Neyman Pearson classifier.

(Refer Slide Time: 42:51)

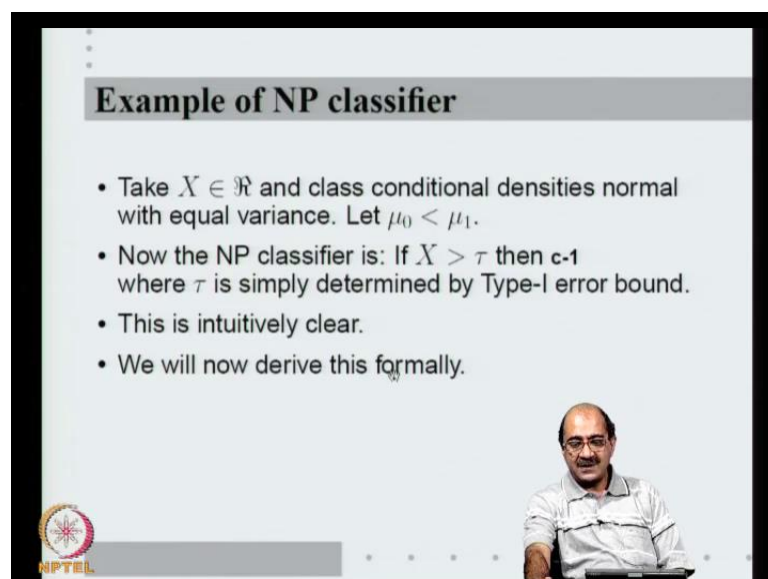


A slide from a presentation with a black border. At the top, there is a grey header bar. The main content area is white and contains a bulleted list of four points. The first point states that the Neymann-Pearson classifier needs knowledge of class conditional densities. The second point compares it to the Bayes classifier, noting it is based on the ratio  $\frac{f_1(\mathbf{X})}{f_0(\mathbf{X})}$ . The third point defines the decision rule for the Bayes classifier:  $c-1$  if  $\frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} > \frac{p_0 L(0,1)}{p_1 L(1,0)}$ . The fourth point states that in the NP classifier, the threshold  $K$  is set based on the allowed Type-I error. In the bottom right corner, there is a small inset image of a man in a white shirt sitting at a desk. In the bottom left corner, there is a circular logo with a star and the text 'NPTEL' below it.

- Neymann-Pearson classifier also needs knowledge of class conditional densities.
- Like Bayes classifier, it also is based on the ratio  $\frac{f_1(\mathbf{X})}{f_0(\mathbf{X})}$ .
- In Bayes classifier we say c-1 if  $\frac{f_1(\mathbf{X})}{f_0(\mathbf{X})} > \frac{p_0 L(0,1)}{p_1 L(1,0)}$ .
- In NP, this threshold,  $K$ , is set based on the allowed Type-I error.

So, Neyman Pearson, classifiers also needs a knowledge of class conditional densities because, you have to calculate  $f_0$  by  $f_1$ . Like Bayes classifier, it is also based on the ratio  $f_1$  by  $f_0$ , in Bayes classifier, we say  $c-1$ , if  $f_1$  by  $f_0$  is greater than some threshold, which is which happens to be  $p_0, l(0,1)$  by  $p_1, l(1,0)$ . In n p this is some other threshold  $K$ , which is set based on the allowed type 1 error, so both of them essentially threshold the ratio of the 2 class conditional densities.

(Refer Slide Time: 43:20)



A slide from a presentation with a black border. At the top, there is a grey header bar with the text 'Example of NP classifier'. The main content area is white and contains a bulleted list of four points. The first point says to take  $X \in \mathcal{R}$  and class conditional densities normal with equal variance, and let  $\mu_0 < \mu_1$ . The second point says the NP classifier is: if  $X > \tau$  then  $c-1$  where  $\tau$  is simply determined by Type-I error bound. The third point says this is intuitively clear. The fourth point says they will now derive this formally. In the bottom right corner, there is a small inset image of a man in a white shirt sitting at a desk. In the bottom left corner, there is a circular logo with a star and the text 'NPTEL' below it.

### Example of NP classifier

- Take  $X \in \mathcal{R}$  and class conditional densities normal with equal variance. Let  $\mu_0 < \mu_1$ .
- Now the NP classifier is: If  $X > \tau$  then c-1 where  $\tau$  is simply determined by Type-I error bound.
- This is intuitively clear.
- We will now derive this formally.

So, let us quickly look at a simple example, of Neyman pearson classifier, let us take a one dimensional feature space normal class conditional densities with equal variance, and suppose, by now you if you are, if you have been following all the lectures, you know that this is always a simplest case. A 2 class problem, one dimensional feature space, normal class conditional densities equal variance.

For let us assume that  $\mu_0$  is less than  $\mu_1$ ,  $\mu_0$  is the mean of the class 0 and  $\mu_1$  is the mean of class 1. So, what is the N p classifier, if  $X$  greater than  $\tau$  then  $c_1$ , where how where how do, I choose  $\tau$ ,  $\tau$  is simple taken by the type 1 error bound.

(Refer Slide Time: 44:07)

**Example**

Now (assuming  $\mu_1 > \mu_0$ ),

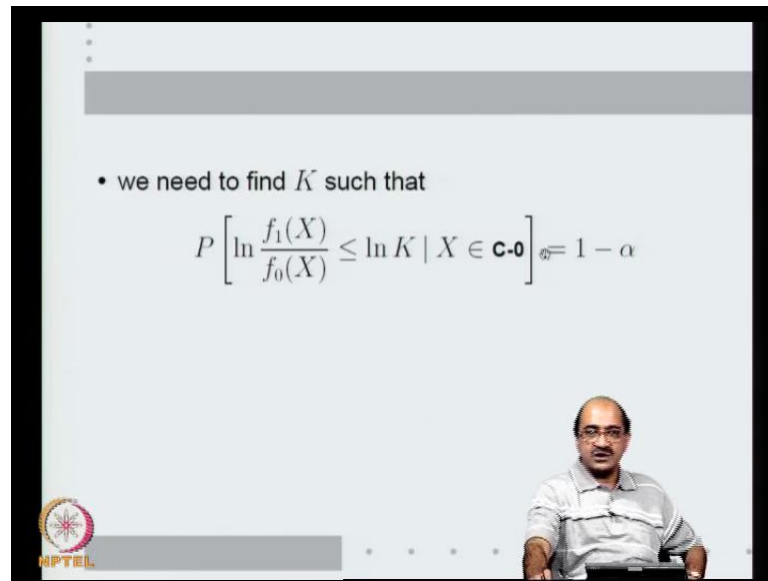
$$\begin{aligned} \frac{f_1(X)}{f_0(X)} &= \exp\left(-\frac{(X - \mu_1)^2}{2\sigma^2} + \frac{(X - \mu_0)^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{1}{2\sigma^2}[\mu_1^2 - \mu_0^2 - 2X(\mu_1 - \mu_0)]\right) \\ &= \exp\left(\frac{\mu_1 - \mu_0}{2\sigma^2}[2X - (\mu_1 + \mu_0)]\right) \end{aligned}$$

NPTEL PR NPTEL course - p.59128

So let us derive this formally, how do how do I get this term, because both  $f_1$  and  $f_0$  are normal, it is easy to see  $f_1$  by  $f_0$  exponential minus  $X$  minus  $\mu_1$ , whole square by  $2\sigma^2$  plus  $X$  minus  $\mu_0$  square by  $2\sigma^2$ , the other factors  $1$  by  $\sigma\sqrt{2\pi}$  will cancel right. Now, we can expand this the  $X$  square term will cancel.

So, what will get minus  $1$  by  $2\sigma^2$   $\mu_1^2$  square from here, sorry  $\mu_1$  square from here minus  $\mu_0$  square from here right and  $2X$  into  $\mu_1 - \mu_0$ , I can absorb this minus sign. So, I can write it as  $\mu_1 - \mu_0$  by  $2\sigma^2$  into  $2X - \mu_1 + \mu_0$ . So, this is the ratio of  $f_1$  by  $f_0$  for the case of normal densities with equal variance.

(Refer Slide Time: 45:00)



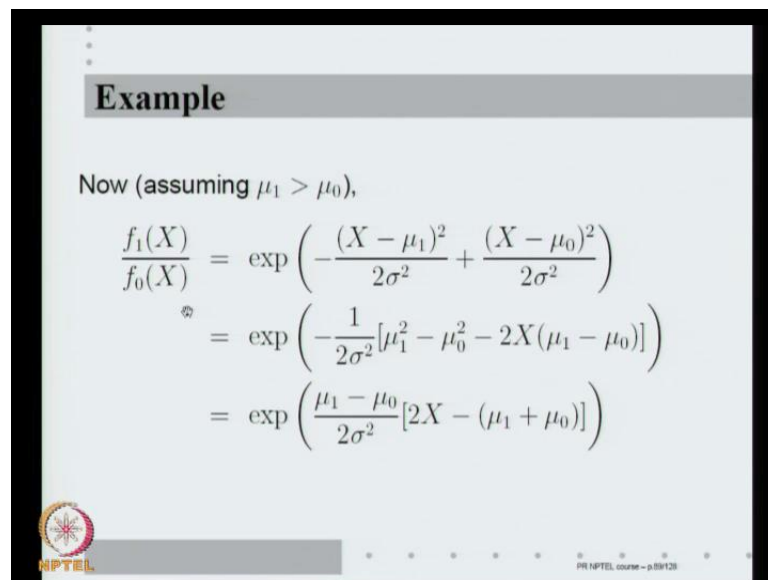
• we need to find  $K$  such that

$$P \left[ \ln \frac{f_1(X)}{f_0(X)} \leq \ln K \mid X \in \mathbf{c-0} \right] = 1 - \alpha$$

The slide also features the NPTEL logo in the bottom left corner and a small inset image of a man in the bottom right corner.

We need to find a  $K$ , such that probability  $f_1$  by  $f_0$  less than or equal to  $K$  is  $1$  minus  $\alpha$ , because  $\log$  is a monotone function, which is same as probability  $\log$  of  $f_1$  by  $f_0$  less than or equal to  $\log$  of conditioned on  $c_0$  that is also good enough.

(Refer Slide Time: 45:20)



**Example**

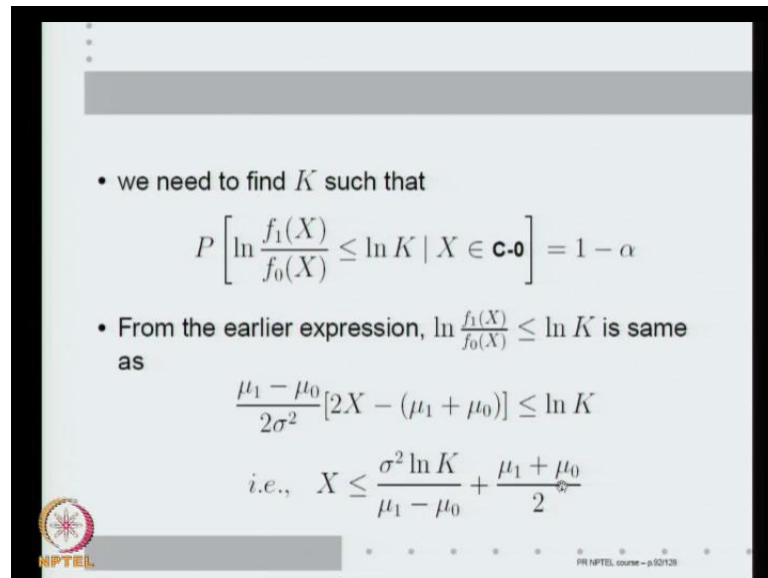
Now (assuming  $\mu_1 > \mu_0$ ),

$$\begin{aligned} \frac{f_1(X)}{f_0(X)} &= \exp \left( -\frac{(X - \mu_1)^2}{2\sigma^2} + \frac{(X - \mu_0)^2}{2\sigma^2} \right) \\ &= \exp \left( -\frac{1}{2\sigma^2} [\mu_1^2 - \mu_0^2 - 2X(\mu_1 - \mu_0)] \right) \\ &= \exp \left( \frac{\mu_1 - \mu_0}{2\sigma^2} [2X - (\mu_1 + \mu_0)] \right) \end{aligned}$$

The slide also features the NPTEL logo in the bottom left corner and a small inset image of a man in the bottom right corner.

Now this is  $f_1$  by  $f_0$ . So,  $\log$  of  $f_1$  by  $f_0$  will be just simply, what is inside the exponent.

(Refer Slide Time: 45:26)



• we need to find  $K$  such that

$$P \left[ \ln \frac{f_1(X)}{f_0(X)} \leq \ln K \mid X \in \mathbf{c-0} \right] = 1 - \alpha$$

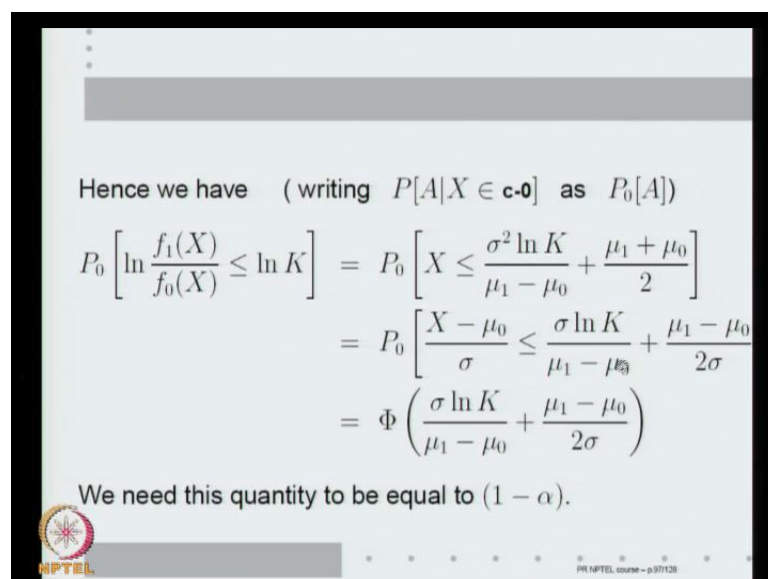
• From the earlier expression,  $\ln \frac{f_1(X)}{f_0(X)} \leq \ln K$  is same as

$$\frac{\mu_1 - \mu_0}{2\sigma^2} [2X - (\mu_1 + \mu_0)] \leq \ln K$$
$$\text{i.e., } X \leq \frac{\sigma^2 \ln K}{\mu_1 - \mu_0} + \frac{\mu_1 + \mu_0}{2}$$

NIPTEL PR/NPTEL course - p.52/28

So, what we get is  $\mu_1 - \mu_0$  by  $2\sigma^2$  into  $2X - \mu_1 + \mu_0$  should be less than or equal to probability  $K$ , so I can use this. So this is some expression involving random variable  $X$  right. I want the probability of this event to be equal to  $1 - \alpha$ , I have to choose  $K$  like that. So we can do it like this, so what will this give me, I can first take  $2\sigma^2$  by  $\mu_1 - \mu_0$  this side right. Then bring  $\mu_1 + \mu_0$  on this side and then, I will get  $2X$  here divided by  $2$ . So, this inequality is same as  $X$  less or equal to  $\sigma^2 \ln K$  by  $\mu_1 - \mu_0$  plus  $\mu_1 + \mu_0$  by  $2$ .

(Refer Slide Time: 46:10)



Hence we have (writing  $P[A \mid X \in \mathbf{c-0}]$  as  $P_0[A]$ )

$$\begin{aligned} P_0 \left[ \ln \frac{f_1(X)}{f_0(X)} \leq \ln K \right] &= P_0 \left[ X \leq \frac{\sigma^2 \ln K}{\mu_1 - \mu_0} + \frac{\mu_1 + \mu_0}{2} \right] \\ &= P_0 \left[ \frac{X - \mu_0}{\sigma} \leq \frac{\sigma \ln K}{\mu_1 - \mu_0} + \frac{\mu_1 - \mu_0}{2\sigma} \right] \\ &= \Phi \left( \frac{\sigma \ln K}{\mu_1 - \mu_0} + \frac{\mu_1 - \mu_0}{2\sigma} \right) \end{aligned}$$

We need this quantity to be equal to  $(1 - \alpha)$ .

NIPTEL PR/NPTEL course - p.57/28

So, what we want is let us say  $P_0$  is the probability conditioned on  $X$  belongs to class 0. So,  $P_0(1 - \alpha) = P(X \leq c_0)$  is same as  $P_0(X \leq c_0)$  of  $X$  less than or equal to that quantity. Now because, this is a  $P_0$  probability means, I am taking this probability under the class 0 class conditional density that is  $\mu_0$  mean and  $\sigma$  variance. So, I can write this as  $X \leq c_0$  same as,  $X - \mu_0 \leq c_0 - \mu_0$  by  $\sigma$  less than or something else.

So, just subtract  $\mu_0$  and divide by  $\sigma$ , I get this expression, why did I do that, because I know the distribution of  $X - \mu_0$  by  $\sigma$ . Because  $X$  is under this probability  $X$  belongs to  $c_0$ ,  $c_0$  is normal with mean  $\mu_0$ , variance  $\sigma^2$   $X - \mu_0$  by  $\sigma$  is standard normal. So, this probability is given in terms of the standard normal function,  $\Phi$  of this quantity  $\frac{c_0 - \mu_0}{\sigma} = \frac{\mu_1 - \mu_0}{\sigma} + \frac{\mu_1 - \mu_0}{2\sigma}$ , where  $\Phi$  is the density of the standard normal. So, ultimately to get  $K$ , I have to equate this quantity to  $1 - \alpha$ .  $\alpha$  is given in Neyman Pearson criteria what I am given is  $\alpha$  that is the allowed type 1 error bound. So, to get  $K$ , I have to equate this to  $1 - \alpha$ .

(Refer Slide Time: 47:33)

Thus we want

$$\Phi\left(\frac{\sigma \ln K}{\mu_1 - \mu_0} + \frac{\mu_1 - \mu_0}{2\sigma}\right) = (1 - \alpha)$$

This gives us an expression for  $\ln K$  as

$$\frac{\sigma \ln K}{\mu_1 - \mu_0} = \Phi^{-1}(1 - \alpha) - \frac{\mu_1 - \mu_0}{2\sigma}$$

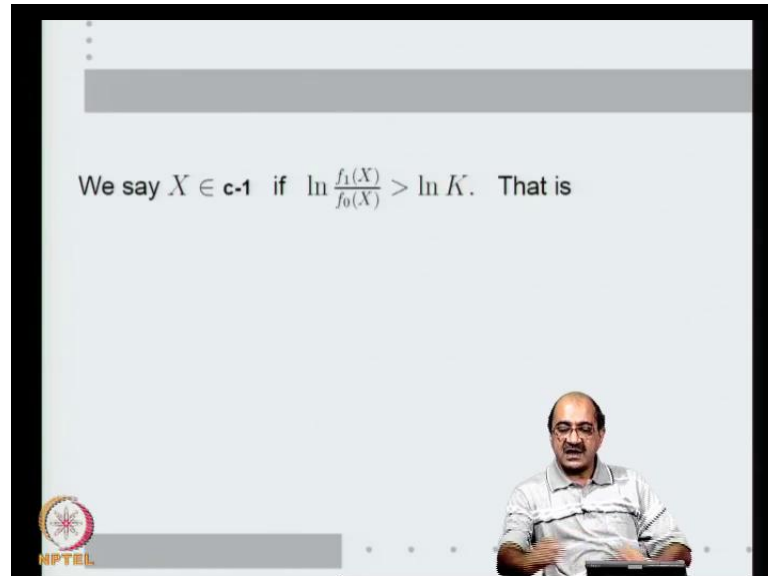
or

$$\ln K = \frac{\mu_1 - \mu_0}{\sigma} \Phi^{-1}(1 - \alpha) - \frac{(\mu_1 - \mu_0)^2}{2\sigma^2}$$

So, let us equate that  $\Phi$  of this is equal to  $1 - \alpha$ , now I can solve this for  $K$  or  $\ln K$  that gives me  $\frac{\sigma \ln K}{\mu_1 - \mu_0} + \frac{\mu_1 - \mu_0}{2\sigma} = \Phi^{-1}(1 - \alpha)$  minus  $\mu_1 - \mu_0$  by  $2\sigma$ . Now multiply by  $\mu_1 - \mu_0$  divided by

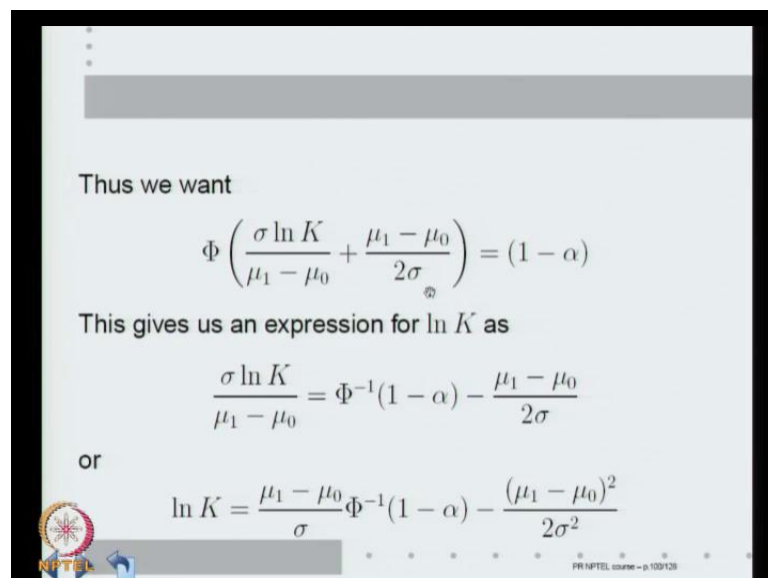
sigma that gives me  $1/NK$  is this much, from this, I can get  $K$ . So this is the threshold, I want for Neyman Pearson classifier.

(Refer Slide Time: 48:06)



Let us ask what does this classifier mean so, in Neyman Pearson classifier, we will put  $X$  in class 1, if  $f_1$  by  $f_0$  is greater than  $K$ , which is same as  $\ln f_1$  by  $f_0$  greater than  $\ln K$ .

(Refer Slide Time: 48:19)



And the  $\ln K$  is given by this expression that, we have just now derived.

(Refer Slide Time: 48:23)

We say  $X \in c-1$  if  $\ln \frac{f_1(X)}{f_0(X)} > \ln K$ . That is

$$\frac{\mu_1 - \mu_0}{2\sigma^2} [2X - (\mu_1 + \mu_0)] > \frac{\mu_1 - \mu_0}{\sigma} \Phi^{-1}(1 - \alpha) - \frac{(\mu_1 - \mu_0)^2}{2\sigma^2}$$

i.e.,  $2X - (\mu_1 + \mu_0) > 2\sigma \Phi^{-1}(1 - \alpha) - (\mu_1 - \mu_0)$

i.e.,  $X > \sigma \Phi^{-1}(1 - \alpha) + \mu_0$

NPTEL logo and course ID PR NPTEL course - p104123 are visible at the bottom.

So, I will put  $X$  in  $c-1$ , if  $\ln \frac{f_1}{f_0}$ , that is this expression is greater than  $\ln K$ , that is this expression. Now I can simplify this expression  $2X$  minus  $\mu_1$  plus  $\mu_0$  is multiplied by  $2\sigma^2$  and divided by  $\mu_1 - \mu_0$ , I get this right.  $2\sigma \Phi^{-1}(1 - \alpha) - (\mu_1 - \mu_0)$ , remember, we are assuming  $\mu_0$  less than  $\mu_1$ . So,  $\mu_1 - \mu_0$  is a positive quantity. So, when I divide by it the inequality does not change. Now, if being  $\mu_1 - \mu_0$  on this side and divided by 2, this is same as  $X$  greater than this.

(Refer Slide Time: 49:09)

Thus NP classifier puts  $X \in c-1$  if

i.e.,  $X > \sigma \Phi^{-1}(1 - \alpha) + \mu_0$

i.e.,  $\Phi\left(\frac{X - \mu_0}{\sigma}\right) > (1 - \alpha)$

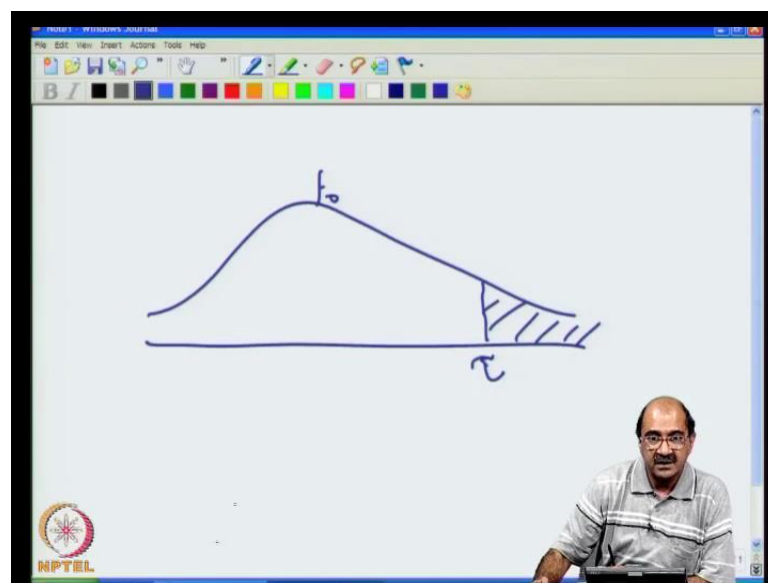
This means the NP classifier puts  $X$  in  $c-1$  if  $X > \tau$   
where  $\int_{\tau}^{\infty} f_0(X) dX = \alpha$ .

NPTEL logo and a video frame of a man are visible at the bottom.



Which, so the Neyman Pearson classifier, we will put  $X$  in class 1, if  $X$  is greater than this. Which is same as  $X$  minus  $\mu_0$  by  $\sigma$  and bring  $\Phi$  inverse this side  $\Phi$  of  $X$  minus  $\mu_0$  by  $\sigma$  is greater than  $1 - \alpha$  right. What is  $\Phi$  of  $X$  minus  $\mu_0$  by  $\sigma$  that is the the distribution of the standard normal density. So, what does this mean, this means this threshold  $X$  greater than  $\tau$ , if I think of this as  $\tau$ , the integral of the density function of class 0, starting from this  $\tau$  to infinity will exactly be equal to  $\alpha$  right. Because  $\Phi$  of this is greater than  $1 - \alpha$ , what is remaining in the integral will be equal to  $\alpha$ .

(Refer Slide Time: 50:14)





So let us let us look at this what it means, so if this is my  $f_0$ . I am choosing a  $\tau$  as my threshold, such that this area. Area from  $\tau$  to infinity is equal to  $\alpha$  that is the allowed type 1 error.

(Refer Slide Time: 50:41)

Thus NP classifier puts  $X \in c-1$  if

$$i.e., \quad X > \sigma \Phi^{-1}(1 - \alpha) + \mu_0$$
$$i.e., \quad \Phi\left(\frac{X - \mu_0}{\sigma}\right) > (1 - \alpha)$$



This means the NP classifier puts  $X$  in  $c-1$  if  $X > \tau$   
where  $\int_{\tau}^{\infty} f_0(X) dX = \alpha$ .

So, the Neyman Pearson classifier ultimately that that ratio being some greater than  $K$ , turns out to be same as,  $X$  greater than  $\tau$ , where  $\tau$  is chosen, so that  $\tau$  to infinity  $\int_{\tau}^{\infty} f_0(X) dX = \alpha$ , this is what, we want. Because, in in that normal and equal variance ultimately, the classifier is a threshold and type 1 error, because we are assuming  $\mu_0$  is less than  $\mu_1$ , type 1 error is simply integral from  $\tau$  to infinity of the class 0 density function.

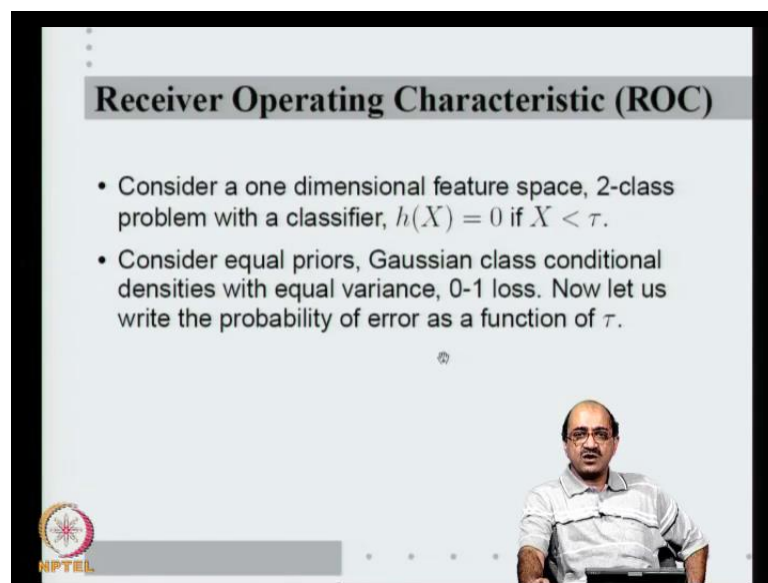
(Refer Slide Time: 51:26)

- Like the Bayes classifier, the NP classifier also needs knowledge of class conditional densities.
- NP classifier is only for the 2-class case.
- It is actually more important in hypothesis testing problems. (Likelihood ratio test)

So, like the Bayes classifier, the Neyman Pearson classifier also needs knowledge of class conditional densities and this classifier is only for 2 class case, we did not define it for multiclass case and is often difficult to define it, for extensive multiclass case. Just as general information, it is more important in certain statistics problems called hypothesis testing problems. More than classifier though for 2 class case is also used, especially when, you do not want to trade 1 kind of error with another rather than, that you want to put a bound on one kind of error and given that, bound is satisfied minimize the other kind of error.

(Refer Slide Time: 52:04)

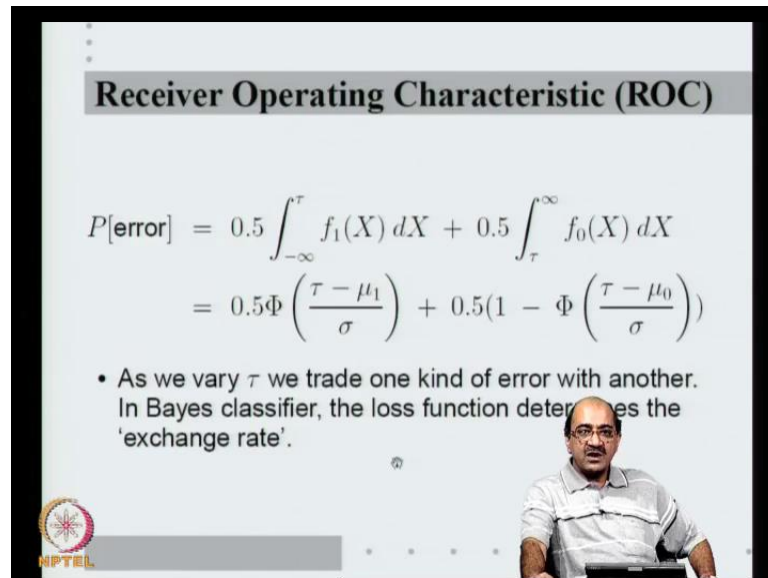


**Receiver Operating Characteristic (ROC)**

- Consider a one dimensional feature space, 2-class problem with a classifier,  $h(X) = 0$  if  $X < \tau$ .
- Consider equal priors, Gaussian class conditional densities with equal variance, 0-1 loss. Now let us write the probability of error as a function of  $\tau$ .

Actually, as we have seen in the Neyman Pearson classifier is for trading one kind of error with another kind of error. So this kind of thing happens in many other ways a good way of looking at it is what is called the receiver operating characteristic. If you consider a one dimensional feature space, 2 class problem with  $h(x)$  equal to 0 with a particular threshold  $\tau$ . Now, if I think of class conditional densities once again as normal with equal variance 0 1 loss function.

(Refer Slide Time: 52:47)



**Receiver Operating Characteristic (ROC)**

$$P[\text{error}] = 0.5 \int_{-\infty}^{\tau} f_1(X) dX + 0.5 \int_{\tau}^{\infty} f_0(X) dX$$
$$= 0.5 \Phi\left(\frac{\tau - \mu_1}{\sigma}\right) + 0.5(1 - \Phi\left(\frac{\tau - \mu_0}{\sigma}\right))$$

- As we vary  $\tau$  we trade one kind of error with another. In Bayes classifier, the loss function determines the 'exchange rate'.

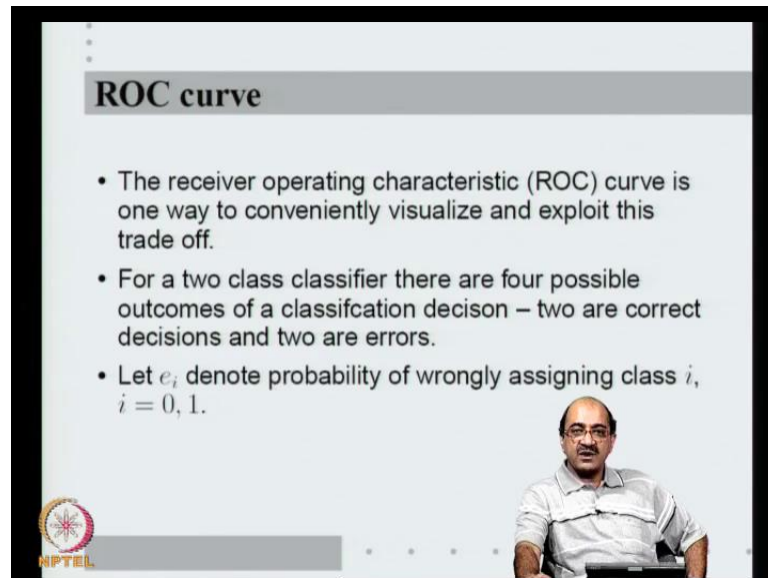
The slide also features a small logo in the bottom left corner and a photograph of a man in a light-colored shirt in the bottom right corner.

We can write the error as a function of tau, I am below tau, I am putting class 1 above tau, I am putting class 2. So, this is my error integral as, we have already written earlier right. So, my error integral is point, because priors are equal 0.5 times probability that class 1 pattern comes below tau and 0.5 times probability that class 0 pattern comes above tau right.

Once again I can if because,  $f_1$  and  $f_0$  are normal. I can write the standard normal as, we vary tau essentially, we are trading 1 kind of error with another, when I change tau 1 kind of error may increase and other kind of error will decrease. So, varying tau allows us to trade 1 error with another and hence atleast a threshold Bayes classifier, we can actually sit and decide, how we want to do the trade of right, risk function is 1 way of doing the trade of where the loss function gives me.



So, to say the exchange rate how much of 1 kind of error, I can trade for how much of other kind of error that these the relative values of the 2 losses. Neyman pearson criteria gives me another way of trading these errors right. I want this error below some alpha and and minimize the other error, but I can choose my own trade of right by using what is called a receiver operating characteristic.

(Refer Slide Time: 54:04)



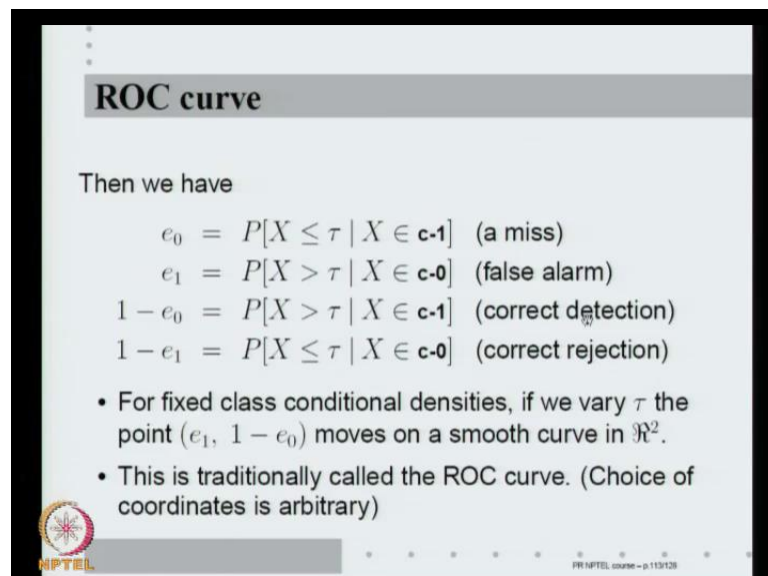
**ROC curve**

- The receiver operating characteristic (ROC) curve is one way to conveniently visualize and exploit this trade off.
- For a two class classifier there are four possible outcomes of a classification decision – two are correct decisions and two are errors.
- Let  $e_i$  denote probability of wrongly assigning class  $i$ ,  $i = 0, 1$ .

The receiver operating characteristic curve is 1 way of conveniently visualizing, this trade off for a 2 class classifier. There are 4 possible errors right, 2 are correct 4 possible outcomes of a classification decision. 2 are correct and the other 2 are wrong, let  $e_i$  denote the probability of wrongly assigning class  $i$ , that is  $e_0$  is wrongly assigning class 0, at calling 0, when it is actually 1, let us say  $e_1$  is wrongly assigning class 1, which means calling 1, when it is actually class 0.

(Refer Slide Time: 54:38)


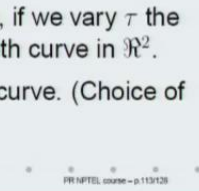


**ROC curve**

Then we have

$$e_0 = P[X \leq \tau | X \in \mathbf{c-1}] \quad (\text{a miss})$$
$$e_1 = P[X > \tau | X \in \mathbf{c-0}] \quad (\text{false alarm})$$
$$1 - e_0 = P[X > \tau | X \in \mathbf{c-1}] \quad (\text{correct detection})$$
$$1 - e_1 = P[X \leq \tau | X \in \mathbf{c-0}] \quad (\text{correct rejection})$$

- For fixed class conditional densities, if we vary  $\tau$  the point  $(e_1, 1 - e_0)$  moves on a smooth curve in  $\mathcal{R}^2$ .
- This is traditionally called the ROC curve. (Choice of coordinates is arbitrary)

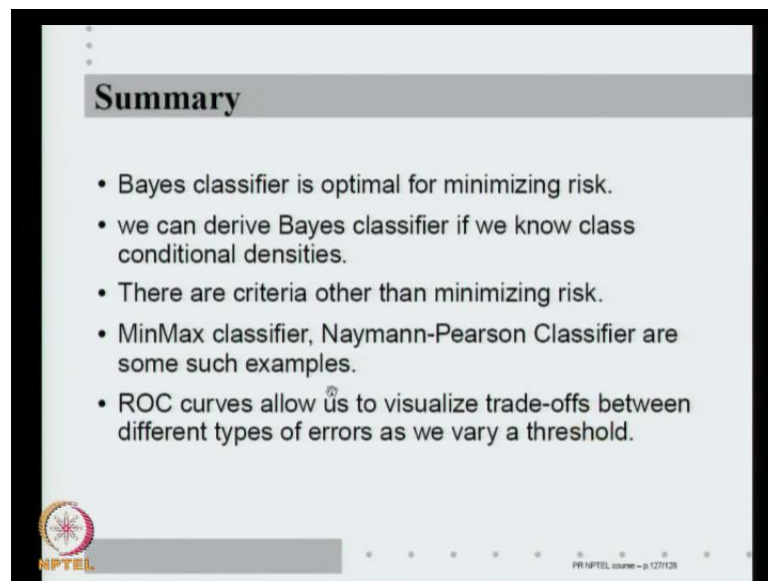
 

PR NPTEL course - p.113/128

So, actually I can define these like this  $e_0$  is the probability that, when  $X$  is actually in class 1,  $X$  is also less than  $\tau$ . So, I call class 0, it is called a misse,  $e_1$  is wrongly assigning class 1, that is even though,  $X$  is actually in for comes from  $c_0$ , because  $X$  is greater than  $\tau$ . I call class 1 it is called a false alarm right, the  $1 - e_0$  and  $1 - e_1$  probabilities are called the correct detection and correct rejection.

For fixed class densities, If we vary  $\tau$ , the point  $e_1$  comma  $1 - e_0$  moves on a smooth curve in  $r^2$ , ofcourse, I could have chosen any 2 of these numbers the choice of co ordinates is arbitrary. But, this curve, which plots the false alarm rate verses the the correct detection rate that is  $e_1$  verses  $1 - e_0$  right. For various values of  $\tau$  right, that is on the  $e_1$   $1 - e_0$  plane, for each  $\tau$  there is 1 value of  $e_1$  1 value of  $1 - e_0$ . So, as I vary  $\tau$  it becomes a smooth curve and such a curve is called the receiver operating characteristic. This is another way of trading 1 kind of error with the other kind of error ok.

(Refer Slide Time: 56:04)



So, let us close down today is class with a summary, the Bayes classifier is optimal for minimizing risk, we have seen that in last class and we have seen a couple of more examples, this class. We can derive Bayes classifier, if we know class conditional densities and for various kinds of loss functions, we can derive. There are criteria other than minimizing risk as we have seen minimizing risk is not the only way, we can we can run this problem.

For example, minmax classifier, Neymann pearson classifier are some examples of criteria other than minimizing risk. All of these are essentially trying to trade of errors in a way different from the trade of that to the Bayes classifiers does, Bayes classifier trades of errors, using the loss function of the exchange rate. Where as, there are other ways of trading of errors and receive a operating curve characteristic curves allow us to visualize the trade of between different types of errors, as we vary a threshold. We will once again briefly look at the receiver operating characteristics next class.

Thank you.