

**Statistical Methods for Scientists and Engineers**  
**Prof. Somesh Kumar**  
**Department of Mathematics**  
**Indian Institute of Technology – Kharagpur**

**Lecture – 26**  
**Multivariate Analysis – XI**

In the last few lectures, I have introduced the problem of classifying an observation into one of the 2 populations. I discussed various procedures. In particular, I showed that one can define Bayesian classification rules or what you can say as good classification rules, so we call them the admissible rules or minimal complete class, that means the rules beyond which you need not discuss. In particular, we consider applications to the classification for an observation into 2 multivariate normal populations.

The first case was when all the parameters are known and then second case we discussed when the parameters are unknown. Now, I will generalize this concept to the problem of classification of one observation into several populations. So, let me introduce the concept of optimal rules here and how to derive these rules.

**(Refer Slide Time: 01:20)**

Lecture - 26

Classifying an Observation into one of several populations  
 $\pi_1, \dots, \pi_m \rightarrow m$  populations  
with associated density functions  $p_1(x), \dots, p_m(x)$  respectively.  
We want to specify mutually exclusive and exhaustive regions of  
the sample space  $R_1, \dots, R_m$ .  
If  $x \in R_i$  we classify  $x$  into  $\pi_i$ ,  $i=1, \dots, m$ .  
The cost of misclassifying an observation from  $\pi_i$  into  $\pi_j$  as  
 $C(j|i)$ .  
 $C_R(j|i) = \int_{R_j} p_i(x) dx \quad \dots (1)$

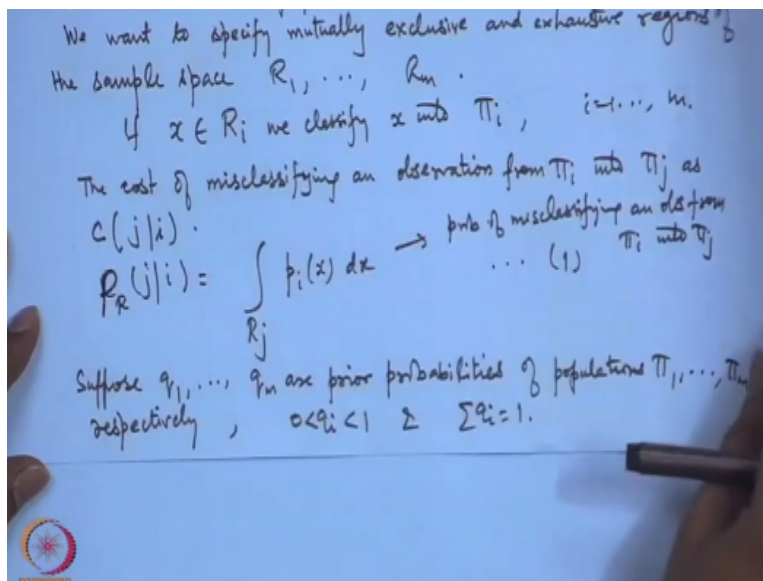
So, classifying an observation into one of several populations. So, supposed we are having populations  $\pi_1, \pi_2, \dots, \pi_m$ , these are  $m$  populations and we are considering the associated density functions say  $p_1(x), p_2(x), \dots, p_m(x)$ , etc., okay. So, we wish to classify or we can find out

m mutually exclusive regions. We want to specify say mutually exclusive and exhaustive regions of the sample space, say  $R_1, R_2, \dots, R_m$ .

So, if the observation  $x$  belongs to say  $R_i$ , we classify  $x$  into the  $i$  population for  $i = 1 \dots m$ . We can also consider the cost of misclassification as we have done earlier. The cost function we can introduce the cost of misclassifying an observation which is actually from say  $\pi_i$  but we classify into  $\pi_j$ . Then, we call this function as  $C_{ji}$ . Now, we can define it like this  $C_{rj}$  given  $i$  = see the observation is initially from the  $i$ -th 1 but we have classified it as into  $J$ th 1.

So, the probability of misclassification or the cost of misclassification can be considered like this. Now, if you remember the case of classification into 2 populations, I had considered one particular case when we fix the initial proportions of the population as  $q_1$  and  $q_2$  where  $q_1 + q_2 = 1$ . In a similar way, if I have  $m$  populations I may consider the case when the initial proportions of these populations are known. We call them prior probabilities say  $q_1, q_2, \dots, q_m$ .

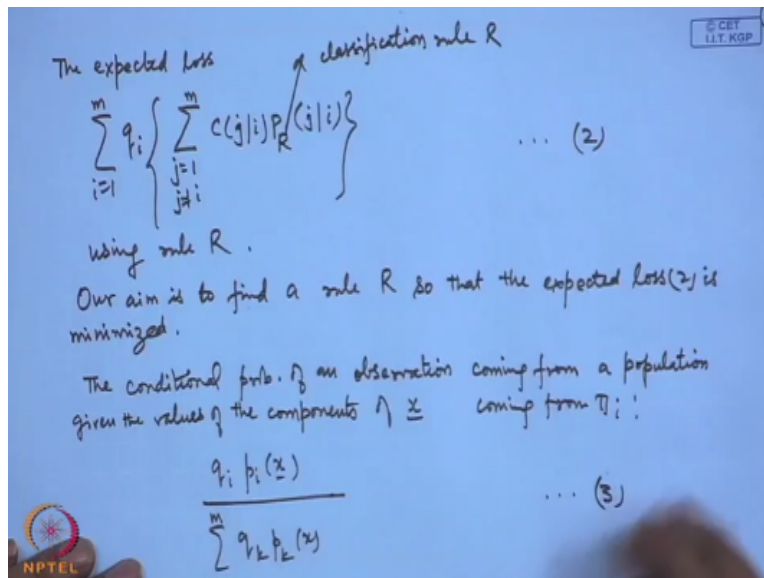
**(Refer Slide Time: 04:35)**



Suppose  $q_1, q_2, \dots, q_m$  are prior probabilities of populations by  $\pi_1, \pi_2, \dots, \pi_m$  respectively, that means  $0 < q_i < 1$  and  $\sum q_i = 1$ . If I am considering the expected loss of classifying  $i^*j$  then I can consider based on the prior probabilities the total expected loss. So, the total expected loss can be defined, so let us consider see we are having. Okay, I just made a small mistake here. This should be  $P_r$ , that is the probability of misclassification. This is the cost.

So, this is actually the probability of misclassifying an observation from  $\pi_i$  as  $\pi_j$ . So,  $P_{Rj}$  given  $i$  and the cost of misclassifying is  $C_{ji}$ . So, if I consider  $C_{ji} * P_{Rj}$ , then this will become the expected cost.

**(Refer Slide Time: 06:10)**



Now, let us consider the total expected loss. So, let us consider this,  $C_j$  given  $i$ ,  $P_{Rj}$  given  $i$ . This  $R$  denotes the classification rule. This is the classification rule. Now, what we have done here is that observation is from  $i$ -th population and we are classifying it as into the  $J$ th population. So, this is the expected cost or expected loss you can consider. Now, you may put it into any of the remaining populations. So, we vary  $j$  from 1 to  $m$  where  $j$  is  $\neq i$ . Now, this term which is written here.

So, now the observation is from the  $i$ -th population and we can put into any other population, other than the  $i$ -th population. So, this is the total expected loss that is coming there. Now, the probability or the proportion of the  $i$ -th population that is  $q_i$ . So, I multiplied by that and then I sum over all  $i$  from 1 to  $m$ , then this is the expected loss using classification rule  $R$ , so this is the expected loss that we can consider.

You can remember the value which I wrote for the case of 2 populations. In the case of 2 populations, let me show you the expression here which we discussed earlier. I will show you the

thing and then you can compare, so it will be clear that how this has been obtained here. If you remember the case of 2 populations, we had only 2 values  $C_{21}$  and  $C_{12}$ . In this case, we have several values  $C_{ji}$  where  $i$  and  $j$  both can vary from 1 to  $m$  where  $i \neq j$ .

So, this is the difference that is coming here. In place of 2 values  $C_{21}$  and  $C_{12}$ , now I have  $C_{ji}$  for all  $j \neq i$  and for all  $i$ . So, total number of values will be  $m \cdot (m-1)$  that you will get here. Let me also show you the expected loss that we had considered here. So, the term that I wrote here  $Pr_j$  given  $i$ . In the case of 2 populations, we had only 2 values  $Pr_2$  given 1 and  $Pr_1$  given 2.

Now, I have  $m \cdot (m-1)$  values once again that will be there and the expected loss of misclassification was only  $C_{21} Pr_2 + C_{12} Pr_1$ . Now, you compare this with the value that I wrote just now because any observation from the  $i$ -th 1 can get into any of the other than the  $i$ -th 1 one, then you consider all such cost then you add them. Then you look at the  $i$ -th 1 and multiply by the prior probability of that and then you sum over.

So, this is the expression that you will be getting. So, this is the full explanation of the expected cost as compared to the case of 2 populations. So, as you can see here the expression becomes much more complex here. However, our aim or the motive remains the same, that is to minimize the expected cost of classification, expected loss by misclassification. So, our aim is to find a rule  $R$  so that the expected loss is minimized.

As again in the case of 1, we had considered  $q_1 p_{1x} / (q_1 p_{1x} + q_2 p_{2x})$ , in a similar way, we can consider the conditional probability of an observation coming from a population given the values of the components of  $x$ . So, given that it is coming from  $\pi_i$ . So, that is defined as  $q_i p_{ix} / \sum_k q_k p_{kx}$ . Earlier it was  $q_1 p_{1x} / (q_1 p_{1x} + q_2 p_{2x})$  or  $q_2 p_{2x}$  divided by the same term, but here now I will have all the  $m$  terms in the denominator.

**(Refer Slide Time: 11:52)**

© CET  
I.I.T. KGP

If we classify the observation as from  $\pi_j$ , the expected loss is

$$\sum_{i=1}^m \frac{q_i p_i(x)}{\sum_{k=1}^m q_k p_k(x)} C(j|i) \quad \dots (4)$$

We minimize the expected loss if we choose  $j$  so that (4) is minimized.

i.e. we consider

$$\sum_{i=1}^m q_i p_i(x) C(j|i) \quad \dots (5)$$

for all  $j$  and select that  $j$  that gives the minimum

**Theorem:** If  $q_1, \dots, q_m$  are prior probs. of  $\pi_i$ , classify  $x$  into  $\pi_k$  if

$$R_k: \sum_{i=1}^m q_i p_i(x) C(k|i) < \sum_{i=1}^m q_i p_i(x) C(j|i), \quad \begin{matrix} j=1, \dots, m \\ j \neq k \end{matrix} \quad \dots (6)$$

If we classify the observation as from  $\pi_j$ , then the expected loss is  $\sum_{i=1}^m q_i p_i(x) / \sum_{k=1}^m q_k p_k(x)$  multiplied by  $C(j|i)$ ,  $i = 1$  to  $m$ ,  $i \neq j$ . We minimize the expected loss if we choose  $j$ , so that (4) is minimized, that is we consider the term  $\sum_{i=1}^m q_i p_i(x) C(j|i)$  given  $i$  because the denominator is common here for all  $j$  and select that  $j$  that gives the minimum. So, in principle if you look at this is a direct extension of the case of 2 populations.

If of course there may be a case when 2 different  $j$  give you the same value, in that case you can choose the one, well it does not matter because then whichever you choose it will give the same minimizing constant. So, now I consider this procedure assigns the value. So, we are assigning towards a  $j$ . So, that is region  $R_j$  here.

So, we consider then the following result then that if  $q_1, \dots, q_m$  are prior probabilities of  $\pi_i$  and the cost function is given here, then the region of classifying  $R_k$  is given, so  $R_k$  region is  $\sum_{i=1}^m q_i p_i(x) C(k|i) < \sum_{j=1}^m q_j p_j(x) C(j|i)$  given  $i$ , here  $i = 1$  to  $m$ ,  $i \neq j$  and here it is  $i = 1$  to  $m$ ,  $i \neq k$ , that means for the  $k$ th 1 if this is the minimum, then you are getting the rule, that is you should classify  $X$  into  $\pi_k$  if this is happening.

I will not get into the proof of this. In fact, the proof is almost the generalization of the proof for the 2 population, that means if I consider any other rule which is minimizing, then I can consider the expected loss from the 2 given rules and write down the difference and as in the case of 2

populations, you can consider the conditions for greater than or equal to thing. So, it will come immediately. So, I am skipping the proof here.

**(Refer Slide Time: 15:47)**

The conditional expected loss of the observation comes from  $\pi_i$

$$\sum_{\substack{j=1 \\ j \neq i}}^m C(j|i) P_R(j|i) = r_R(i)$$

A procedure  $R$  is at least as good as procedure  $R^*$  if

$$r_R(i) \leq r_{R^*}(i) \quad , \quad i=1, \dots, m$$

If strict inequality holds for at least some  $i$ , then  $R$  is said to be better than  $R^*$ .

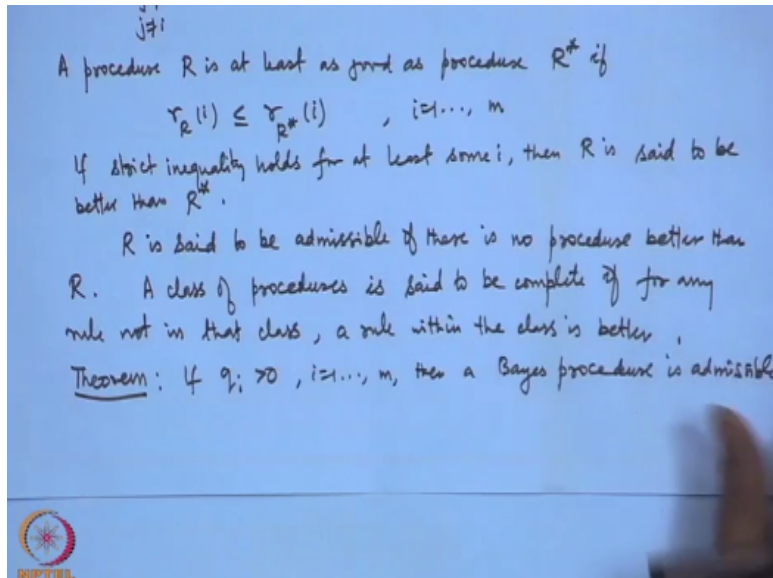
$R$  is said to be admissible if there is no procedure better than  $R$ .

A class of procedures is said to be complete if for any rule not in that class, a rule within the class is better.

Like the case of 2 populations, we can consider the optimality criteria like admissibility, Bayes rules etc. Let me just formally define that here. We can consider the conditional expected loss if the observation comes from  $\pi_i$ , that is  $\sum C_j$  given  $i$ ,  $P_{Rj}$  given  $i$ , this term I wrote earlier, this one basically. So, I am writing this one. So, this is for  $j = 1$  to  $m$ ,  $j \neq i$ . I will use the notation  $r_R i$  here.

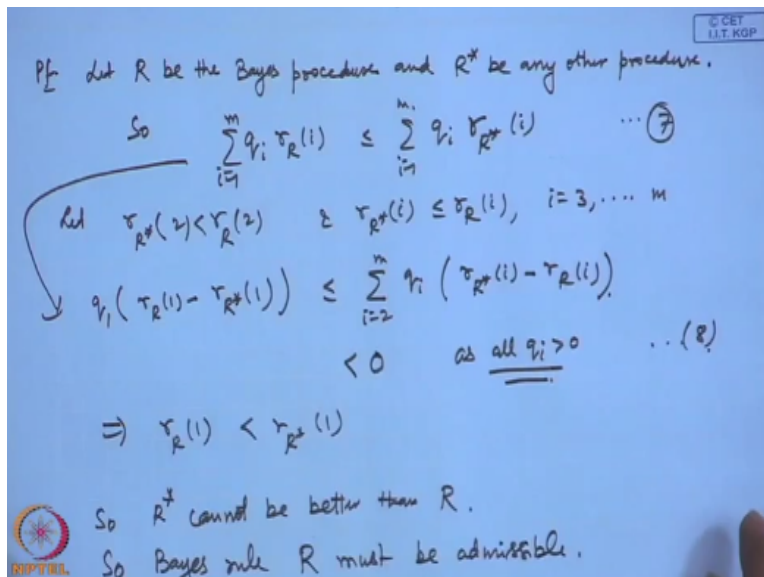
So, a procedure  $R$  is at least as good as procedure  $R^*$  if and only if we are having  $r_R i \leq r_{R^*} i$  for  $i = 1$  to  $m$ . If strict inequality holds for at least some  $i$ , then  $R$  is said to be better than  $R^*$ .  $R$  is said to be admissible if there is no procedure better than  $R$ . A class of procedures is said to be complete if for any rule not in that class, a rule within the class is better. So, these definitions are similar to the one which I gave for the case of 2 populations.

**(Refer Slide Time: 18:49)**



We can consider that this result is also similar as we had earlier that if  $q_i$  is positive, then a Bayes procedure is admissible. Once again, the proof is almost the same. Let me exhibit at least this proof here.

**(Refer Slide Time: 19:29)**

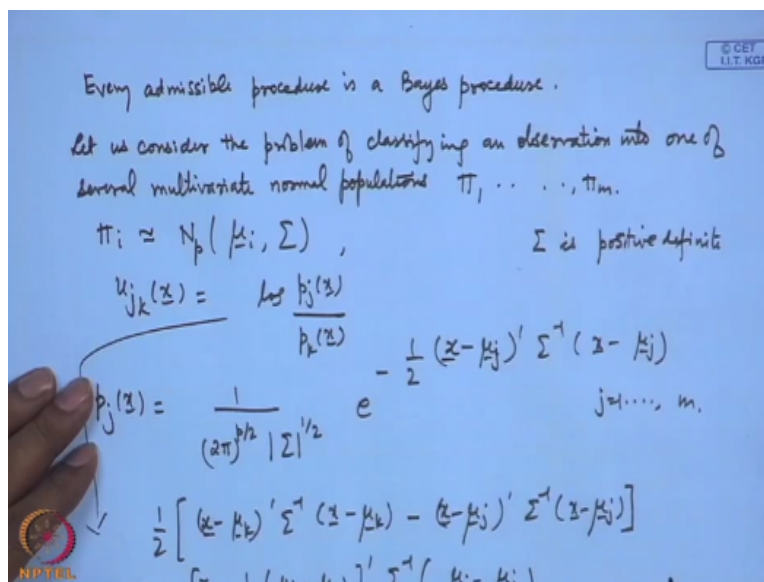


Let  $R$  be the Bayes procedure and  $R^*$  be any other procedure. Now, Bayes procedure will give you the expected loss  $\leq$  the expected loss corresponding to the rule  $R^*$ , okay. Now, let us take say one of the components to be strictly smaller than the corresponding component for the rule  $R$ . So, what I am saying is that actually if I want to show that  $R^*$  is better than  $R$ , then we should have  $r_{R^*}(i) \leq r_R(i)$  for all  $i$  and strict inequality for at least one value.

So, for the time being I am assuming less than or equal to for 2...m and strict inequality at least for the 2, then let us see what happens. So, if we substitute it here, so what we get  $q_1 r_{R1} - r_{R^*1}$ , this will be  $\leq \sum_{i=2}^m q_i r_{Ri} - r_{R^*i}$ . Now, let us look at these 2, I am writing for 2 to m, so for 3 to m, it is less than or equal to and at least there is one strict inequality and all  $q_i$  are positive.

Therefore, this will be strictly  $< 0$  as all  $q_i$  are positive. If this is happening, this is implying that  $r_{R1}$  is strictly  $< r_{R^*1}$ , so  $R^*$  cannot be better than  $R$ . So, that means the Bayes rule  $R$  must be admissible. Actually, you can see that the proof is similar to the case for the case of 2 populations, that is  $m = 2$ . There I had taken only strict inequality for 2 and then for 1 I got the reverse one and similarly for the other case. Now, if the cost functions are given then also the Bayes procedure are admissible.

**(Refer Slide Time: 22:55)**



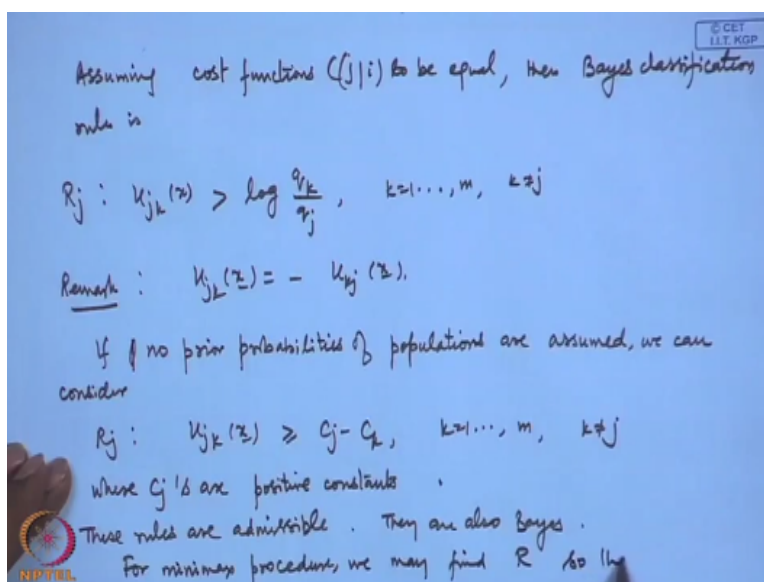
The converse of this result is also true, that is every admissible procedure is also a Bayes procedure. I will not give the proof of this also as the proof is similar to the case of 2 populations. Now, what we want to do is that let us consider the classification of multivariate normal population. So, let us look at this problem which had considered for the case of  $m = 2$ . So, let us consider the problem of classifying an observation into one of several multivariate normal populations.



The populations are  $\pi_1, \pi_2, \dots, \pi_n$  where  $\pi_i$  is the population normal  $\mu_i \sigma_i$ . So, let us define say  $U_{jk}$  function that is equal to log of  $P_{jk}/P_k$ . Actually hear,  $P_{jk}$  will be  $1/2\pi_i$  to the power  $P/2$  determinant of  $\sigma_i$  to the power  $1/2$  e to the power  $-1/2(x-\mu_j)$ ,  $j$  I am writing so it should be  $j$  here, prime  $\sigma_i$  inverse  $x-\mu_j$ . So I am assuming  $\sigma_i$  is positive definite because I am writing down the existence of the density function for  $j = 1$  to  $m$ . So, this quantity that is the ratio here log of  $P_{jk}/P_k$ .

You can consider here, see this term will get cancelled out. So, you will get e to the power  $1/2$  and the term will be corresponding to  $k$  in the numerator and  $j$  in the denominator and then you take log. So, e will go way, so you will get basically  $1/2$  of  $x-\mu_k$  prime  $\sigma_i$  inverse  $x-\mu_k$   $x-\mu_j$  prime  $\sigma_i$  inverse  $x-\mu_j$ , that is equal to after simplification  $x-1/2 \mu_j + \mu_k$  prime  $\sigma_i$  inverse  $\mu_i - \mu_j$ .

**(Refer Slide Time: 26:49)**



Then, the regions of classification, if we apply this method that is we are considering  $\sigma_i$   $q_i R_i \leq$  this. Basically, we have mentioned here that the one which is minimizing this. So, the procedure will become, if we assume the cost functions that is  $C_{ji}$  to be say equal, then the Bayes classification rule is  $R_j$  that is classify the observation into  $j$ th population if  $U_{jk}$  is  $>$  a log of  $q_k/q_j$  for  $k$  equal to  $1$  to  $m$ ,  $k \neq j$ .

We can notice here that this function  $U_{jk}$  they satisfy the symmetry property. So, actually what

kind of regions are these. If you see it carefully, they are nothing but the  $R_i$  are actually bounded by the hyperplanes type of region because  $x - \frac{1}{2} \mu_j + \mu_k$  prime sigma inverse  $\mu_j - \mu_j$ . So, what kind of regions you will be getting. You will be getting the regions of the type of the hyperplane.

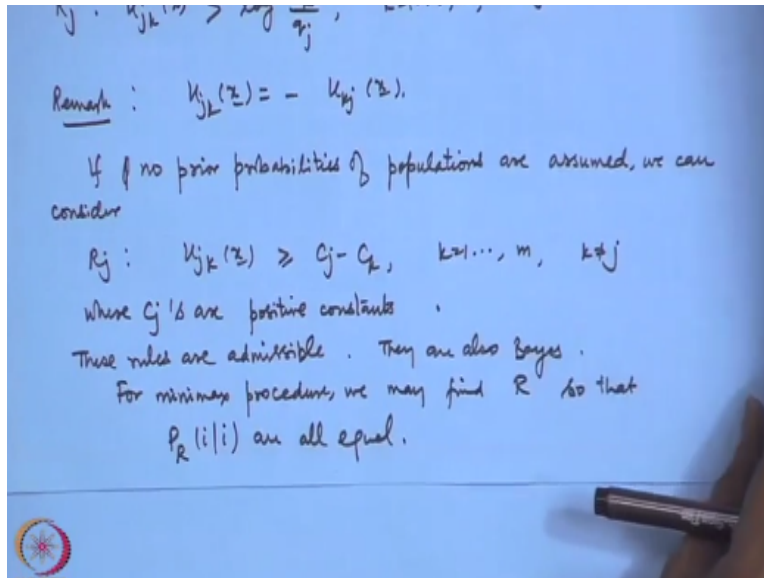
So, if the mean  $(\bar{x})$  (28:55)  $m-1$  dimensional hyperplane, then  $R_i$  is bounded by  $m-1$  hyperplanes that you will be getting, because the  $X$  value that will give you a hyperplane,  $X$  greater than something or less than something and if the prior probabilities are not given, then in place of  $\log$  of  $q_k/q_j$ , you can put some value and in order to maintain some sort of symmetry of representation, actually what is the value of  $\log q_k - \log$  of  $q_j$  because these are probabilities.

So, basically you are getting them to be negative because they are lying between 0 and 1. So, we can write  $\log$  of  $q_k - \log$  of  $q_j$ . So, with minus signs were are getting, so we can put a  $\log$  of  $q_j$  before because there is a minus sign, so we can put it in terms of the non-negative values. If no prior probabilities of populations are assumed, we can consider  $R_j$  as  $U_{jk} x \geq C_j - C_k$  for  $k=1$  to  $m$ ,  $k \neq j$ , where  $C_j$ 's are positive constants.

Actually, any rule of this type is a Bayes rule. So, in case of prior probabilities, if we are putting some other numbers we can actually define respective probabilities in such a way that they will be equal to something. So, all such procedures, they will be giving you the Bayes admissible rules. So, basically, this is you can say minimal complete class of the classification procedures for classifying into one of several populations. These rules are admissible rules.

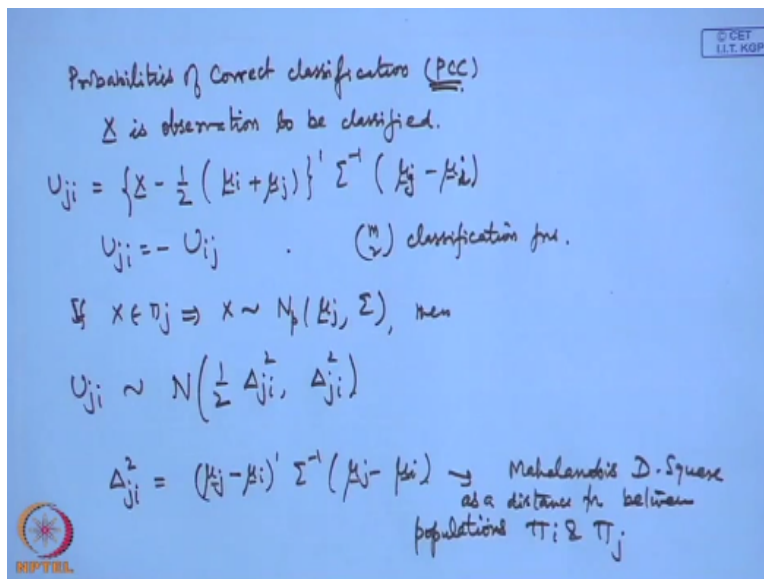
They are also Bayes because the class of Bayes rules and admissible rules is the same here. Now, if you want to find out a minimax procedure, then we can consider say probability of the correct classification and we can make them to be equal.

**(Refer Slide Time: 32:00)**



For minimax procedure, we may find R so that  $P_{Ri}$  given  $i$  are all equal.

**(Refer Slide Time: 32:13)**



Let us look at what are the probabilities of correct classification which we can also call PCC. So,  $x$  is the observation to be classified, then we can consider say  $U_{ji} = x - \frac{1}{2}(\mu_i + \mu_j)$  transpose sigma inverse  $\mu_i - \mu_j$ . This is the classification function that we got and of course  $U_{ji} - U_{ij}$  they are related here, so that means basically we can consider MC2 classification functions, because we do not have to consider both here, MC2 classification functions are there.

Of course, this is because if they span in  $m-1$  dimensional hyperplane. So, now if  $x$  belongs to  $\pi_j$ , that means  $x$  is having  $N_p(\mu_j, \Sigma)$  distribution, then what is the distribution of  $U_{ji}$ . See

this is normal and we can apply the linearity properly. So, this will become actually  $\mu_j - \mu_i$  and here also you are having this thing here. So, this is  $\mu_j - \mu_i$  prime  $\Sigma^{-1}$   $\mu_j - \mu_i$   $1/2$  can be taken outside.

If I define the term say  $\Delta_{ji}^2$  which is a generalization of the Mahalanobis D square function which I wrote in the case of 2 populations, then this is equal to  $\mu_j - \mu_i$  prime  $\Sigma^{-1}$   $\mu_j - \mu_i$ . So, in terms of this, this is actually Mahalanobis D square as a distance function between populations  $\pi_i$  and  $\pi_j$ . So, then this is equal to normal  $1/2 \Delta_{ji}^2$ ,  $\Delta_{ji}^2$ . Also, we can look at the covariance of  $U_{ji}$  and  $U_{jk}$ .

**(Refer Slide Time: 35:36)**

The covariance ~~matrix~~ between  $U_{ji}$  &  $U_{jk}$

$$\Delta_{jk,ji} = (\mu_j - \mu_k)' \Sigma^{-1} (\mu_j - \mu_k).$$

$$P_R(j|i) = \int_{c_j - \infty}^{\infty} \dots \int_{c_j - \infty}^{\infty} f_j dx_{j1} \dots dx_{j,j-1} dx_{j,j+1} \dots dx_{jm}$$

where  $f_j$  is the density of  $U_{ji}$ ,  $i \neq j$

We can choose  $c_j$ 's so that  $P_R(j|i)$  is equal for all  $j$ .

If the parameters are not known, then we can substitute estimates

$$\hat{\mu}_i = \bar{x}_i, \quad \hat{\Sigma}_i = \frac{1}{n_i} \sum_{j=1}^m \sum_{l=1}^{n_i} (x_{ij} - \bar{x}_i) (x_{il} - \bar{x}_i)'$$

when we have random sample  $(x_{i1}, \dots, x_{in_i})$  from  $\pi_i$

So, this is a scalar because  $U_{ji}$  is a scalar function. The covariance matrix between  $U_{ji}$  and  $U_{jk}$ , so that is equal to I use a notation  $\Delta_{jk,ji} = \mu_j - \mu_k$  prime  $\Sigma^{-1}$   $\mu_j - \mu_k$ . In the classification rule when the prior probabilities are not fixed in advance, then we have to determine constant  $C_j$ ,  $C_k$ , etc which I mention that we can choose them to be nonnegative. So, we can consider the probability of classifying in  $2_j$  when it is from  $j = F_j$  that is the observation is from the  $j$ th 1,  $d_{j1}$  and so on...  $d_{j,j-1}$ ,  $d_{j,j+1}$  and so on  $d_{jn}$  and these are  $C_j - C_1$  and so on...  $C_j - C_n$  the because upper side is infinity here.

Where  $F_j$  is the density of  $U_{ji}$  for  $i \neq j$ . So, we can choose  $C_j$  so that  $P_{Rj}$  given  $j$  is equal for all  $j$ . Now, the another situation arises if the parameters are not equal, they are not known then we can

substitute estimates. For example,  $\mu_i$  head can be  $\bar{x}_i$  and similarly you can have  $\sigma_i^2 = 1/n_i - m$ ,  $i = 1$  to  $m$  and  $X_{ij} - \bar{X}_i$ ,  $j$  is equal 1 to say  $n_i$ ,  $i = 1$  to  $m$ . When we have random samples  $x_{i1}, x_{i2}, \dots, x_{in_i}$  from  $\pi_i$  then we can consider these estimates.

**(Refer Slide Time: 39:36)**

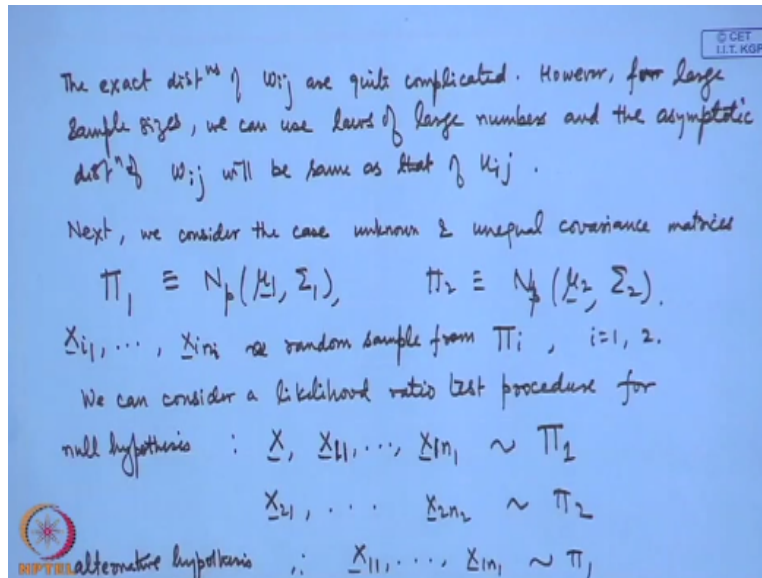
$$P_R(j|i) = \int_{c_j - \epsilon}^{\infty} \dots \int_{c_j - \epsilon}^{\infty} f_j dA_{j,1} \dots dA_{j,j-1} dA_{j,j+1} \dots dA_{j,m}$$

where  $f_j$  is the density of  $U_{ji}$ ,  $i \neq j$   
 We can choose  $c_j$ 's so that  $P_R(j|i)$  is equal for all  $j$ .  
 If the parameters are not known, then we can substitute estimates  
 $\hat{\mu}_i = \bar{x}_i$ ,  $\hat{\sigma}_i^2 = \frac{1}{n_i - m} \sum_{i=1}^m \sum_{j=1}^{n_i} (X_{ij} - \bar{x}_i)^2 (X_{ij} - \bar{x}_i)'$   
 when we have random sample  $(x_{i1}, \dots, x_{in_i})$  from  $\pi_i$   
 $W_{ij} \rightarrow W_{ij}(x) = \left( x - \frac{1}{2}(\bar{x}_i + \bar{x}_j) \right)' S^{-1} (\bar{x}_i - \bar{x}_j)$

The analog of  $U_{ij}$ , this will become  $U_{ij}$  that will be  $x - 1/2 \bar{x}_i + \bar{x}_j$  prime  $S$  inverse  $\bar{x}_i - \bar{x}_j$  bar. Now, as we discussed the case of 2 normal populations, the distribution theory for this part is somewhat more complicated; however, the exact distribution theory is not very difficult because strong of large numbers will hold and if I take here a large sample sizes and then you can consider here that  $\bar{x}_i$  will converge to corresponding  $\mu_i$ ,  $\bar{x}_j$  will go to  $\mu_j$ ,  $S$  inverse will go to  $\sigma$  inverse in probability and so on.

Therefore, the asymptotic distribution of  $U_{ij}$ ,  $W_{ij}$  will be almost the same as the  $U_{ij}$ . So, the problem can be handled.

**(Refer Slide Time: 40:54)**



So, the exact distributions of  $W_{ij}$  are quite complicated; however, for large sample sizes we can use laws of large numbers and the asymptotic distribution of  $W_{ij}$  will be same as that of  $U_{ij}$ . So, this problem can be solved. Now let us also go back to one of the problems that I discussed earlier that is classifying into 2 multivariate populations when the variance covariance matrices were unequal.

I discussed the rule when the 2 populations had known parameters, so if you remember the rule that I had mentioned here. It was given by this that is when  $\sigma_1$  and  $\sigma_2$  are different, I mentioned that in place of hyperplanic regions, you are actually getting much more complex regions because I mention that one of them becomes central chi-square distribution. So, this region becomes much more complicated.

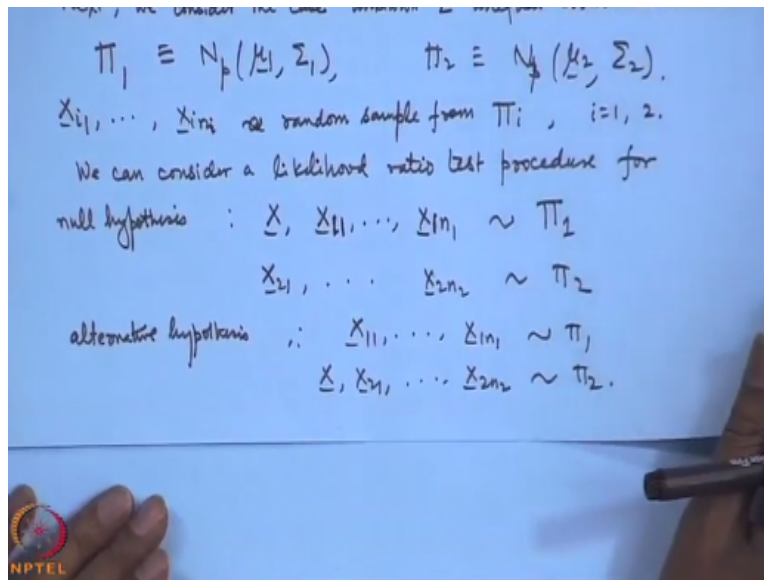
Now, we also consider this case for the unknown  $\sigma_1$  and  $\sigma_2$  case. In that case, we have to substitute the estimators. So, let me briefly discuss this case also. Next, we consider the case of unknown and unequal covariance matrices. So, in particular let us consider say  $\pi_1$  as the population  $N_p(\mu_1, \sigma_1)$  and  $\pi_2$  is the population  $N_p(\mu_2, \sigma_2)$ . So, let me go back to the expression that I derived earlier.

The expression that we obtained was in fact there was a power here which I would have missed at that time. It should be power half here and power half here also and it is e to the  $1/2 \times \mu_2$

prime  $\Sigma_2^{-1} (x - \mu_2)'$   $\Sigma_1^{-1} (x - \mu_1)$ . So, we can consider a likelihood ratio procedure. We are having the samples here say  $x_{i1}$  and so on... $x_{in_1}$ . This is from a random sample from  $\pi_i$ ,  $i$  is equal 1 to 2.

We can consider a likelihood ratio test procedure for null hypothesis, that is the observation  $x$ ,  $x_{11}, x_{12}, \dots, x_{1n_1}$ , they are from  $\pi_1$  and  $x_{21}, x_{22}, \dots, x_{2n_2}$  this is from  $\pi_2$  and the alternative hypothesis will be that is  $x_{11}, \dots, x_{1n_1}$  this is from  $\pi_2$ .

**(Refer Slide Time: 46:02)**



$x_{21}$  and so on... $x_{2n_2}$  this is from  $\pi_2$ . Now, in the likelihood ratio procedure, I have to consider the maximization of the likelihood function under both null and alternative hypothesis. So, in the null hypothesis, I will have  $n_1+1$  observations from  $\pi_1$  and  $n_2$  observations from  $\pi_2$ . In the alternative hypothesis, I will have  $n_1$  observations from  $\pi_1$  and  $n_2+1$  observations from  $\pi_2$ . Since all the parameters are unknown and unequal, this is simply reducing to the problem of finding out maximum likelihood estimators for these cases.

**(Refer Slide Time: 46:57)**

In the likelihood ratio procedure, we are required to determine the maximum value of the likelihood fn. under null as well as alternative hypothesis. Thus we find MLE's under both cases:

Under  $H_0$ : MLE's are  $\hat{\mu}_1(1) = \frac{n_1 \bar{X}_1 + X}{n_1 + 1}$ ,  $\hat{\mu}_2(1) = \bar{X}_2$

$\hat{\Sigma}_1(1) = \frac{1}{n_1 + 1} \left[ A_1 + \frac{n_1}{n_1 + 1} (X - \bar{X}_1)(X - \bar{X}_1)' \right]$

$\hat{\Sigma}_2(1) = \frac{1}{n_2} A_2$

$A_i = \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)(X_{ij} - \bar{X}_i)'$

Under  $H_1$ :  $\hat{\mu}_1(2) = \bar{X}_1$ ,  $\hat{\mu}_2(2) = \frac{n_2 \bar{X}_2 + X}{n_2 + 1}$

$\hat{\Sigma}_1(2) = \frac{1}{n_1} A_1$ ,  $\hat{\Sigma}_2(2) = \frac{1}{n_2 + 1} \left[ A_2 + \frac{n_2}{n_2 + 1} (X - \bar{X}_2)(X - \bar{X}_2)' \right]$

So, we can easily write down the maximum likelihood estimators. In the likelihood ratio procedure, we are required to determine the maximum value of the likelihood function under null as well as alternative hypothesis. Thus, we find MLE's under both cases. So, this we can write easily because the procedure for finding out the MLE is known in the case of multivariate normal distribution.

We know actually that the sample mean and the sample covariance matrix, they are the maximum likelihood estimators. So, under the null hypothesis MLE's are, so we write it as  $\mu_1$  head  $1 = n_1 \times 1$  bar  $+ X/n_1 + 1$  because this is a sum of all the observations from the first sample plus  $X$  because we are saying that it is coming from the first population,  $\mu_2$  head  $1$ , so  $1$  means basically under the null hypothesis.

This is given by  $\bar{X}_2$  and  $\sigma_1$  head  $1$  that will be  $1/n_1 + 1$  sigma, so we call it say  $A_1 + n_1/n_1 + 1$   $(X - \bar{X}_1)(X - \bar{X}_1)'$  and  $\sigma_2$  head  $1 = 1/n_2$   $A_2$ . Here  $A_i$  are  $\sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)(X_{ij} - \bar{X}_i)'$ ,  $j = 1$  to  $n_i$ . Under  $H_1$  that is the alternative hypothesis, actually this null hypothesis I am calling  $H_0$  and this alternative hypothesis I am calling  $H_1$ . So, under  $H_1$  this will turn out to be  $\mu_1$  head  $2$  that will become  $\bar{X}_1$  and  $\mu_2$  head  $2$ . Now here I will have  $n_2 + 1$  observations from the second one, so it will be  $n_2 \times 2$  bar  $+ X/n_2 + 1$ .

For  $\sigma_1$  head, this will become equal to  $1/n_1$   $A_1$  and for  $\sigma_2$  head, this will become



$1/n_2+1 A_2+n_2/n_2+1 \bar{x}-x_2 \bar{x}-x_3 \bar{x}$  transpose.

(Refer Slide Time: 51:25)

The likelihood ratio will give

$$\frac{|\hat{\Sigma}_1(z)|^{n_1/2} |\hat{\Sigma}_2(z)|^{(n_2+1)/2}}{|\hat{\Sigma}_1(1)|^{(n_1+1)/2} |\hat{\Sigma}_2(1)|^{n_2/2}}$$

$$= \frac{[1 + (\bar{x} - \bar{x}_2)' A_2^{-1} (\bar{x} - \bar{x}_2)]^{\frac{n_2+1}{2}} (n_1+1)^{\frac{1}{2}(n_1+1)p} n_2^{\frac{n_2 p}{2}} |A_2|^{1/2}}{[1 + (\bar{x} - \bar{x}_1)' A_1^{-1} (\bar{x} - \bar{x}_1)]^{\frac{(n_1+1)}{2}} n_1^{\frac{1}{2}n_1 p} (n_2+1)^{\frac{1}{2}(n_2+1)p} |A_1|^{1/2}}$$

We classify  $\bar{x}$  into  $\pi_1$  if the ratio is more than 1,  
else classify  $\bar{x}$  into  $\pi_2$ .

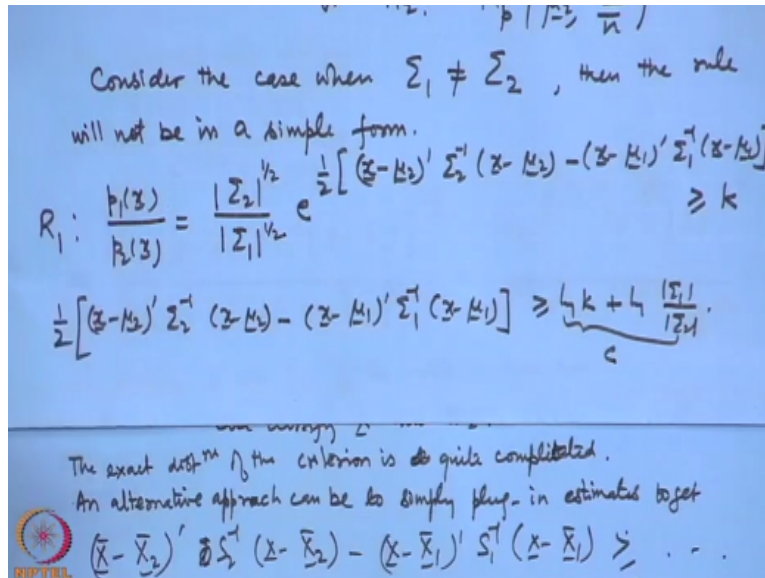
The exact dist<sup>n</sup> of the criterion is quite complicated.  
An alternative approach can be to simply plug-in estimates

So, if we consider the likelihood ratio criteria that will give us the likelihood ratio. So, the exponent term will get cancelled out and we will get  $\sigma_1$  head to the power  $n_1/2$   $\sigma_2$  head 2 to the power  $n_2+1/2$  divided by  $\sigma_1$  head 1 to the power  $n_1+1/2$  and  $\sigma_2$  head 1 to the power  $n_2/2$  which we can also write after simplification as  $1 + (\bar{x} - \bar{x}_2)' A_2^{-1} (\bar{x} - \bar{x}_2)$  whole to the power  $n_2+1/2$  /  $1 + (\bar{x} - \bar{x}_1)' A_1^{-1} (\bar{x} - \bar{x}_1)$  whole to the power  $n_1+1/2$  \*  $n_1+1$  to the power  $1/2$   $n_1+1$   $p$   $n_2$  to the power  $n_2 p/2$ , determinant of  $A_2$  to the power  $1/2$  /  $n_1$  to the power  $1/2$   $n_1+1$   $p$   $n_2+1$  to the power  $1/2$ ,  $n_2+1$   $p$   $A_1$  to the power  $1/2$ .

If we consider the costs of misclassification to be the same and the prior probabilities to be equal, we can consider this ratio to be  $> 1$ , then you classify into  $\pi_1$ . So, we classify  $X \in \pi_1$  if the ratio is more than 1. Else classify  $X \in \pi_2$ . Now, this is one criteria that is the likelihood ratio criteria. Let me also again come back to this original observation that I got here. Another thing could be that I substitute direct estimates, that means here I put  $\bar{x}_2$ , here I put  $S_2^{-1}$ , here I put  $S_1^{-1}$ , here I put  $\bar{x}_1$ .

So, in both the cases, the exact distributions of the criteria are not easy. The exact distribution of the criteria is quite complicated.

(Refer Slide Time: 55:04)



An alternative approach can be to simply plug-in estimates into this function, that is  $\bar{x} - \bar{x}_2$  prime sigma, that is  $S_2$  inverse  $\bar{x} - \bar{x}_2$  bar  $\bar{x} - \bar{x}_1$  bar prime S inverse  $\bar{x} - \bar{x}_1$  bar. So, you put this as some greater than or equal to something and less than something. So, this could be another one. Once again, the distribution of the criteria is quite complicated.

Even if I look at the asymptotic distribution for the large sample sizes by applying the laws of large numbers, I will get this one. I have already discussed that if  $x$  belongs to  $\pi_1$ , then this is non-central chi-square. This will be something like a central chi-square. So, the difference of 2 and it is going to be quite complicated. In case it is from  $\pi_2$ , then this one is central and this one is non-central. Once again the exact distribution of these things are difficult to obtain.

So, in particular we are saying is that when the variance covariance matrices are unequal, the classification rules no doubt can be easily found out but in order to obtain desirable rules such as a minimax procedure among them is a difficult task. Because the probabilities of correct classification are the probabilities of misclassification will be quite complicated. Friends, so we have actually discussed so many classification rules.

In fact, I framed a general decision theoretic approach to the classification problem by considering the Bayes decision rule and the criteria of admissible rules, the minimax classification procedure, the minimal complete class, etc and in particular, I showed applications

to the classification of multivariate normal populations. We have considered 2 normal populations and multivariate normal population.

So, I actually wind up the discussion on the problem of classification now. One can consider some other classification procedures which are available nowadays but that can be a subject of full-fledged discussion. I will move over to another topic that is the problem of principal components. So, in the next lecture I will briefly introduce the problem how to determine the principal components and also maybe I will touch up on the canonical correlations. So, that will be the topic of next lecture.