

Applied Multivariate Analysis

Prof. Amit Mitra

Prof. Shramishtha Mitra

Department of Mathematics and Statistics

Indian Institute of Technology, Kanpur

Lecture No. #33

Discriminant Analysis and Classification

So, in the last couple of lectures, what we were looking at was we try to derive classification rules, classification functions; that is partition of the sample space for two populations problem. We had looked at various concepts like what happens to when we look at a deriving a classification rule, when we have got a say for example, a rule which tries to minimize total probability of misclassification. What type of rule to be get when we look at a functions like expected cost of misclassification, and the partition corresponding to one that would minimize an expected class of cost misclassification. We had also seen in the last lecture for specific type of populations namely, normal populations, and we had seen how these optimum rules actually look like for these special type of cases of multivariate normal populations. So, first thing, that we are going to look at today will be performance measures of various classification functions.

So, that gives us some way of comparing various partitions, various classification functions, and then we will move on to looking at a classification problem in a multi population problem. So, we will generalize, whatever we have been drink for two population cases, we will look at generalizing that general c population problem.

(Refer Slide Time: 01:38)

Performance measures for comparing different classification

Recall

$$TPM = p_1 P(2|1) + p_2 P(1|2)$$
 i.e.
$$TPM = p_1 \int_{R_2} f_1(x) dx + p_2 \int_{R_1} f_2(x) dx$$
 optimum error rate (OER) \equiv min TPM rule
 optimum TPM rule : $(R_{1(opt)}, R_{2(opt)})$
 Then
$$OER = p_1 \int_{R_{2(opt)}} f_1(x) dx + p_2 \int_{R_{1(opt)}} f_2(x) dx \quad (1)$$
 Actual error rate (AER) $= p_1 \int_{\hat{R}_{2(opt)}} f_1(x) dx + p_2 \int_{\hat{R}_{1(opt)}} f_2(x) dx \quad (2)$
 $\hat{R}_{1(opt)}$ & $\hat{R}_{2(opt)}$ are based on the learning sample \mathcal{L} .

Let us start today, looking at this performance measures thus, we had looking at performance measures for comparing different classifications, we looking at comparing **comparing** such different classification functions. Let us recall, what we had define earlier, we had a total probability of misclassification, where we are considering still two population problem. So, that was given by p_1 . So, it is coming from first population getting misclassified into the second population this class p_2 times p_1 given 2.

Now this, in terms of the regions that we have a obtain. So, this is integration over the R_2 region, the object is coming from the first population. So, it has got density like this, this plus p_2 into integral over R_1 . And since it is coming from the second population, it has got a density f_2 this term.

Now, based on this total probability of misclassification, we can look at the following measures, which is called the optimum error rate or O E R. Now, that is what is corresponding to the total probability of misclassification rule, which minimizes that total probability of misclassification. So, this is what is corresponding to the minimum T P M rule.

Now, suppose we have got this T P M minimum optimum T P M rule, suppose that is given by this partition R_1 optimum R_2 optimum. So, this gives us the minimum total probability of misclassification, then this O P R expression is given by O E R rather is

given by p_1 into integral over R_2 opt, this into $f_1(x)$, this plus p_2 times integral over this R_1 opt region of the density for the second population $f_2(x)$.

So, this is one such measure for classification, you look at the optimum partitions that you obtain then calculate what is the total probability of misclassification corresponding to that. Now note, that in this particular situation here, when we talk about R_1 opt and R_2 opt, these quantities as we had seen in the earlier lectures they depend on population quantities like unknown mean vector, unknown covariance matrix in case of multivariate normal populations. And hence, for any practical purpose the given data one based on the learning sample, one would be estimating this quantities, and the measure that is based on such quantities is what is call the actual error rate actual error rate or A E R.

Now, this is given by the following expression, that instead of writing it as R_1 opt and R_2 opt, we write it as R_2 opt cap that is a estimate of that particular region based on the given sample data, which is there in the learning sample. So, this is R_1 opt hat $f_2(x)$ this is a measure. Now this R_1 opt R_2 opt caps and this R_2 opt estimate are based on the learning sample, which was denoted by L . So, these are based on the learning sample which is say script l .

Now once again, if you look carefully at this O E R or A E R. Although A E R has got this estimate here, we still keep it as $f_1(x) f_2(x)$.

(Refer Slide Time: 06:20)

Note: OER/AER depend on unknown quantities p_1, p_2 and/or $f_1(x), f_2(x)$.

Apparent error rate (APER): Fraction of observations in the training sample that are misclassified by the classification function

(Confusion matrix based on $\hat{\alpha}$)

		Predicted class		
		π_1	π_2	
Actual class	π_1	(n_{1c})	(n_{1H})	$n_1 = n_{1c} + n_{1H}$
	π_2	(n_{2H})	(n_{2c})	$n_2 = n_{2c} + n_{2H}$

So, these quantities are still unknown, now note that O E R or A E R both of this depend on unknown quantities say for example, p_1 p_2 and or as may be the case with O E R or A E R, $f_1 \times f_2 \times$. And hence, they as such cannot be applied to compute for a given data what is the O E R or A E R corresponding to partition rule.

Now, a measure that is defined as apparent error rate or A P E R. This is the measure that is going to be based on the sample data. And after the classification has been done, one actually would look at the proportion of misclassifications that is being done based on such measure. Now this is going to be defined, as the fraction of observation **fraction of observations** in the training sample, because that is what we have in the training sample, that are misclassified by the classification rule classification rule or the classification function. Now, how is that going to be defined, it is based on something which is called a confusion matrix.

So, one first constructs the confusion matrix after the classification rule is put forward. So, one has got a first start with a learning sample based on the learning sample, learning sample the cases are classified. Then based on those pre classified cases, one has constructed a classification rule. So, the sample classification rule is in place, and using the sample classification rule, one has classified those pre classified examples which where there in the learning sample. And then one looks at this matrix, which is called the confusion matrix. It looks like the following. So, we will be having on one hand predicted class memberships.

So, there can be two such predicted classes p_1 and p_2 , because we have looking at two class problem. And each observation what we have has got another tag, which is there in the learning sample. So, actual class membership. So, there are two possibilities p_1 and p_2 . Now, this confusion matrix based on the learning sample that is script 1.

Now, suppose in this learning sample, these are pre classified example. And hence, after we use the classification function based on the sample data, when we are going to now classify the feature vectors. We are going to come up with some predicted class memberships. Now, this predicted class membership are going to lead us to this numbers say this is n_{1c} n_{1m} .

Now, this **this** n_{1c} is the number of observations coming from the first population p_1 and being correctly classified by the classification rule. So, the predicted class

membership is π_1 corresponding to those observations coming from the first population. So, n_{1c} is the number of correctly classified observations, for in the learning sample coming from the first population. Now n_{1m} is the number of observations, that are coming from the first population there in the learning sample and by the classification function they are misclassified into this π_2 population.

(Refer Slide Time: 12:19)

$$\Rightarrow APER = \frac{n_{1M} + n_{2M}}{n_1 + n_2}$$

So, this is this notation. Now similarly, we can define this as n_{2m} , n_{2m} is the number observations which are actually belonging to the second population π_2 , and by the classification function are getting classified into the π_1 population. And hence, these numbers of this number observations n_{2m} are wrongly classified examples, coming from the second population. And similar to this one, we will have a number n_{2c} here, which is the number of observation cases coming from the second population, and also being classified to belong to the second population using our classification function.

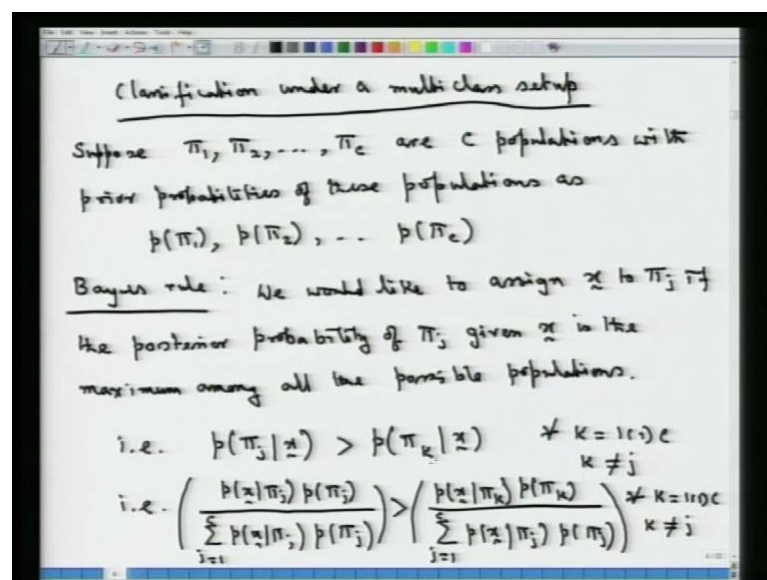
Now, the total cardinality suppose this learning sample is n . Then, we will be having this, suppose is the sum of this two numbers. So, this is n_{1c} plus n_{1m} . So, these this n_1 is the number of observations belonging to the first population in the learning set l , and this is n_2 which is the number of observations belonging to the second population in the learning set. And this is the sum of n_{2c} plus n_{2m} .

Now, if this is what we get as the confusion matrix which is derived, when we have a particular classification function in place. Then these observations are misclassified coming from the first population, and these observations n_{2m} are the number of observations misclassified coming from the second population. And hence when we talk about apparent error rate A P E R, then its fraction of observations in the training sample that misclassified by the classification function. So, this would imply that what we have this A P E R apparent error rate is the fraction of misclassified observations.

So, it is n_{1m} plus n_{2m} , these are the observations which are misclassified that divided by n_1 plus n_2 . So, this is in a perfectly implementable form. So, based on any classification function and the given learning sample I , one can compute this confusion matrix and hence one can get to this apparent error rate.

So, that can be computed for any practical data. So, we end this particular small section on looking at performance measures of for comparing classification functions. Now, we look at the important problem of extending, what we were trying to learn in classification in terms of looking at a multiclass problem.

(Refer Slide Time: 13:19)



So, let us look at that now. So, we are looking at classification under a multiclass setup. So we have, we make provision for more than two populations, suppose we have got now c populations, suppose $\pi_1, \pi_2, \dots, \pi_c$ are c populations c possible populations with prior probabilities of these populations as say $p(\pi_1), p(\pi_2), \dots, p(\pi_c)$.

So the problem is simple, that we have got a multivariate normal, multivariate in general not necessarily normal. We have multivariate population and where in the multivariate population, there are c such possible populations. They may differ in the mean vector, they may differ in the covariance matrix or any other measure characterizing that particular population. And then, these populations has got this prior probabilities p_{π_1} , p_{π_2} , p_{π_c} , and given an a multidimensional observation, we will have to look at in which class this is going to belong to. So, we are trying to a classify a multivariate observation if is a vector into one of these populations.

Now, let us look at the three type of classification rules that one can think of or the type of classification rules that we have derive or two population problem. Let us first look at the base rule. Now, the base rule is going to choose that particular population which has got the highest posterior probability.

Now let us write that, we would like to assign an multivariate observation, a multivariate observation x a particular π_j , if the posterior probability of π_j given this x is the maximum among all the possible populations. That is, what we are trying to do here is that let us denote by this quantity p_{π_j} given x to be posterior probability of the population π_j given x is observed. If this is greater than p_{π_k} given x , if this is true for every k equals to 1 to up to c with k not equal to j , then we will assign x the multivariate observation to the population π_j . Because the posterior probability of π_j given x which is this, if that is the maximum among all possible such posterior probabilities for other populations.

So, we have got a k , which is not equal to k . So, this is basically what we have, now this can be return alternatively in the following form that we have got, the using base theorem what we can write straight away is that probability of x given π_j this into probability of π_j . This divided by summation j equal to 1 to up to c $p(x$ given $\pi_j)$, this a multivariate observations p_{π_j} into this multiplied by p_{π_j} . If that is greater than the corresponding quantity on the right hand side. So, this is probability of x given π_k , this into the prior probability of π_k , that divided by summation j equal to 1 to up to c $p(x$ given $\pi_j)$, this into the prior probability of the j th population. If that is true for every k equal to 1 to up to c with k not equal to j .

So, we have got this two give us the base rule, which looks at which population has got the maximum posterior probability given the observation x . So, this is in a nice form, that we have got this to be the base rule, wherein we can infer that this left hand side here what we have is the posterior probability of π_j given x . And the right hand side is posterior probability of π_k for k not equal to j for all other populations.

(Refer Slide Time: 18:19)

The image shows a handwritten derivation on a whiteboard. The title is "TPM minimizing classification rule". The derivation starts with the equation:
$$TPM = \sum_{i=1}^c p(\pi_i) P(\text{error} | \pi_i)$$
Below this, it says "Let (R_1, R_2, \dots, R_c) be the c classification partition". Then it defines:
$$P(\text{error} | \pi_i) = \int_{R_i^c} p(x | \pi_i) dx$$

$$= \int_{\Omega - R_i} p(x | \pi_i) dx$$
Finally, it concludes:
$$\Rightarrow TPM = \sum_{i=1}^c p(\pi_i) \int_{\Omega - R_i} p(x | \pi_i) dx$$

Now, let us look at the total probability of misclassification approach, T P M minimizing classification rule. (No audio from 18:56 to 19:06) Now how is that going to look like in this multiclass problem, the total probability of misclassification the type of concept that we had introduced for the multiclass problem. This will take the following form that it is summation i equal to 1 to up to c , then probability the a priory probability of the π_i population into the probability of committing, and error given an observation is coming from π_i .

Let us see, what are the terms here, because if we are looking at the term by term here. So, the first term is $p(\pi_1)$ into the probability of classifying that observation which is coming from the first population π_1 , and then putting it into any other population. So, what is this probability of error by the way, this probability of error given π_i , this is going to be given in terms of the partitions that we have suppose we say that, let $R_1 R_2 R_C$ be the classification partition; that means, that if x belongs R_i , we are going to put it into population number i . That is, for every x that is belonging to this region the class

membership is π_1 , for every x that is belonging to R_2 the class membership is going to be π_2 , and for the rest of this also.

So, it is basically that. So, we can write this probability of error, given an observation is coming from π_1 , π_i here. So, that has got a probability which we had denoted earlier $p(x|\pi_j)$. So, this is the density, when we are looking at an observation coming from π_j . Now, this we are actually putting it into some other set here, other than the i th set. So, this is x given π_i , and then this is integral over the complimentary region of R_i . Because if x belongs to R_i , then we are going to correctly classify into π_i , otherwise if we are having x belonging to any other R_i not equal to that particular term for which we are looking at this. And hence, this term can be written as ω minus this R_i . So, this is $p(x|\pi_i) dx$.

(Refer Slide Time: 22:39)

The image shows a handwritten derivation of the Total Probability of Misclassification (TPM) for a multiclass problem. The derivation is as follows:

$$\begin{aligned} \text{i.e. TPM} &= \sum_{i=1}^c p(\pi_i) \left(\int_{\Omega} p(x|\pi_i) dx - \int_{R_i} p(x|\pi_i) dx \right) \\ &= \sum_{i=1}^c p(\pi_i) \left(1 - \int_{R_i} p(x|\pi_i) dx \right) \\ &= \sum_{i=1}^c p(\pi_i) - \sum_{i=1}^c p(\pi_i) \int_{R_i} p(x|\pi_i) dx \\ &= 1 - \sum_{i=1}^c p(\pi_i) \int_{R_i} p(x|\pi_i) dx \end{aligned}$$

\Rightarrow Minimizing TPM w.r.t. the partition (R_1, \dots, R_c) is equivalent to maximizing

$$\sum_{i=1}^c p(\pi_i) \int_{R_i} p(x|\pi_i) dx \quad \text{w.r.t. } (R_1, \dots, R_c).$$

Now, if this error is given by this particular expression, then this would imply that the total probability of misclassification, for this multiclass problem expression is given by $1 - \sum_{i=1}^c p(\pi_i) \int_{R_i} p(x|\pi_i) dx$.

Now, let us see what this term is equal to, this term would be equal to this total probability of misclassification, thus is equal to $\sum_{i=1}^c p(\pi_i) \int_{\Omega} p(x|\pi_i) dx - \sum_{i=1}^c p(\pi_i) \int_{R_i} p(x|\pi_i) dx$. Now this expression is equal to 1. So, what we will be having is $\sum_{i=1}^c p(\pi_i) \int_{\Omega} p(x|\pi_i) dx - \sum_{i=1}^c p(\pi_i) \int_{R_i} p(x|\pi_i) dx$. This is multiplied by 1 minus integral, we leave it as it is. The

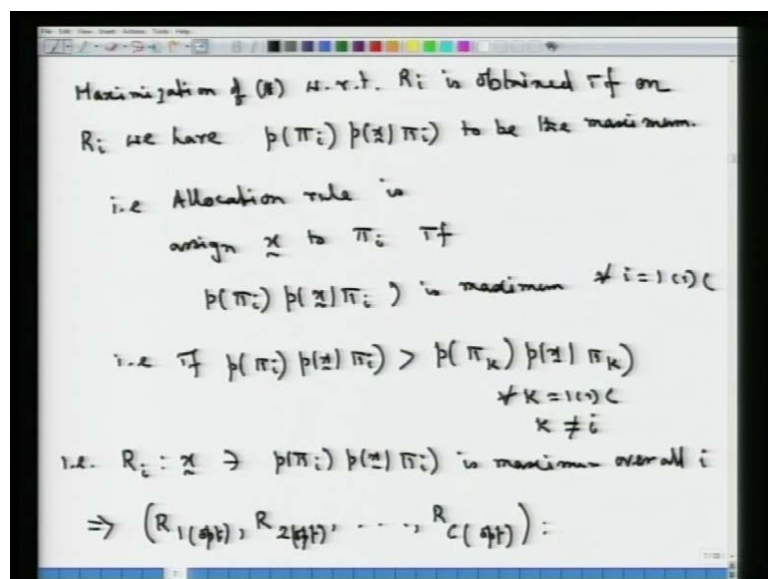
second term $\int p(x) \pi_i dx$. So, the first term, that we have here after we open the bracket is $\sum_{i=1}^c p(\pi_i) \int p(x) \pi_i dx$, this minus $\sum_{i=1}^c p(\pi_i) \int p(x) \pi_i dx$ into $\int p(x) \pi_i dx$.

Now, we note that what is this term equal to, the first term in this expression is equal to 1, because it looks at the prior probability, the sum of the prior probability of all possible c population. And hence, summation of this $p(\pi_i)$ terms will be equal to 1. So, it is 1 minus this term here that it is $\sum_{i=1}^c p(\pi_i) \int p(x) \pi_i dx$ into the term, which will be have their $\int p(x) \pi_i dx$. So, this would imply that the rule or rather minimizing the total probability of misclassification is equivalent to maximizing this expression, which is second expression.

So, minimizing total probability of misclassification with respect to the partition R_1, R_2, R_C is equivalent to maximizing the quantity, which is $\sum_{i=1}^c p(\pi_i) \int p(x) \pi_i dx$. This with respect to the partition R_1, R_2, R_C . Now the **minimize** this partition, that is going to lead us to minimum total probability of misclassification, thus is same as the partition which would maximize the expression which is given by this.

Now, what is that? So, the optimizing partition, which is going to maximize this particular expression star is going to be, that it is going to be the setup all x 's, for which we will have the corresponding term $p(\pi_i)$ into $\int p(x) \pi_i dx$ to be the maximum.

(Refer Slide Time: 25:52)

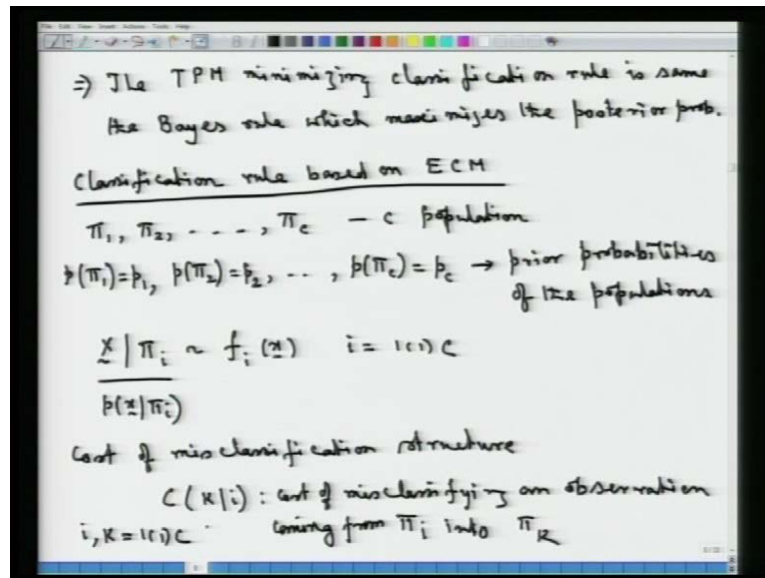


So, what we have here is finally, that the maximization of star with respect to this R_i is obtained, if on R_i we have the corresponding expression $p_{pi i}$ into p_x given $pi i$ to be the maximum this is. So, because we are trying to maximize this particular quantity with respect to R_i . So, with respect to each of this R_i terms here, we are trying to find out $p_{pi i}$ into p_x given $pi i$. If that is maximum, then the expression is going to be maximized. That is, allocation rule is going to be given by the following, allocation rule is assigned x to $pi i$, if we have $p_{pi i}$ into p_x given $pi i$ is maximum for every i equal to 1 to up to c . That is, if we have got $p_{pi i}$ into p_x given $pi i$ to be greater than $p_{pi k}$ into p_x given $pi k$ this is true for every k , equal to 1 to up to c where k is not equal to i .

So, this is what is going to give us a partition, that is we have got in other words R_i is the region of x 's, such that this term $p_{pi i}$ into p_x given $pi i$ is maximum over all I , overall say overall I . So, that is basically what is going to lead us to this particular region, and hence if we have this R_i , we can construct all the other region. So, we will be having the optimum partition, minimizing the total probability of misclassification. So, this would lead us to this $R_1^{opt} R_2^{opt}$, opt in the sense that it is leading us to the partition which is going to be the partition which would minimize the total probability of misclassification.

Now, what we observe is the following that, when we have obtained this classification rule under the par dime of say the total probability of misclassification, and we had a earlier looked at the rule which was looking at the base rule, which also was essentially the same. So, for a base rule what we had was this condition, this posterior probability to be greater than this posterior probability. However, we have got the both the, a denominators in the left hand side and the right hand side to be the same. And hence that is going to have the base rule is going to be based on this expression being greater than the numerator of the right hand side. And hence, the two rules are basically equivalent, because we have got exactly the same rule here.

(Refer Slide Time: 29:33)



So this implies the total probability of misclassification minimizing partition or minimizing classification rule is same as the base rule, which maximizes the posterior probability. (No audio from 30:00 to 30:11)

So, this is another justification of looking at either of these two rules. Now next, we are going to look at classification rule based on E C M. Still on a multiclass problem, so we will now look at classification rule or construction of the classification partition based on the expected cost of misclassification. Let us recall, the structure that is what we have, we have got π_1, π_2, π_c . These are the c populations, c possible populations into which an observation can belong to, and then we have got these as the prior probabilities, say let us write that is to be equal to p_1, p_2, \dots, p_c which is prior probability of the c th population given by P_C . So, these are prior probabilities of the populations.

Now, prior probabilities of the populations, we will also say that say for example, this the density of x , the density of x given π_i say is given by $f_i(x)$ this is for i equal to 1 to up to c . In our earlier notation, we had perhaps denoted this by $p(x | \pi_i)$. So, that density we are denoting by $f_i(x)$.

Now, we have got multiclass problem is c class problem. So, let us also look at the cost of misclassification structure, cost structure or cost of misclassification structure. We define that $C(k|i)$. So, this is the cost of misclassifying an observation, coming from the

The i th population that is π_i into the k th population, that is π_k . So, we have got this c_{ki} is for i, k equal to 1 to up to c .

(Refer Slide Time: 32:58)

Handwritten notes on a whiteboard:

$$C(i|i) = 0 \quad \forall i = 1, \dots, c.$$

(R_1, R_2, \dots, R_c) - partition of classification

$P(k|i)$: prob of an obs from π_i getting misclassified into π_k .

$$P(k|i) = \int_{R_k} f_i(x) dx$$

$$P(i|i) = \int_{R_i} f_i(x) dx = \int_{\Omega - \bigcup_{k \neq i} R_k} f_i(x) dx$$

$$= \int_{\Omega} f_i(x) dx - \int_{\bigcup_{k \neq i} R_k} f_i(x) dx$$

$$= 1 - \sum_{k=1, k \neq i}^c \int_{R_k} f_i(x) dx.$$

Now, we of course, will be having this term with this c_{ki} , that cost of misclassifying an observation from the π_i into π_k . So, there is no misclassification as such and hence to be equal to 0. This is for every i equal to 1 to up to c . So, once we have this particular term in place, then we will look at constructing the region which is going to look at E C M minimizing rule, but before that we need to actually define how the E C M rule E C M looks like under the presence situation. So, this is the type of partition that we are trying to get. So, this is the partition of this our classification. Then we can also have a similar notation as to what we had for a two class problem under a cost structure.

So, this is what, this is the probability of let me write that, this is the probability of an observation from π_i getting misclassified into π_k getting misclassified into the population k , that is π_k . Now, what is this equal to p_{ki} given i . So, this is going to be the integral over the region R_k , because we are putting into the k th population, where in the observation as such is coming from the i th population. So, it has got density $f_i(x) dx$.

So, we can also see what is p_{ii} equal to. So, this is equal to integral over the region R_i of $f_i(x) dx$. So, that you can write this as Ω the full space minus union of all other R_i 's, This i let me write this as k , this i is not equal to k . So, we are looking at the complementary region of R_i , that is Ω minus union of all other regions, regions

which are other than i of this quantity $f_i(x) dx$. So, this term, the first term would be integral over ω of $f_i(x) dx$, which is going to be equal to 1. And the next term is over union $i \neq k$ of this r_k regions $f_i(x) dx$. So, the first term here is going to be equal to 1 and the second term is summation over $k = 1$ to up to c with $k \neq i$, because we have taken out that R_i region from here. Of these term here, integrals $f_i(x) dx$. So, these are the quantities which would be required as such to define the expected cost of misclassification.

(Refer Slide Time: 36:19)

The conditional ECM of x from π_1 to π_2 or π_3 or ...
 π_c is

$$ECM(i) = c(2|i)P(2|i) + c(3|i)P(3|i) + \dots$$

$$\dots + c(c|i)P(c|i)$$

$$ECM(i) = \sum_{i=2}^c c(i|i)P(i|i)$$
 The ECM(i) is w.r.t. p_1

$$ECM = p_1 ECM(i) + p_2 ECM(2) + \dots + p_c ECM(c)$$

$$= p_1 \sum_{i=2}^c c(i|i)P(i|i) + p_2 \sum_{\substack{i=1 \\ i \neq 2}}^c c(i|i)P(i|i) + \dots$$

$$\dots + 1$$

Now, let us build up that expected cost of misclassification in the following way, that we first look at the conditional ECM, the conditional expected cost of misclassification of an observation x from π_1 say π_2 or π_3 or any of the other populations that is π_1 , π_2 , π_3 , π_c , any of these is going to be given by say ECM 1. So, this is the conditional expected cost of misclassification of an observation x , which is coming from π_1 into any other population other than the first population π_1 .

So, what is this going to be equal to, this is going to be equal to the probability, the cost say first off all suppose let us consider the case that a x is misclassified into π_2 given it is coming from π_1 . So, what is the cost of that? This is the cost that we are going to incur from the general terminology what we had said was p_{ki} to be **I am sorry** this c_{ki} the cost of misclassifications structure, c_{ki} is the cost of misclassifying an observation coming from π_i into π_k . And hence, here we are looking at an observation coming

from π_1 and we are putting it into π_2 . And hence, this is cost that we are going to incur and what is probability of misclassifying an observation coming from the first population into the second this is given by p_{21} , where p_{21} is given by the expression that we have written for a general k_i, k situation here.

So, c_{11} or p_{11} would in particular be integral over R_2 of $f_1(x) dx$. That is how, this term is defined, this plus suppose that observation is still from first. And then, it is classified into the third population. So, what will be having is c_{31} given 1. These are all mutually exclusive cases. So, that we will be having this as summation as the expected cost of misclassification, it is a conditional expected cost of misclassification given that it is from 1.

So, that this would be given by this term, this plus the last term would be in this summation will be c_{c1} given 1. So, this is a c th population. So, this is the misclassifying an observation from the first population into the c th population, this multiplied by the probability of misclassifying this. So, this basically is equal to summation c_i into p_i summation of i from 2 to up to c . So, this is the $E C M_1$ term, now this $E C M_1$ is with a probability p_1 . This is $E C M_1$. So, the expected cost of misclassification $E C M_1$ is with probability p_1 , because that is the prior probability for the first population.

So, if we have this conditional $E C M$ of x from π_1 getting misclassified into π_2 or π_3 or π_c to be given by the expression that we written here with a probability p_1 , because it is from the first population with a priory probability as p_1 . We will have the expected cost of misclassification to be given by the expected cost of misclassification is $E C M_1$ the conditional 1, that with probability p_1 . That times the expected cost of misclassification condition on the observation coming from the second population would similarly be $E C M_2$ with a probability p_2 . And thus we will have to look at all such conditional probabilities, condition cost of misclassification, expected cost of misclassifications, which is going to be given by this $E C M_c$ for the c th population.

So, we can write these terms without much of difficulty that this term would be equal to summation i equal to 2 to c . What we have written there, i given 1 into p_i given 1, this plus p_2 would be a term. Now, what would be the type of terms that would be there in this conditional $E C M$ of x from π_2 . That is going to be a similar sum, with the sum a little just write it as in the first term here c_i given 1 into p_i given 1. This summation i is

from 1 to up to c, wherein this c is not equal to 2, because we are looking at expected cost of misclassification, their conditional expected cost of misclassification of x from 2. And hence, this would be the expression when we talk about E C M to which is this expression.

Now in a similar manner, we can write the other terms. So, the last term would be equal to summation i equal to 1 to up to c, not c because c is going to be left out here. So, this is up to c minus 1, of the terms which are same as what we have in here.

(Refer Slide Time: 42:10)

$$\Rightarrow ECM = \sum_{i=1}^c p_i \sum_{\substack{k=1 \\ k \neq i}}^c c(k|i) P(x|i) \quad \text{--- (*)}$$

The classification partition (R_1, \dots, R_c) minimizing ECM as given in (*) is given by

Allocate x to π_k ($k=1(c)c$) for which

$$\sum_{\substack{i=1 \\ i \neq k}}^c p_i f_i(x) c(k|i)$$
 is smallest

(*)

$$ECM = \sum_{i=1}^c \sum_{\substack{k=1 \\ k \neq i}}^c \int_{R_k} p_i c(k|i) f_i(x) dx \quad \text{--- (**)}$$

So, we can write this E C M under the present multiclass situation. In the following compact form, that this E C M is equal to summation i equal to 1 to up to c, say that is p i times summation k equal to 1 to up to c, where k will not be equal to c of expressions which are c k given i into p k given i. That is basically, the term that what we have out here. So, the summation is over i for this terms here and then summation over k equal to 1 to c with k not equal to **sorry** k is not equal to i, because outer sum here is i. So, if we have p 1, then this some k is not equal to 1. So, the summation here for p equal to 1 star from k equal to 2, if we have i equal to 2 that is p 2, we will be looking at summation over k equal to 1 to c without including k equal to 2. And similarly, for the last term if we have here p c, then this particular some would run from k equal to 1 to up to c minus 1.

So, in terms of the partition that we are looking at the classification partition which is R_1, R_2, \dots, R_C minimizing (No audio from 43:45 to 43:53) **minimizing** the expected cost of misclassification, expected cost of misclassification as given in this star 1 here, as given in star 1 is given by. We are now looking at regions. So, that we will be allocating x to a particular population p_i , this there are c such possible populations, k equal to 1 to c for which we will be having the inner sum here, that i equal to 1 to c say with i not equal to k , because we were looking at this p_i for which we will be having this p_i into $f_i(x)$ c given i is smallest, why is that. So, because this particular term if you recall, that that p_i is in terms of this, and hence this expected cost of misclassification. I will just write this expression which would lead one to get to this particular term here, which is summation i equal to 1 to c . So, this is the that is star 1 expression, which is equal to i equal to 1 to c . I will take this p_i also inside and write this expression as k equal to 1 to up to c with k not equal to i .

So we will have a p_i let me also take it inside. So, we will have any integral here $\int_{R_k} f_i(x) dx$. So, rather than writing it in this form, I will also take this Constants inside the integral. I will write it as p_i into c given I , this into $f_i(x) dx$. So, this is an alternate expectation of star 1, in terms of breaking of this particular term here p_i , because this p_i here as we have noted earlier this is integral over the region R_k of $f_i(x) dx$. And hence this star 1 expression is going to be given by this.

(Refer Slide Time: 47:08)

i.e. R_k is the region of all $x \Rightarrow$
 $\sum_{i=1}^c p_i f_i(x) c(k|i)$ to be the smallest
 $i \neq k$
 \downarrow
 i.e. $\left(\sum_{i=1}^c p_i f_i(x) c(k|i) \right) < \sum_{i=1}^c p_i f_i(x) c(k|i)$
 $i \neq k$ $i \neq j$ $j = 1 \dots c$
 $j \neq k$
 Note: Under the equal cost setup, the above classification partition reduces to:
 Allocating x to π_k if $\sum_{i=1}^c p_i f_i(x)$ is the smallest
 $i \neq k$

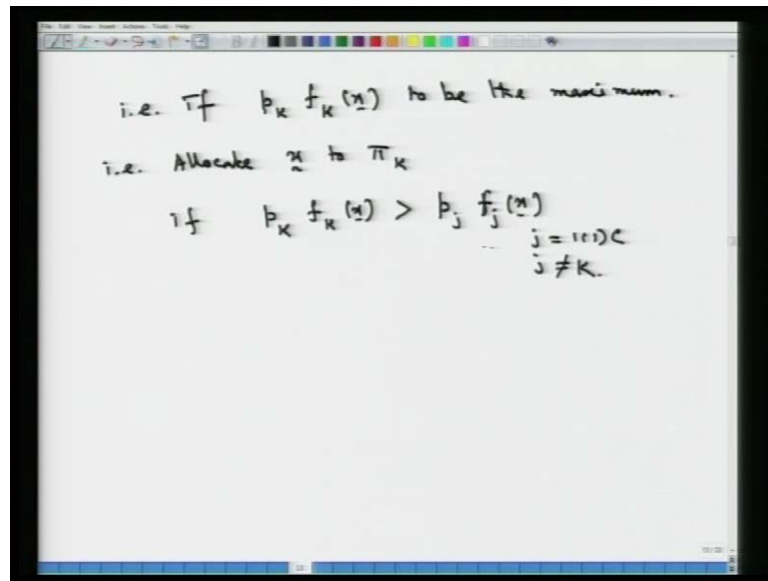
Now, if we have looking at the rule the partition R_1, R_2, \dots, R_C of the sample space which would lead us to minimizing this particular expression. We will look at choosing that particular region are k such that, we will have this term here which is given here to be the smallest among all possible case. And hence, this is what is leading us to the region r_k , that is r_k is the region of all x 's, such that we have got this summation that we have written in the previous slide, that i equal to 1 to c i not equal to k of $p_i f_i(x)$ into $c k i$ to be the smallest. That is the expression here which we have is going to be less than the expression for every other k other than the region, wherein we have got this k here. So, that is summation i equal to this is the left hand side only, this is i equal to 1 to up to c with i not equal to k of the expression $p_i f_i(x)$ given i . If that is less than summation over i equal to 1 to up to c i not equal to j . So, you take out c other term in this here and that would $p_i f_i(x) c k i$, this is for every j equal to 1 to up to c and j is not equal to k . So, we will be looking at all such sums, deleting a particular i , index i equal to 1 to c , and the region r_k is going to be the region of all such x 's for which this left hand side is less than the right hand side out here.

So, this is what is giving us the rule, which minimizes the expected cost of misclassification. And the allocation rule is what we have written here that to allocate x to $p_i k$, if we have got the **the** summation here to be the smallest among all possible such k terms

Now, we will note that this particular rule here under the equal cost setup. Under the equal cost setup, the above classification partition reduces to what if we look at this particular partition here. It is that we are going to have this to be smallest. Now if all the cost of misclassification are same that is under the equal cost setup. This cost of misclassification is not going to play any rule this can be taken outside.

So, if we have got equal cost setup, the this partition here is going to be reduced to allocating x to $p_i k$. If we have the summation i equal to 1 to c , i not equal to k $p_i f_i(x)$ is the smallest, smallest among all such c summations for which we will be having this i not equal to that particular chosen index there.

(Refer Slide Time: 51:05)



So, this is what we are going to get if we assume that we have got equal cost structure. Now, when is this going to be smallest if the term among all the c terms i equal to 1 to c , that is taken out is the largest. That is if we have got the term which is taken out, now what is the term which is taken out if we have looking at i equal to 1 to c , it is a term corresponding to k . That is if we have got p_k times $f_k(x)$ to be the maximum.

So, if we have got to be the maximum, we are essentially looking at this rule that we are going to allocate x to π_k . If we have got I said that, this $p_k f_k$ is a maximum that is if you have p_k times $f_k(x)$, this is greater than p_j times $f_j(x)$ for every j equal to 1 to up to c with the j not equal to k . So, we are not looking at that k th product on the right hand side, we are looking at all other c minus 1 products. And hence, we are looking at this all these c terms and then finding out for which of those c possible products. We have got the product to be maximum and then x going to be assigned to π_k corresponding to the population for which this is going to be the maximum.

Now, recall that when we were looking at the total probability of misclassification rule, what was a rule that we had got let us look back. And see what we had got earlier, when we were looking at the total probability of misclassification rule. We had said that the total probability of misclassification rule is one that is going to assign a x to π_i . If the probability product here like this is greater than this type of probability product for all

other products other than I . So, this essentially in terms of our notation, we have in the present case denoted this to be equal to p_i and we have denoted this to be f_i .

So, what we were trying to say is that, in total probability of misclassification classification minimizing rule, that x is going to be allotted to p_i . If $p_i f_i x$ is the maximum over all possible such i 's, i equal to 1 to see leaving out that particular i which is on the left hand side. Now if we look at the expected cost of misclassification minimizing rule, which we are derived just now, which was this? And that under equal cost setup, we was reduce to allocating x to p_k if the some like this is the smallest or in other words the terms it is left out in this particular sum here is the largest, that is was being allocated to p_k . If the $p_k f_k x$ is maximum among all possible such products. So, this would imply that under the equal cost, we can just remove this bracket actually. Under the equal cost setup this would imply that under the equal cost setup the ECM minimizing rule is same as that of the TPM minimizing rule.

So, this ECM minimizing rule is same as the TPM minimizing rule. This is what we expect also, because if we considering in the ECM setup, there is no nothing special about the ECM minimization or rather any special structure of the cost is not assumed. Then the rule which would minimize the ECM would naturally be same as the rule which would minimize the total probability of misclassification. And hence, we have also seen that, that is what is happening that if we take in the ECM minimizing rule the cost structures to be identical without having any special a preference about misclassification costs. Then the total probability of misclassification rule can be a derived from the ECM minimizing rule.

So, will stop at this particular point in this lecture. In the next lecture, we will first look at some examples of how to apply for a multiclass problem. This type of concepts of ECM minimizing rule or TPM minimizing rule, and then we will talk about some other important concepts in this classification problem. Thank you.