

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

Lecture - 9
Bayes Decision Theory – Binary Features

Good morning so, we are going to have, going to discuss today the Bayes decision theory for binary features. So far what we have discussed is, assuming that the feature vectors are distributed following some normal distribution of the form.

(Refer Slide Time: 00:41)

Bayes Decision Theory -
Binary Features

$$p(x|\omega_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp\left[-\frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i)\right]$$

$\Sigma_i = \sigma^2 I_{d \times d}$

$\omega_i \quad \omega_j$

μ_i μ_j

$P(\omega_i) = P(\omega_j)$

P of X give omega i is equal to so, this is the expression for a multivariate normal distribution. And we have seen that, in this expression where sigma i represents the covariance matrix, that for different conditions of covariance matrices, we can have different types of classifier.

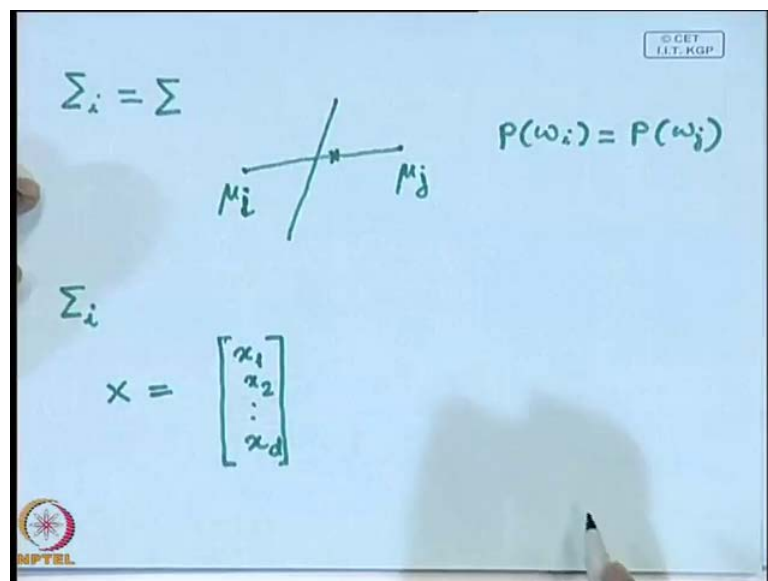
So, the first case we have discussed is, if the covariance matrix for every class i is of the form, sigma square into identity matrix I, where both this covariance matrix is of dimension d by d. Where our feature vectors are of dimension d and this identity matrix is also of dimension d by d. So, because this is identity matrix so, this simply says that covariance matrix is a diagonal matrix, where every diagonal element is of value sigma square. That means every component have the same variance or else, off diagonal

elements are equal to 0 so, because off diagonal elements are 0's so, the components, different components of the feature vectors are statistically independent.

So, in such cases we have seen that, the classifier is nothing but a linear classifier or when we talked about the discriminant function, the discriminant function for, function for individual classes, they are also linear functions. And because they are linear functions so, the classifier which employs this linear functions to classify and unknown feature vector, that is a linear machine. And in particular, if we want to find out the decision boundary two different classes say ω_i class, ω_i and the j th class, ω_j the decision boundary between these two different classes is a hyper plane. Which is orthogonal to the line joining μ_i and μ_j , where μ_i and μ_j is at the centers of the classes ω_i and ω_j .

So, effectively I have a situation something like this, that, if I have μ_i somewhere over here, which is the mean of the class ω_i and mean μ_j is somewhere over here. Then, the decision surface is orthogonal is a hyper plane, which is orthogonal to the line joining μ_i and μ_j . And if the apriori probabilities p of ω_i , is equal to p of ω_j then, this decision boundary or the decision surface becomes an orthogonal bisector of the line joining μ_i and μ_j . So, this was our simplest case.

(Refer Slide Time: 04:20)



In the second case, we have seen that if Σ_i is equal to Σ_j , that is every class or the samples belonging to every class have the same covariance matrix, but otherwise the

covariance matrix is arbitrary. Unlike in the first case where, covariance matrix has its specific form like this. In the second case the covariance matrices are arbitrary, but every class have the same covariance matrix. Which ideally means that, the points belonging to the same class or the points belonging to different classes, they are clustered in hyper ellipsoidal spaces of same shape and same size. And in such case we have seen that, the discriminating function that we get for different classes they are also linear.

So, in both the cases in the first case, as well in the second case the discriminating functions becomes linear so, the classifier is nothing but, a linear machine. However there is some difference between the decision surfaces, that we get between two different classes ω_i and ω_j . In the first case, the decision surface was orthogonal to the line joining μ_i and μ_j , in the second case the decision surface is not in general orthogonal to the line joining μ_i and μ_j .

So, here again if μ_i and μ_j so, this is μ_i and this is μ_j , which are the centers of the two classes ω_i and ω_j then, the decision surface between these two classes will be something like this. In the previous case, it was orthogonal to the line joining μ_i and μ_j , in this case in general it is not orthogonal. However the decision surface is a hyper plane or it represents the linear equation and here again, if the apriori probability is p of ω_i is same as, p of ω_j . Then, this decision surface though it is not orthogonal to the line joining μ_i and μ_j , but it will pass through the point which is midway between the points μ_i and μ_j .

And the third case we have said that, the covariance matrices of different classes are totally arbitrary. So, for i th class I will have one covariance matrix, for j th class I will have another covariance matrix. So, that effectively means that the clusters or the points, vectors belonging to different classes they are clustered into hyper ellipsoidal spaces. In this, in the first, second case the hyper ellipsoidal spaces were of same shape and same size. In this case, the points belonging to different classes will also form hyper ellipsoidal spaces, but these hyper ellipsoidal spaces may not have same shape or may not have same size. So, they will have different shapes as well as different sizes, but the points belonging to the same class they, form an hyper ellipsoidal spaces.

However, in all these three different cases that we have discussed, we have assumed that the feature vector x is continuous or the individual components of the feature vector x .

Because, our feature vector x is, nothing but a d dimensional vector having the components x_1, x_2 up to x_d . So, it is a d dimensional feature vector so, in all this three different cases our basic assumption was that, the feature vectors are continuous or in otherwise, individual components are also continuous. That effectively means, that if I consider a d dimensional feature space then, the feature vector can be represented by any point, is represented by any point within that d dimensional feature space. I do not have any specific set of points, from which the feature vectors are drawn.

However, in most of the practical applications and particularly in these days as we are walking with digital computers, all the data that we get are digital data. And the moment we get digital data, the vectors that we generate are no more continuous rather, they are discrete vectors or every component of the feature vector every x_i will have a discrete values. Discrete values means, it will assume one half a set of specific values so, instead of the continuous variable it becomes a discrete variable.

So, when it is a discrete variable, in that case in all our previous lectures wherever we have talked about integration, the integration is to be replaced by summation. And the summation has to be carried out, over the discrete space. So, we will take a specific case of this discrete feature vectors. So, let us consider a case where, so, we will have a feature vector x which will be discrete.

(Refer Slide Time: 10:11)

$x \rightarrow$ discrete

Two class problem $\rightarrow \omega_1$ & ω_2

Binary Feature Vectors

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \quad \begin{array}{l} x_i = 0/1 \\ p_i = P_r[x_i=1 | \omega_1] \\ q_i = P_r[x_i=1 | \omega_2] \end{array}$$

$p_i > q_i \Rightarrow x_i$ is more likely to have value 1 if $x \in \omega_1$

NPTEL

IIT KGP

And we consider a specific case of a two class problem and the feature vectors, we will assume to be binary feature vectors. Binary feature vector means, every component of the feature vector will assume binary value either 0 or 1. So, when I have this feature vector x , which is given by x_1, x_2 up to x_d . Again we assume that, we have d dimensional feature vectors. Every component of the feature vector x_1, x_2 or x_3 can assume either a value equal to 0 or a value equal to 1.

So, that means that, if a feature vector is taken from a particular class it will say whether a particular feature is present or it is absent. So, when I have this sort of binary feature vectors so, we will also assume the different components of this feature vectors are conditionally independent. To have something similar to statistical independence of different components inventing most feature vectors. So, here also you will assume that different components of the feature vectors are conditionally independent.

So, every feature value in this feature vector x_i , every feature value can have a value either 0 or 1. So, it will say whether the feature is present or the feature is absent and correspondingly, the probabilities will be something like this. So, every feature component will be represented by a probability value where, the probability is something like this, that I will represent p_i to represent the probability of the i th component. So, the p_i will be equal to probability that, component x_i is equal to 1 given that, the true state of nature or the true class is ω_1 . We are concentrating a two case problem.

So, we have classes ω_1 and ω_2 . So, p_i represents probability that x_i is equal to 1 given that, the true state of nature is ω_1 similarly, q_i is probability that, x_i the same component is equal to 1, given the true state of nature is ω_2 . So, given this type of probability measures, it simply means that if I have a situation that p_i is greater than q_i . In such cases, it simply means that i th component x_i is more likely to have value 1 if, x belongs to class ω_1 .

Because, it is the probability of assuming a value equal to 1, when the true state of the nature of the two classes ω_1 or when the true state of the nature of two classes ω_2 . So, if p_i is greater than q_i that simply means that, if the sample is taken from class ω_1 . Then, it is more likely that the i th component x_i will have value equal to 1 or x_i will have a value equal to 1, more frequently if the vector is taken from class

omega 1. And if it is taken from class omega 2 then, it is less frequent that the i th component x_i will have a value equal to 1.

So, when I have this kind of situation now, let us see that what are the cases in which these kind of feature vectors are more useful. Say for example, if we want to find out the health of a plant say, power plant, if I want to determine the health of a power plant. Then what, we normally do is there are number of sensors which are used to monitor different parameters of the plant. And after monitoring those different parameters, you decide whether the plant is or there is some danger in the plant.

And when you sensor, when you monitor the sensor outputs you just see that, whether the sensor output is above a threshold level or below a threshold level. So, if it is above a threshold level, we set a value equal to 1, if it is below a threshold level we set a value equal to 0. Let me take a more obvious example, when I go to the market to purchase oranges, you must have noticed that even in our tech market you usually get two types of oranges. One type of oranges which are produced at Nagpur and one type of orange which are coming from Darjeeling. Have you noticed any difference in appearance between these types of oranges.

Student: Color is different

Color is different, if it is from Darjeeling color is more attractive. It is really orangy color, it is more yellowish whereas, oranges from Nagpur they are more greenish and if it becomes quite old, it becomes more reddish. If you look at the surface texture, the oranges which are coming from Darjeeling, they are smooth whereas, the oranges which are taken from Nagpur, they are rough. So, if simply based on these two features I want to determine, I want to have an automated machine which will simply classify, tell me that whether these are Darjeeling orange or a Nagpur orange. So, it will try to take the decision based on the color in the simplest case or it will try to take the decision based on the feature, the texture feature.

So, if I keep some of the feature vectors like, whether the color is yellowish answer will be either yes or no, whether it is greenish either it will be, answer will be either yes or no whether it is smooth, answer will be either yes or no. However, when I get an orange from Darjeeling and I take an orange from Nagpur, there is no guarantee that all the oranges from Darjeeling, will always have smooth texture or will always have orangy or

yellowish color. Or if take oranges from Nagpur there is no guarantee that, I will always have greenish color, Nagpur even may produce some oranges which will have yellowish color or which will have smooth textures.

So, there is always a finite probability that an orange produced at Nagpur will have yellowish color. So, it is not necessary if that color to be yellowish, I put that as a binary feature, for all the oranges coming from Nagpur that binary value will always be equal to 0. It is not guaranteed, for some of them I may get values which are equal to 1, but that is less frequent than, the value equal to 1 when the oranges are taken from Darjeeling.

So, simply over here if that feature I put as the i th feature x_i then, p_i will be more frequently equal to 1, x_i will be more frequently equal to 1, if this ω_i is Darjeeling. And this will be less frequently equal to 1, if ω_2 is from Nagpur so, coming over here if p_i is greater than q_i . So, it simply explains this particular situation, that is for oranges coming from Darjeeling, I will have more number of oranges having yellowish color than, the number of oranges I get with yellowish color from the Nagpur oranges. So, this is a typical situation like this.

And what does conditional independence mean, the texture and the color they are independent. I mean if the texture is rough, that does not necessarily mean that the color will be greenish or if the texture is smooth, that does not necessarily mean that the color will be yellowish. So, coming to the plant if I take two features say, I want to monitor the health of a boiler. If I take two features, one is pressure inside the boiler and temperature inside the boiler.

They are not independent because, if temperature increases the pressure will increase. I mean just from the basic laws of physics so, they are not independent they are dependent. So, when I try to select the features I should select in such a way that the features are independent because, that solves many of the mathematical problems, I do not have to look for complicated mathematics. So, if I assume that the features are conditionally independent.

(Refer Slide Time: 21:17)

$$P(x|\omega_1) = \prod_{i=1}^d p_i^{x_i} (1-p_i)^{1-x_i}$$
$$P(x|\omega_2) = \prod_{i=1}^d q_i^{x_i} (1-q_i)^{1-x_i}$$

Likelihood Ratio.

$$\frac{P(x|\omega_1)}{P(x|\omega_2)} = \prod_{i=1}^d \left(\frac{p_i}{q_i}\right)^{x_i} \left(\frac{1-p_i}{1-q_i}\right)^{1-x_i}$$

Then, the same probability function I can write as p of x given ω_1 , where x is the feature vector, which has d number of binary valued feature components. So, this p of x given ω_1 , I can simply write as p_i , p_i is the probability that x_i equal to 1, given the true state of nature is ω_1 . But it also has a finite probability that, it may have a value equal to 0 so, I have to concentrate on that as well.

So, it is p_i to the power x_i into $1 - p_i$ to the power $1 - x_i$ obviously, if x_i equal to 1, $1 - x_i$ equal to 0, if x_i equal to 0, $1 - x_i$ equal to 1. And because, the components are conditionally independent so, the overall probability will be product of independent probability values. So, this I have to take for i is equal to 1 to d , as I have d number of components so, this is the class conditional probability of a feature vector x , if the state of nature is ω_1 .

So, in the same manner I can write p of x given ω_2 , that is class conditional probability if the feature vector belongs to class ω_2 is nothing but. Now, the probability of x_i equal to 1, then the true state of nature is ω_2 is q_i . So, I will have q_i to the power x_i into $1 - q_i$ to the power $1 - x_i$, take the product from i equal to 1 to d .

Student: Sir where d is dimension of the.

D is the dimension of the feature vector so, I have d number of components in the feature vector. So, these are the two class conditional probability values now, from here I can find out, what is called likelihood ratio. So, the likelihood ratio is given by p of x given ω_1 upon p of x given ω_2 which is, nothing but if I simply multiply these two, it becomes p_i upon q_i to the power x_i into $1 - p_i$ upon $1 - q_i$ to the power $1 - x_i$. Take the product where, i varying from one to d so, this is what is the likelihood ratio right.

Student: Sir x_i could be a vector or X_i is a single component?

x_i is the single component, it is a scalar, x_i is the single component, the capital x is the vector, having d number of components so, but the value of i will vary from 1 to d .

(Refer Slide Time: 25:56)

The image shows a whiteboard with handwritten mathematical equations. At the top right, there is a small box containing the text '© CET I.I.T. KGP'. The main equation is:

$$g(x) = \ln \frac{P(x|\omega_1)}{P(x|\omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

Below this, an arrow points to a more detailed expression:

$$\Rightarrow g(x) = \sum_{i=1}^d \left[x_i \ln \frac{p_i}{q_i} + (1-x_i) \ln \frac{1-p_i}{1-q_i} \right] + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

At the bottom left of the whiteboard, there is a logo for 'MPTEL'.

Now, you know that the decision surface or our decision function between the two classes given, two classes is given by $g(x)$ is equal to \log of $p(x)$, given ω_1 upon $p(x)$, given ω_2 plus \log of $P(\omega_1)$ upon $P(\omega_2)$. This we derived earlier now, if we put this $p(x)$ given ω_1 upon $p(x)$ given ω_2 , which is equal to this expression, if I put this expression into this function.

So, it will give us because we are taking logarithm so, the product term over here will be converted to sum of logarithmic functions. So, what I simply get is $g(x)$ is equal to summation of x_i because, here x_i was a power. So, this will become a product of the, in

the logarithmic term so, it is x_i then, \log of p_i upon q_i plus $1 - x_i$ \log of $1 - p_i$ upon $1 - q_i$, where i will vary from 1 to d plus \log of p of ω_1 .

Where this p of ω_1 is the apriori probability upon p of ω_2 . So, this is the decision function and for a two category case we have already seen that, if g of x becomes greater than 0 then, our decision was that x belongs to class ω_1 . If g of x becomes less than 0 then, our decision was that x belongs to class ω_2 . If g of x is equal to 0, that actually tells us that what is the decision boundary between the classes ω_1 and ω_2 .

Now, if you notice that this equation is a linear equation, isn't it? Because this simply says, the linear combinations of different components x_i of the feature vector x , I do not have any x_i term, x_i square term or x_i cube term. So, the equation is a linear equation and this linear equation can simply be written in the form, if I just rearrange this particular linear equation I can write it in the form.

(Refer Slide Time: 29:31)

$$g(x) = \sum_{i=1}^d w_i x_i + w_0$$

$$w_i = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)} \quad ; \quad i=1 \dots d$$

$$w_0 = \sum_{i=1}^d \ln \frac{1-p_i}{1-q_i} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

$$g(x) > 0 \Rightarrow x \in \omega_1$$

$$g(x) < 0 \Rightarrow x \in \omega_2$$

$g(x)$ is equal to sum of $w_i x_i$ plus w_0 , do not confuse w with ω . So, it becomes sum of $w_i x_i$ plus w_0 , where i varies from 1 to d . Because, I have d number of components in the feature vector. So, you find that this is, nothing but a linear combination of the different components, of the feature vector x_i plus a threshold term, which is w_0 . And this w_i is for different values of i , this represents a wide vector and this is, nothing but a dot product or inner product of the wide vector with the feature vector.

So, when I write it like this, over here this w_i is simply $\log \frac{p_i}{1 - q_i}$ that is quite obvious. Because, you find that this becomes $x_i \log \frac{p_i}{1 - q_i}$ into $\log \frac{p_i}{1 - q_i}$ plus $x_i \log \frac{1 - q_i}{p_i}$. So, this $\log \frac{p_i}{1 - q_i}$ term that goes to the numerator and $\log \frac{1 - q_i}{p_i}$ term that comes to the denominator.

So, it simply becomes $\log \frac{p_i}{1 - q_i}$ upon $\log \frac{p_i}{1 - q_i}$ upon q_i into $1 - p_i$ where, i varies from 1 to d . Because, I have d number of components in the wide vector as well and w naught, that is the threshold is given by $\log \frac{1 - p_i}{q_i}$ upon $1 - p_i$, sum of this for i is equal to 1 to d plus $\log \frac{p}{\omega_1}$ upon p of ω_2 .

Now if you analyze this, as I said that our decision will be that if, g of x is greater than 0 then, we decide that x belongs to class ω_1 . If g of x is less than 0 then, we decide that x belongs to class ω_2 , if this is equal to 0 that is a boundary case. So, if you analyze this expression or what do we get, what does this different components of w , that is w_i , that effectively tell us. You find that x_i , that is the i th component of the feature vector x is the binary value feature vector. It can have a value equal to 0, it can have a value equal to 1.

Now if it has a value equal to 1 then, the contribution of the term w_i into x_i for that particular component x_i , to this function g_x is, nothing but equal to the magnitude of w_i . Because, x_i is equal to 1 so, it is nothing but equal to the value, the magnitude of that particular component of w_i . So, effectively this w_i , magnitude of it simply tells you that what is the importance or what is the relevance of the component x_i in decision making that, whether the sample will belong to class ω_1 or the sample will belong to class ω_2 . If value of w_i is very large then, x_i has more weightage to decide about the class, if value of w_i is small then x_i has less weightage to decide about the class.

And in other case, if p_i is equal to q_i that is value of x_i to be equal to 1 is more likely, is same equally likely, even if the x belongs to class ω_1 or the feature vector x belong to class ω_2 . So, p_i is equal to q_i , p_i as we said it is the probability that x_i will be equal to 1, if the true state of nature is ω_1 . And q_i is the probability that, x_i will be equal to 1 if the true state of nature is ω_2 . So, if p_i is equal to q_i that simply indicates that, whether x belongs to class ω_1 or x belongs to ω_2 , x_i is equally

likely to have value equal to 1. So, that simply means that x_i has no relevance in deciding the class.

So if i has, if x_i has no relevance in deciding the class then, why should the corresponding vector w_i be there. I can make w_i equal to 0, without hampering my decision. So, if you come to this w_i , the expression for this w_i you have find that if p_i is equal to q_i then, this expression $p_i \ln \frac{1-p_i}{1-q_i} + q_i \ln \frac{1-q_i}{1-p_i}$. This expression will be equal to 0 so, the corresponding w_i is equal to 0. And that is quite obvious because, if p_i is equal to q_i then x_i , the particular feature vector x_i has no relevance in deciding the class of the feature vector x .

On the other hand if p_i is greater than q_i , if p_i is greater than q_i then, having value of x_i equal to 1 should tell me that, the sample is more likely to belong to class ω_1 than, to belong to class ω_2 . Whereas, if p_i is less than q_i then, the sample x_i equal to 1 tells me that it is more likely to belong to class ω_2 than, its likelihood to belong to class ω_1 .

So, again coming to this particular case, if p_i is greater than q_i , if p_i is greater than q_i then obviously $1 - q_i$ will be greater than $1 - p_i$. So, in this expression the numerator becomes larger than the denominator so, this value is greater than 1. And when this value is greater than 1, value of w_i is positive, if value of w_i is positive what happens to my g_x .

Student: Positive.

Not necessary, it depends on other values of i as well. So, effectively I can say that if x_i is equal to 1, that particular component then, this component x_i gives a vote of value w_i to g_x , to decide that this feature vector x belongs to class ω_1 . So, it is the, it gives the vote equal to the corresponding w_i in favor of class ω_1 .

On the other hand if p_i is less than q_i , p_i is less than q_i so, $1 - q_i$ will be less than $1 - p_i$ so, numerator becomes less than the denominator. So, this term is a fraction which is less than 1, log of this will be negative that means, w_i in that case is negative. If w_i is negative that means, the corresponding x_i into w_i is trying to make g_x less than 0, trying to make, whether it will be less than 0 or not, that depends upon other w_i into x_i term. But what x_i is trying to do in this case, it is trying to give a vote equal to the

modulus of w_i , in favor of class ω_2 . Because, it is subtracting so, it is giving a vote which is equal to modulus of w_i in favor of class ω_2 .

So if p_i is greater than q_i , the component x_i gives a vote equal to w_i in favor of class ω_1 , if p_i is less than q_i then component x_i , gives a vote equal to modulus of w_i in favor of class ω_2 . That means this component is trying to push the decision surface of the decision boundary, either towards ω_1 or towards ω_2 .

(Refer Slide Time: 40:15)

Example
3-dimensional feature vectors
Two class problem $\rightarrow w_1$ & w_2
 $P(w_1) = P(w_2) = 0.5$
 $p_i = 0.8$ and $q_i = 0.5$ for $i=1,2,3$
 $w_i = \ln \frac{0.8(1-0.5)}{0.5(1-0.8)} = 1.3863$
 $w_0 = \sum_{i=1}^3 \ln \left(\frac{1-0.8}{1-0.5} \right) = -1.2$

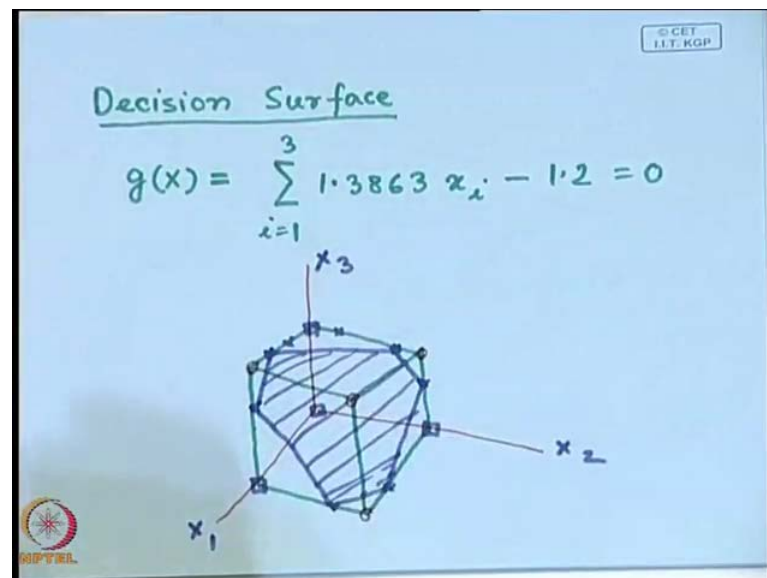
Let us take an example so, let us consider a three dimensional space or three dimensional feature vectors. And let us so, and let us consider a two class problem that is, I have classes ω_1 and ω_2 . So, these are the two classes I have and every feature vector is a three dimensional feature vector where, every individual component can be either 0 or 1.

And let us also assume that, the apriori probabilities p of ω_1 is same as p of ω_2 , which is equal to 0.5. And let us assume that value of p_i is equal to 0.8 and q_i is equal to say, 0.5 for all values of i , that is for i varying from 1, 2 and 3. So, it says that every component of the feature vector has a probability of being equal to 1 is 0.8 if the feature vector is taken from class ω_1 . And every component has a probability of 0.5 of being equal to 1, if the feature vector is taken from class ω_2 .

So given this p_i and q_i values a probability values, I can compute the corresponding wide vectors. So, simply from this expression that, w_i is equal to \log of p_i into $1 - q_i$ upon q_i into $1 - p_i$ for different values of i , I get different components of the wide vector w_i . So, I get w_i is equal to \log of p_i , that is 0.8 into $1 - q_i$, that is 0.5 upon q_i that is 0.5 into $1 - p_i$ that is 0.8. And if I compute this, this becomes a value 0.3863 and the threshold w_{naught} , which is given by this expression, w_{naught} is equal to \log of $1 - p_i$ upon $1 - q_i$, take the summation for i is equal to 1 to d plus \log of p_{ω_1} upon p_{ω_2} .

Now over here, in this case p of ω_1 and p of ω_2 they are equal to, they are equal and both equal to 0.5. So, this last term \log of p_{ω_1} upon p_{ω_2} that will be equal to 0. So, what I have to compute is simply this term, w_{naught} is equal to \log of $1 - p_i$ upon $2 - q_i$ take the summation over, i is equal to 1 to d . So, this w_{naught} will be simply \log of $1 - p_i$ that is 0.8 upon $1 - q_i$ that is 0.5. This component take the summation for i is equal to 1 to 3 and we will find that, this value will be something like, it will have a value something like this.

(Refer Slide Time: 45:30)



So given this, the decision surface between the classes ω_1 and ω_2 . You verify this values, whether you are truly getting this values or not. So, assuming this our decision surface $g(x)$ will be given by sum of 1.3863 x_i , where this i varies from 1 to 3

because, I have three number of components. So, this is nothing but $1.3863 x_1$ plus $1.383 x_2$ plus $1.383 x_3$ minus 1.2 this is equal to 0 .

So, if I try to plot this decision surface the, in three d one point you might have noticed that, because our features are binary features. Every feature component can assume a value either 0 or equal to 1 so, every feature vector will be represented by a vertex of a hypercube in the d dimensional space. It can be either $0 0 0$ or $0 0 1$ or $0 1 0$ or $1 0 0$ and so on. So, every feature vector will be represented by a vertex, in the d dimensional space of a hypercube.

So, if I try to plot this decision surface, the decision surface will be something like this. So, let us take a cube in this three dimensional space and if you plot the surface, the surface will come out to be something like this. So, this is our decision surface so, you find that, this decision surface says that, these are the points which lie on one side of the hyper plane. And these are the points which lie on the other side of the hyper plane so, it simply says that, if at least two vectors of the feature, if at least two components of the feature vector are equal to 1 . Then, the point is classified to class ω_1 .

Student: Sir it should be at least one because, equation says it will be at least one ((Refer time: 49:34)).

Which one.

The equation from this decision surface that you ((Refer time: 49:40)).

This one.

Student: Yes sir.

Student: If at least one is ((Refer time: 49:52)).

I mean it is the other way, these are the points which are put to class ω_1 , ω_2 and the other side is put to class ω_1 . So, it says if at least one of them is equal to 1 then, the decision will be in favor of class ω_1 . Otherwise, the decision will be in favor of ω_2 .

So, we find that I again get a simple hyper plane in three dimension, it is just a plane, which is boundary between the two classes ω_1 and ω_2 . So, if the probability

values are different, if different π s have different other values then, the position as well as orientation of this plane may be different. But effectively what it does is, the d dimensional space is broke into two halves, one half will be given to class ω_1 , the other half will be given to class ω_2 . The nature will be a bit more complicated when, the number of classes are more than 2 because, then we have to think of more than one decision surfaces and how they combine.

Student: Sir, will axis represents x_1, x_2 ?

Axis represents x_1, x_2, x_3 . So, let us stop here today.

Thank you.