

Deep Learning
Prof. Prabir Kumar Biswas
Department of Electronics and Electrical Communication Engineering
Indian Institute of Technology, Kharagpur

Lecture - 06
Discriminant Function - I

Hello welcome to the NPTEL Online Certification Course on Deep Learning.

(Refer Slide Time: 00:35)



In the previous class we have talked about topics like different types of classifiers like Bayes minimum error classifier and Bayes minimum risk classifier. And, we have also seen that Bayes minimum risk classifier under a specific case of 0 1 loss that is when the loss function for a correct decision is taken to be 0 and the loss function for an incorrect decision is taken to be 1 under that situation Bayes minimum risk classifier and Bayes minimum error classifier, they are identical.

So, in today's lecture we will talk about we will start from those Bayes classifiers and then we will move on to what is known as discriminant function. And, then using the Discriminant Function we will also try to derive and we will also try to demonstrate the decision boundary between different classes.

(Refer Slide Time: 01:43)

Discriminant Function

$$P(\omega_i | X) > P(\omega_j | X)$$
$$P(\omega_i | X) = P(X | \omega_i) \cdot P(\omega_i)$$

So, when we talk about a discriminant function you remember from the previous class in case of Bayes minimum risk classifier or Bayes minimum error classifier what we said is that for Bayes minimum error classifier if P of say ω_i given X is greater than P of ω_j given X where ω_i and ω_j are two different classes and X is the unknown input vector. In that case we classify the X to this class ω_i and if I expand this P of ω_i given X is nothing, but P of X given ω_i multiplied by the a priori probability P of ω_i , where P of X given ω_i is what is known as the class conditional probability density.

P of ω_i is the a priori probability and P of ω_i given X is the a posteriori probability based on which we make the decision that whether this unknown vector X should be classified to class ω_i or it should be classified to class ω_j . So, obviously if P of ω_i given X is greater than P of ω_j given x , then it is more likely or more probable that your unknown feature vector belongs to class ω_i . And, this is what we had derived in our previous lectures using Bayes minimum error classification rule.

(Refer Slide Time: 03:43)

Discriminant Function

$$R(\alpha_i | x) = \sum_{\forall \omega_j} \lambda(\alpha_i | \omega_j) \cdot P(\omega_j | x)$$
$$\lambda(\alpha_i | \omega_j) = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases}$$

$-R(\alpha_i | x)$

And then Bayes minimum risk classification what we had said is for an unknown feature vector X , if we take an action α_i , then the risk involved is given by R of α_i given X . And, which we said that this is nothing, but $\lambda(\alpha_i | \omega_j)$ into P of ω_j given x . You take the summation over this for all the classes ω_j .

So, for every action α_i if I have say c number of actions, so i varies from 1 to c . So, for every such action I have to compute this risk function. And, for whichever action the risk the estimated risk R of α_i given X is minimum, I have to take that corresponding action. And, as we said that under a specific case when this $\lambda(\alpha_i | \omega_j)$ is equal to 0 for i is equal to j and if I take this equal to 1 whenever i is not equal to j . That means, for every correct decision the loss incurred is 0 and for every incorrect decision the loss incurred is 1 under that situation we had shown in the previous class that Bayes minimum risk classifier and Bayes minimum error classifier, they turn out to be identical.

Now, starting from here we can define something called discriminant function because every time you find that whether I go for Bayes minimum error classification or Bayes minimum risk classification, in case of Bayes minimum error classification for every class I am computing P of ω_i given x , where i varies from 1 to c where c is the number of classes that I have and for whichever i P of ω_i given X turns out to be maximum i classify X to that corresponding class.

Similarly, in case of Bayes minimum risk classification for every class I compute R of α_i given X and then for whichever class R of α_i given X is minimum that is for whichever class for whichever action the risk involved is minimum, I am taking the corresponding action or I am classifying X to that to that corresponding class. And, I can say that in each of this case I am taking an action based on certain maximum criteria that is in case of Bayes minimum error classification for whichever class the a posterior ability is maximum, I am taking that corresponding action or classifying X to that corresponding class.

Similarly, for Bayes minimum risk classification for whichever class R of α_i given x , so for whichever value of i are all R of α_i given X turns out to be minimum or in other case I can say that instead of considering R of α_i given X , I will consider minus R of α_i given X . So, if R of α_i given X is minimum, then obviously R of α_i given X minus R of α_i given X will be maximum.

So, for whichever action this negative of the risk of value turns out to be maximum, I am taking that corresponding decision or I am classifying X to that corresponding class.

(Refer Slide Time: 08:00)

Discriminant Function

$$\left. \begin{array}{l} g_1(x) \\ g_2(x) \\ \vdots \\ g_i(x) \\ \vdots \\ g_c(x) \end{array} \right\} \underline{\underline{\text{Max}}} \quad X \rightarrow \omega_i$$

$g_i(x)$ ω_i

So, or in other words I can say that I can define a function say $g_i X$ for class say ω_i . So, here X is the unknown feature vector and for every class I every class ω_i , I am computing a function $g_i X$ and for whichever i this $g_i X$ turns out to be maximum I take decision in favor of that particular class or that particular ω_i .

So, what I am doing is for an unknown feature vector X I will compute g_1 of X , I will compute g_2 of X , I will compute g_i of X and if there are c number of classes I will compute g_c of X . And, then I will try to find out that out of all these functional values whichever is maximum. So, I take maximum of all of this and for whichever i this turns out to be maximum I classify X to that corresponding class ω_i . So, this is a function that $g_i(X)$ I want to design for every class ω_i . So, what are the possible options that I can have $g_i(X)$?

(Refer Slide Time: 09:29)

Discriminant Function

$$g_i(x) = P(\omega_i | x)$$

$$= P(x | \omega_i) \cdot P(\omega_i)$$

$$g_i(x) = -R(\alpha_i | x)$$

One of the options is obviously I can have $g_i(X)$ to be equal to P of ω_i given X that is straightforward which is nothing, but P of X given ω_i into a priori probability P of ω_i . So, this is a straightforward definition of $g_i(X)$ or I can also say that I will use $g_i(X)$ to be minus R of α_i given x .

So, here also if $g_i(X)$ is maximum I take that corresponding decision here also if $g_i(X)$ is maximum I take that corresponding decision. So, out of these two options possible options because there might be other options as well. Out of these two we will try to explore this that is based on Bayes minimum error classification rule.

(Refer Slide Time: 10:37)

Discriminant Function

$$g_i(x) = p(\omega_i|x) \ln$$
$$f(p(\omega_i|x))$$
$$p(\omega_i|x) > p(\omega_j|x)$$
$$\Rightarrow f(p(\omega_i|x)) > f(p(\omega_j|x))$$

So, I will assume that we will use this discriminant function $g_i(x)$ where this discriminant function is defined as P of ω_i given X or it is also possible that instead of P of ω_i given X if I use a function of P of ω_i given X where this function f has to be a monotonically increasing function. That means, if P of ω_i given X is greater than P of ω_j given X , this should imply the f of P of ω_i given X should be greater than f of P of ω_j given x . Let me rewrite this. So, this implies f of P of ω_i given X has to be greater than f of P of ω_j given x .

So, that is f has to be a monotonically increasing function, then this form that is $g_i(x)$ as f of P of ω_i given X that can also be used as a discriminative function because whenever $g_i(x)$ is maximum, the f of P of ω_i given X will also be maximum. So, given this you will find that one very convenient function that can be used is logarithmic function. So, at or I can use natural logarithm \ln . What is the advantage?

(Refer Slide Time: 12:39)

Discriminant Function under Multivariate Normal Distribution
Discriminant Function

$$p(w_i | x) = p(x | w_i) \cdot P(w_i)$$

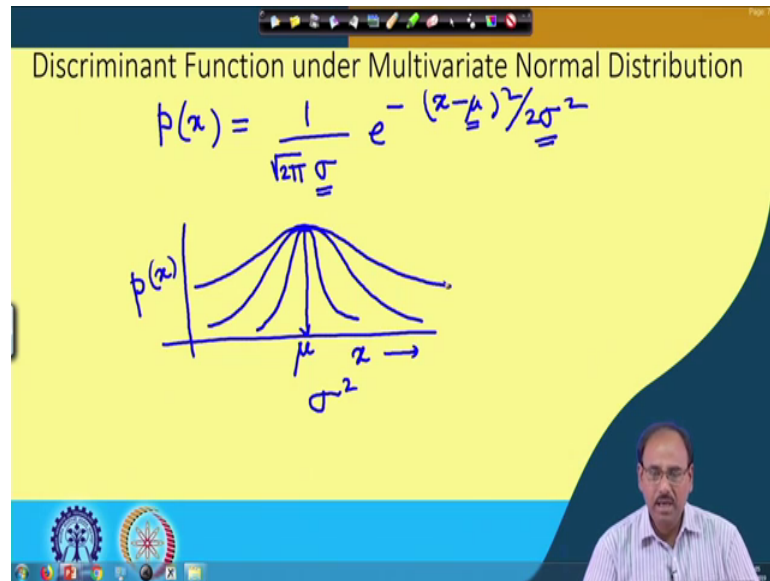
f : ln

$$\ln p(w_i | x) = \ln p(x | w_i) + \ln P(w_i)$$
$$g_i(x) = \ln p(x | w_i) + \ln P(w_i)$$

The advantage is because P of ω_i given X is nothing, but P of X given ω_i which is the class conditional probability density function that can be estimated experimentally into P of ω_i which is the a priori probability of class ω_i that is also pre-computed. And, now if this f the function I used as the logarithmic function \ln , then the advantage that I get is $\ln P$ of ω_i given X that turns out to be $\ln P$ of X given ω_i plus log of P of ω_i . So, this multiplication state way is converted to an addition operation and which is very very advantageous in many computational purposes.

So, I will use this particular form that $g_i(X)$ is nothing, but log of P of X given ω_i plus log of a priori probability P of ω_i . So, this is the discriminant function form that we will use in the remaining part of this lecture, right. So, given this let us try to see that if I assume a particular form or a particular distribution function probability density function which usually we use as a normal probability density function, then what form of expression of the discriminant function that we get or what form of the decision boundary between classes that we get?

(Refer Slide Time: 14:54)

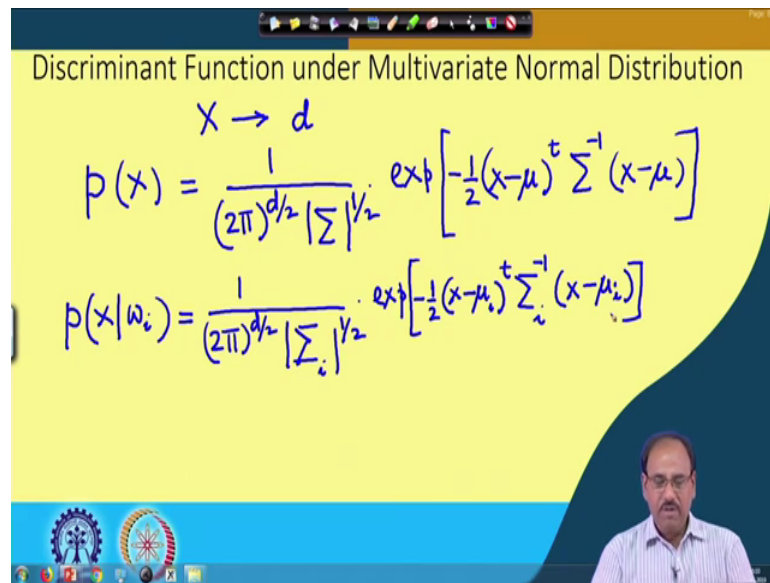


So, let us talk about this discriminant function under multivariate normal distribution. So, we all know that the normal distribution if I have a single variable say $p(x)$ is given by 1 over square root of $2\pi\sigma$ e to the power minus $(x-\mu)^2$ upon $2\sigma^2$. This is the normal density in case of a single variable x or scalar variable x where this σ is nothing, but standard deviation and σ^2 is the variance and μ is the mean.

And, you all know that the typical form of this is if I plot x and $p(x)$, the typical form is like this where this is what is your the mean of x that is μ and the value of this envelope depends upon the value of σ or σ^2 .

So, if σ is lowered, the σ^2 is low then I will have a distribution something like this. If the σ^2 is high, then the distribution will be flat of this form. So, this is the form that I get when x is a single variable or it is a scalar variable.

(Refer Slide Time: 16:41)



Discriminant Function under Multivariate Normal Distribution

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

But in our case since we are talking about feature vector which describes an object and the feature vector consists of multiple number of features where every feature captured some property or some attribute of the object. So, those features may be computed from the shape of the object, they might be computed from the color of the object. They might be computed from the intensity of the object, they might be computed from the texture of the object and various such different properties are put together in the form of a vector or a feature vector. So, the type of distribution that is important in our case is not a single variant distribution, but it is a multivariate distribution.

So, in our case I have a feature vector X and let me assume that the dimension of the feature vector is d . So, it is a d dimensional feature vector right; d is the number of components or the number of features which are packed into this feature vector X . So, given this the multivariate property density is now given by P of X which is 1 over 2π to the power d by 2 , then instead of variants now I have multiple variables. So, what I have is a covariance matrix.

So, σ is the covariance matrix. You take the determinant of that and square root of the determinant into exponential minus half X minus μ transpose σ inverse into X minus μ that is what is the normal distribution form of normal distribution in case of multivariate or in case of vectors.

Now, what we are interested in is the expression that we have that contains X given ω_i . That is the class conditional probability density and we said earlier that we get this class conditional probability density by taking the feature vectors X from class ω_i . So, when I take feature vectors X from class ω_i , so for those feature vectors the mean that I will get is dependent and I will represent that by μ_i . Similarly the covariance Σ , the covariance matrix Σ that I compute will also be on for that particular class ω_i . So, I will also represent this as covariance matrix Σ_i .

So, what I will do is I will put this class conditional probability density function, express it in the form $\frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}}$. Now, this Σ actually becomes Σ_i because this is for class ω_i square root of that into exponent minus half $(X - \mu_i)^T \Sigma_i^{-1} (X - \mu_i)$. Now, this μ becomes μ_i . It is for i th class transpose Σ becomes Σ_i . So, it is Σ_i inverse into $X - \mu_i$. So, this is my multivariate probability density function where, Σ_i is the covariance matrix computed over all the feature vectors which we call as feature vectors because using those vectors, I am computing Σ_i and μ_i .

So, it is the covariance matrix computed using those feature vectors taken from class ω_i μ_i is the mean of those feature vectors taken from class ω_i .

(Refer Slide Time: 21:07)

Discriminant Function under Multivariate Normal Distribution

$$g_i(x) = \ln P(w_i|x) + \ln P(w_i)$$

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

$$g_i(x) = \cancel{\frac{d}{2} \ln 2\pi} - \frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) + \ln P(w_i)$$

Now, given this the way we have defined $g_i(X)$ is equal to \log of $P(\omega_i|X)$ plus \log of a priori probability P of ω_i . Now you find that P of ω_i given X is

nothing, but 1 over 2π to the power d by 2 , then σ_i square root of this exponential minus half X minus μ_i transpose. So, this is what is P of ω_i given X .

So, once I use this logarithm, then my g_i X it simply becomes minus d by 2 log of 2π minus half log of σ_i minus half of X minus μ_i transpose σ_i inverse X minus μ_i plus of course I have this log of P of ω_i . From here you will find that d by 2 log of 2π this particular term is independent of the class because there is no term like subscript i over here.

So, this minus d by 2 log of 2π , this does not differentiate between an i th class and j th class. So, easily I can conveniently ignore this particular term minus d by 2 log of 2π .

(Refer Slide Time: 23:47)

Discriminant Function under Multivariate Normal Distribution

$$g_i(x) = -\frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) + \ln P(\omega_i)$$

$\Sigma_i = \sigma^2 I$ — (1) ✓
 $\Sigma_i = \Sigma$ — (2)
 Σ_i — (3)

So, that simplifies my g_i X as minus half log of P σ_i minus half X minus μ_i transpose σ_i inverse X minus μ_i plus log of a priori probability P of ω_i . Now, here I can have different cases; say for example this covariance matrix σ_i in a particular case in a specific case if all the components the components of the feature vector X , they are statistically independent. Then covariance matrix that we get will be a diagonal matrix and if every component has same variance, then this covariance matrix σ_i will be of the form σ squared I .

So, what I am assuming here that for all the classes the feature vectors that you obtain the components of the feature vectors are statistically independent. So, that means if I try

to compute the variance involving say i th component and j th component because they are statistically independent. So, that variance will be equal to 0 which leads to the covariance matrix to be a diagonal matrix, where only I will have diagonal elements to be non-zero and all the off diagonal elements will be 0. And, then again if I assume that all those components, for all those components the variance is same in that case all the diagonal elements which are non-zero, they will be equal.

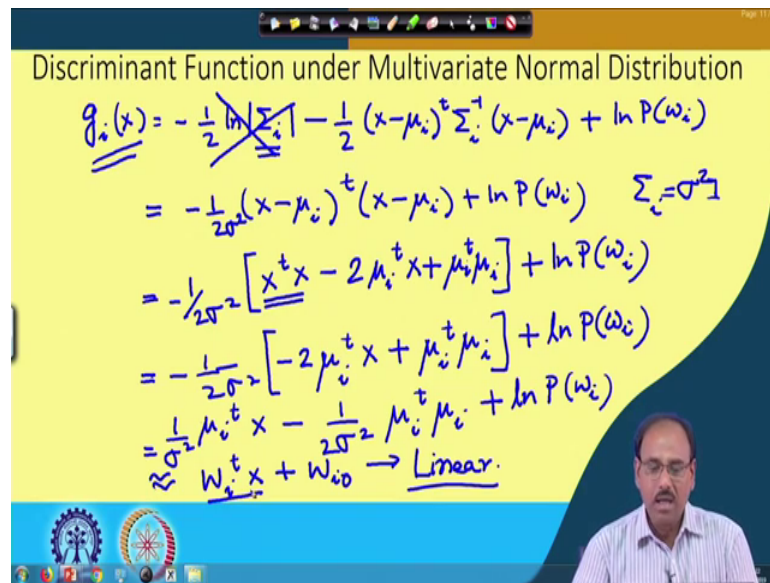
So, that ultimately leads to the covariance matrix to be of the form the sigma square I and I am assuming this to be same for all the classes. That means, for every sigma I , I have this covariance matrix for every omega i , the covariance matrix sigma i is of the form sigma square I . So, the sigma squared is again same for all the features across the classes. So, that is one of the simplified assumption that I can make the other assumption that can be used is where sigma i is of the form sigma.

So, in this case it is not necessary that the different components of the feature vector will be statistically independent, not even necessary that every component will have the same variance, but what I am assuming is that whatever is the covariance matrix, the same covariance matrix is valid for all the classes. So, this is a simplified condition, condition 2 and the third one where I have the most general case that every class will have its own covariance matrix that is the covariance matrix of one class need not be same as covariance matrix of other classes. So, that is the most general case which is case 3.

So, initially I will try to see that how this discriminant function look like when I assume the first case that is covariance matrix of every class is of the form sigma square I .

(Refer Slide Time: 27:58)

Discriminant Function under Multivariate Normal Distribution

$$\begin{aligned}
 g_i(x) &= -\frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x - \mu_i)^t \Sigma_i^{-1} (x - \mu_i) + \ln P(\omega_i) \\
 &= -\frac{1}{2\sigma^2} (x - \mu_i)^t (x - \mu_i) + \ln P(\omega_i) \quad \Sigma_i = \sigma^2 I \\
 &= -\frac{1}{2\sigma^2} [x^t x - 2\mu_i^t x + \mu_i^t \mu_i] + \ln P(\omega_i) \\
 &= -\frac{1}{2\sigma^2} [-2\mu_i^t x + \mu_i^t \mu_i] + \ln P(\omega_i) \\
 &= \frac{1}{\sigma^2} \mu_i^t x - \frac{1}{2\sigma^2} \mu_i^t \mu_i + \ln P(\omega_i) \\
 &\approx W_i^t x + W_{i0} \rightarrow \text{Linear.}
 \end{aligned}$$


So, let us see this. So, what I have is $g_i(x)$ is equal to minus half log of σ_i minus half x minus μ_i transpose σ_i inverse x minus μ_i plus log of P of ω_i . So, here as I am assuming that this σ_i the covariance matrix is same for all the classes. So, this minus half log of determinant σ_i , again this does not have any role in discriminating among different classes. So, I can simply ignore this term from the function of from the expression of the discriminant function.

So, my $g_i(x)$ now simply becomes minus half x minus μ_i transpose and σ_i because it is $\sigma_i^2 I$. So, the σ_i is of the form $\sigma_i^2 I$. So, this σ_i inverse is simply 1 upon σ_i^2 . So, what I will do is I simply put it as 1 upon $2\sigma_i^2$ x minus μ_i transpose into x minus μ_i plus log of P of ω_i which simplify comes minus 1 upon $2\sigma_i^2$. If I expand this it becomes x transpose x minus twice μ_i transpose x plus μ_i transpose μ_i plus log of P of ω_i . In this expression again this x transpose x is class independent, right. So, again this term does not contribute to discrimination.

So, I further simplify this as minus 1 by $2\sigma_i^2$. What I have within the bracket is minus $2\mu_i$ transpose x plus μ_i transpose μ_i plus log of P of ω_i . You will simplify this, it simply becomes μ_i transpose 1 upon σ_i^2 μ_i transpose x minus 1 upon $2\sigma_i^2$ μ_i transpose μ_i plus log of P of ω_i which I can write in the form W_i transpose x plus W_{i0} where this W_i is nothing, but 1 upon

$\sigma^2 \mu_i^T W$ and $W^T \mu_i$ is $\frac{1}{2\sigma^2} \mu_i^T \mu_i + \log P(\omega_i)$.

So, find that the expression that you get is a linear expression. That means, under the simplified case when all the components of the feature vectors are statistically independent, all the components have the same variance σ^2 and this is same for all the classes. Or, in this particular case I am not assuming it is same for all the classes, but I am considering only a particular class ω_i . The $g_i(X)$ is simply of the form of $W^T X + W^T \mu_i + \log P(\omega_i)$ which is a linear expression, right. So, from here I can try to find out that what is the boundary between two different classes ω_i and ω_j .

(Refer Slide Time: 33:19)

Discriminant Function under Multivariate Normal Distribution

$$g(x) = g_i(x) - g_j(x) = 0$$

$$g_i(x) = -\frac{1}{2\sigma^2} (x - \mu_i)^T (x - \mu_i) + \ln P(\omega_i)$$

$$g_j(x) = -\frac{1}{2\sigma^2} (x - \mu_j)^T (x - \mu_j) + \ln P(\omega_j)$$

$$g(x) = g_i(x) - g_j(x) = 0$$

$$W^T (x - x_0) = 0$$

$$W = (\mu_1 - \mu_2)$$

$$x_0 = \frac{1}{2} (\mu_1 + \mu_2) - \frac{\sigma^2}{\|\mu_1 - \mu_2\|^2} \ln \frac{P(\omega_1)}{P(\omega_2)} (\mu_1 - \mu_2)$$

So, in order to do that let me again try to find out. So, that boundary I can simply defined as $g(x)$ and on a boundary I must have $g_i(x)$ is equal to $g_j(x)$. That is the discriminant functional value for i th class and for j th class; they should be same on the boundary.

So, the equation of the boundary can simply be written as $g(x)$ is equal to $g_i(x)$ minus $g_j(x)$ which is equal to 0. So, this is simply the equation of the boundary between two classes ω_i and ω_j . And, what we have seen is that for $g_i(x)$ under this simplified assumption, we have seen that $g_i(x)$ is nothing, but $-\frac{1}{2\sigma^2} \mu_i^T X - \frac{1}{2\sigma^2} \mu_i^T \mu_i + \log P(\omega_i)$. Similarly

for $g_j X$ I will also have the case that it is $\frac{-1 \pm \sqrt{1 - 4\sigma^2}}{2\sigma^2} X$.
Now it will be $\mu_j^T X - \mu_j + \log P(\omega_j)$.

And, if I equate these two if I make $g_i X$ equal to $g_j X$ to be equal to 0, then we will find that by putting $g_i X$, this expression and in place of $g_j X$ this expression you will find that this $g_i X$ equal to 0. This will take a form $W^T X - \mu_1 - \mu_2 X$ is equal to 0 where, you will find that this W is nothing, but $\mu_1 - \mu_2$. And, X naught will be same as $\frac{\mu_1 + \mu_2 - \sigma^2}{\mu_1 - \mu_2}$ upon $\log P(\omega_1) - \log P(\omega_2)$ into $\mu_1 - \mu_2$.

So, this is the expression that I will get for the boundary between the two different classes. So, I will derive the expression of this boundary under this simplified case in our next lecture.

Thank you.