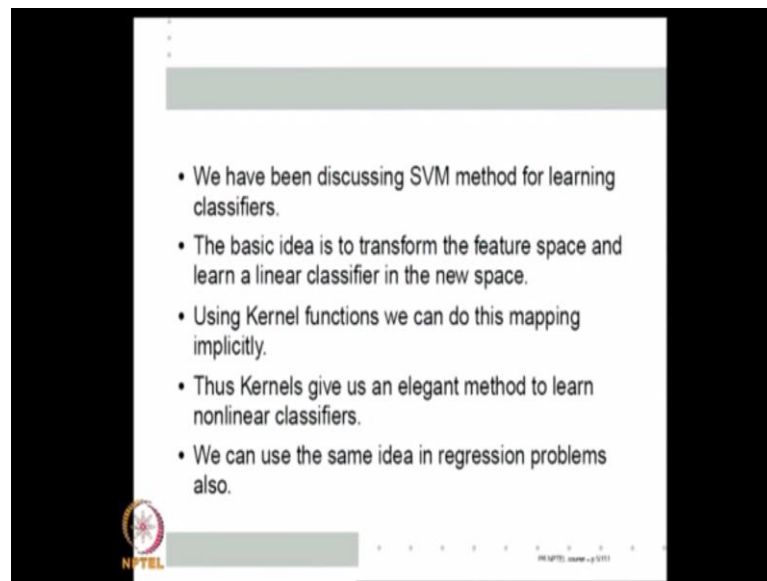


Pattern Recognition
Prof. P. S. Sastry
Department of Electronics and Communication Engineering
Indian Institute of Science, Bangalore

Lecture - 35
Support Vector Regression and E-insensitive
Loss function, examples of SVM learning

Hello and welcome to this next lecture in pattern recognition. We will continue with the support vector machine method. We saw last class the basic support vector classifier for two class case and we have also look at the kernel functions this class. We look at how to poss the regression problem also in the same way. So, we derive a support vector regression algorithm also and then discuss a few issues about how actually we solve for the support vector machine. So, just to briefly recapitulate, we have been discussing the SVM method for learning classifiers.

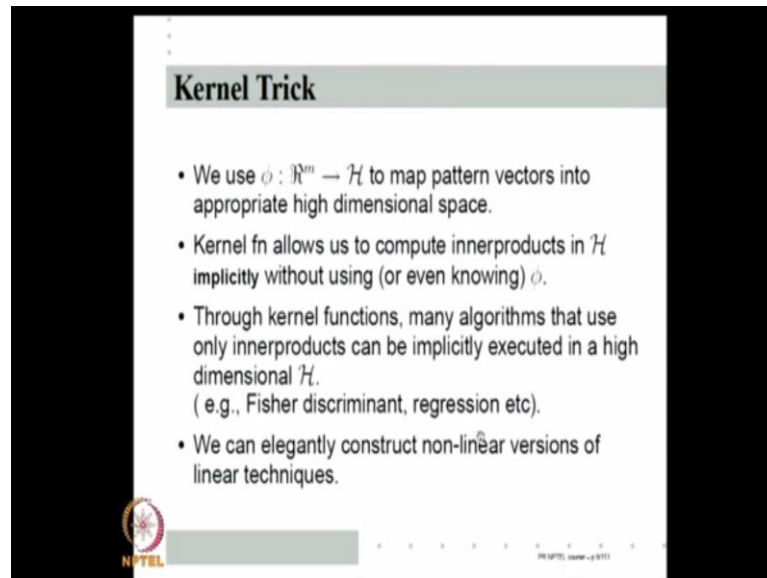
(Refer Slide Time: 01:01)



The basic idea of the SVM method is that we transform the feature space that is with we map the feature vectors into some other space and then learn a linear classifier in the new space. So, basic idea is what is a linear classifier in this new space would be a nonlinear classifier in the old space. So, we can actually learn a nonlinear classifier using techniques of learning a linear classifier. The Kernel functions allow us to do this mapping implicitly and thus Kernels given elegant method to learn nonlinear classifiers.

So, essentially what is a nonlinear classifier in the original space could be a linear classifier in the transform space and the Kernel function allow us to do this mapping implicitly. So, Kernels give us a elegant way to learn nonlinear classifiers the same idea can be use in regression and many other problems. So, let us basically recapitulate what the Kernel Trick is.

(Refer Slide Time: 01:55)



Kernel Trick

- We use $\phi : \mathbb{R}^m \rightarrow \mathcal{H}$ to map pattern vectors into appropriate high dimensional space.
- Kernel fn allows us to compute innerproducts in \mathcal{H} implicitly without using (or even knowing) ϕ .
- Through kernel functions, many algorithms that use only innerproducts can be implicitly executed in a high dimensional \mathcal{H} . (e.g., Fisher discriminant, regression etc).
- We can elegantly construct non-linear versions of linear techniques.

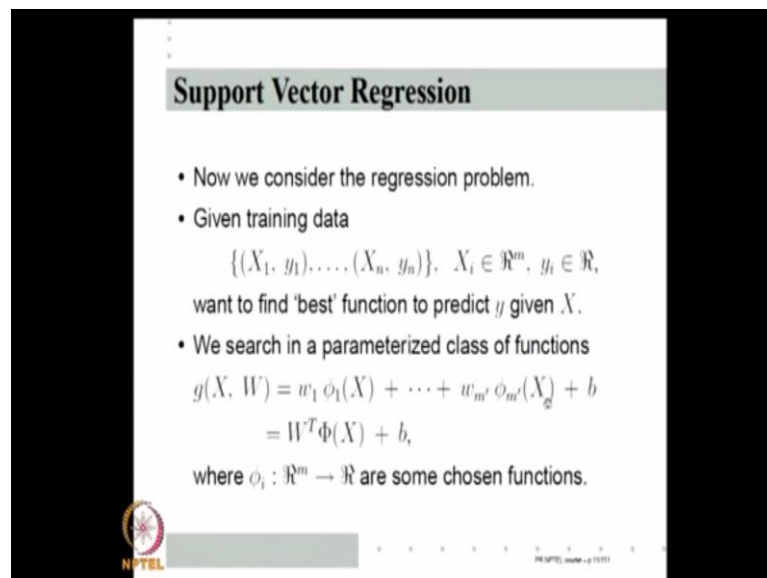
Kernel Trick is if originally the feature space is m dimensional let say \mathbb{R}^m we are mapping the feature vectors from \mathbb{R}^m into some other space I called it \mathcal{H} it could be \mathbb{R}^m prime most of the time or just called it \mathcal{H} which is some appropriate high dimensional space a vector space in which there is an inner product define. What the Kernel function allows us is to compute inner products in the space \mathcal{H} . So, if I have two vectors in \mathbb{R}^m say x and x' I want to calculate the inner product between $\phi(x)$ and $\phi(x')$ I want to calculate $\phi(x)^T \phi(x')$ there is inner product in the \mathcal{H} . I can compute them as a function of just x and x' .

So, Kernel function allows us to compute this inner products implicitly without using ϕ or without even knowing what the ϕ is given x and x' . I can derive a function of x and x' which happens to be the inner product of $\phi(x)$ and $\phi(x')$. So, kernels allow us to do the computation of inner products in the \mathcal{H} space implicitly and with no knowledge of the function ϕ . So, what it means is by use of kernel functions many algorithm that essentially use only inner products can be implicitly execute a high

dimensional space for example, if I look at Fisher discriminant the discriminant which is $w^T x$ as well as the way I learn it using generalized eigenvalue problem all of them only use inner products.

So, I can actually if I think of Fisher discriminant as the finding a linear transform on to one dimensional space where I get the best separation between classes I can ask can I find the nonlinear transform. So, I can have a Kernel Fisher discriminant and so on. So, any algorithm that uses essentially inner product can be implicitly executed in a high dimensional space that is the basic idea of the Kernel Trick it can be using Fisher discriminant and it could be using regression it can be using many other similar technique algorithms. So, basically what it means is using kernels we can elegantly construct nonlinear versions of linear techniques. Now, we going to look at this in the context to regression. So, far you looked at it the context of classification now we will look at it in the context of regression.

(Refer Slide Time: 04:17)



Support Vector Regression

- Now we consider the regression problem.
- Given training data

$$\{(X_1, y_1), \dots, (X_n, y_n)\}, \quad X_i \in \mathbb{R}^m, \quad y_i \in \mathbb{R},$$
 want to find 'best' function to predict y given X .
- We search in a parameterized class of functions

$$g(X, W) = w_1 \phi_1(X) + \dots + w_{m'} \phi_{m'}(X) + b$$

$$= W^T \Phi(X) + b,$$
 where $\phi_i : \mathbb{R}^m \rightarrow \mathbb{R}$ are some chosen functions.

NPTEL

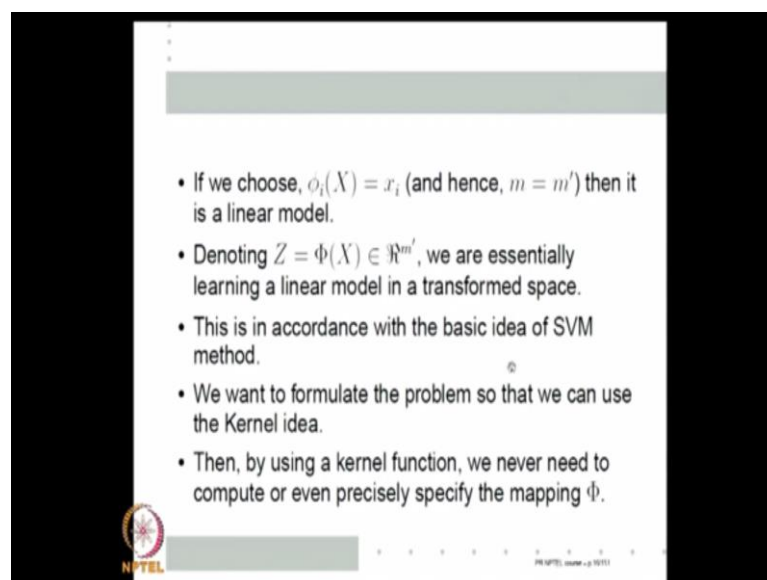
What is the regression problem were given training data $X_1, y_1, \dots, X_n, y_n$ and training samples where each X_i is in \mathbb{R}^m is of like our feature vector only thing the targets now y_1, y_2, \dots, y_n are not plus one minus one but, they are the real value targets. So, given this training data I want to learn the best function to predict y given an X we want to blend a function f such that y had is equal to f of X is a very good prediction for y .

As you seen earlier essentially we can construct such functions in a parameterized manner $X w_1 \phi_1 X w_2 \phi_2 \dots w_{m'} \phi_{m'} X + b$ where m' is some arbitrary number. Lets write this $\phi_1 \phi_2 \dots \phi_{m'}$ as the vector capital Φ and the rest of it as w . So, I can write as $W^T \Phi X + b$ with a bit of use of notation I am putting the parameterization in the left hand side only as W even thou W can as you already seen we can put b inside or outside depending on our convenience.

So, even thou I have a b here I am think of the parameterization is W but, minor point apart this is what we seen when we want when we used the linear least square regression we said as long as the functions $\phi_1 \phi_2 \dots \phi_{m'}$ are fixed this is like linear regression. I can always fix some functions I would have choose what the functions should be and I would have choose what I am trying would be then I can do this. So, Φ are as in general \mathbb{R}^m to \mathbb{R} functions.

Now, essentially if $\phi_i(X)$ is equal to X_i that will give me $w_1 X_1 + w_2 X_2 \dots$. So, one of course, which will mean m' is equal to m . Then, it is a linear model where simply $W^T X + b$. In general, of we have seen this as a generalize linear model earlier but, now we can think of this as if I think of $\phi_1 \phi_2 \dots \phi_{m'}$ as an m' dimensional vector represent that as capital Φ of X which is a transform regression of X .

(Refer Slide Time: 06:32)



The slide contains the following text:


- If we choose, $\phi_i(X) = x_i$ (and hence, $m = m'$) then it is a linear model.
- Denoting $Z = \Phi(X) \in \mathbb{R}^{m'}$, we are essentially learning a linear model in a transformed space.
- This is in accordance with the basic idea of SVM method.
- We want to formulate the problem so that we can use the Kernel idea.
- Then, by using a kernel function, we never need to compute or even precisely specify the mapping Φ .

At the bottom left is the NPTEL logo. At the bottom right, there is a small text: "© NPTEL 2019".

So, I can think of Z is equal to capital ϕ of X which is in $\mathbb{R}^{m'}$ and then we are essentially learning a linear model in a transformed space. This is basically the basic idea of SVM is but, even in our earlier least square regression we said if we fix the function $\phi_1 \phi_2 \dots \phi_m$ for example, we will see how to learn polynomials in 1 variable simply as a linear regression problem. If I fix the functions $\phi_1 \phi_2 \dots \phi_{m'}$ its no difference from linear regression the main thing that you want to gain now is even if I choose the functions ϕ_1 to $\phi_{m'}$ I have to work in the m' dimensional space I have to first transform all the X 's to Z 's and then do a linearly squares regression in a Z space that is not what we want to do. We want to do this problem formulating such a way that you can use the Kernel idea.

So, once again it the same kind of problem as in the classification case. So, I may have a hundred dimensional space in each I want to learn quadratic functions but, if I want to learn quadratic function by explicitly specifying $\phi_1 \phi_2 \dots \phi_{m'}$ I will have ten thousand dimensional linear least square regression problem we do not want to do that we want to use a Kernel function. So, that we do the simplicity we do not even have to precisely specify the mapping ϕ we only specific the function. So, we want to formulate the problem such a way that just like in the classification case we will use the Kernel Trick. So, that we do not have to actually get into the m' dimensional space we do not have to suffer the computational complexity of calculating this ϕ 's and solving the linear least square problem in the m' dimensional space.

(Refer Slide Time: 08:17)



Loss function

- As in a general regression problem, we need to find W to minimize

$$\sum_i L(y_i, g(X_i, W))$$

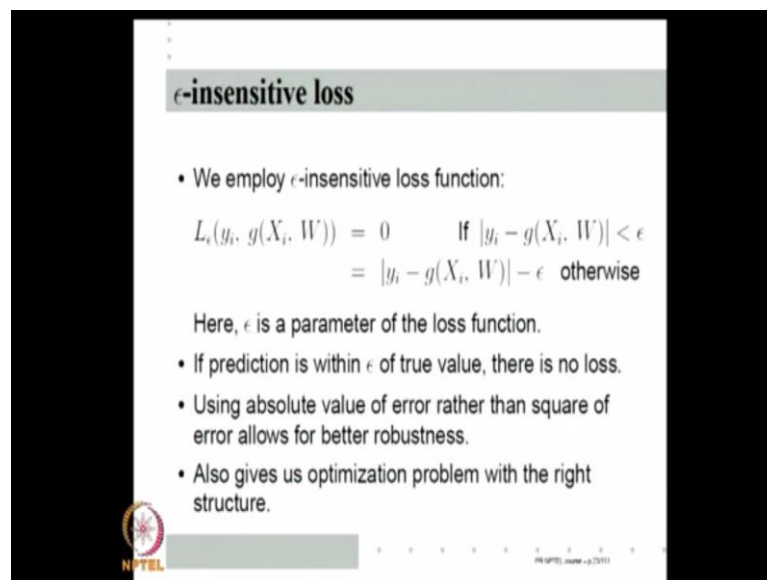
where L is a loss function.

- This is the general strategy of empirical risk minimization.
- We consider a special loss function that allows us to use the kernel trick.

NPTEL

How do we do this in any general regression problem as you already seen like in classification we essentially minimize empirical loss something like the empirical loss. So, if I am learning a model $g(X, W)$. So, $L(y_i, g(X_i, W))$ is the loss. So, given an X_i I would say $g(X_i, W)$ the prediction that is given to me in the training data as y_i . So, $L(y_i, g(X_i, W))$ is the loss. So, I summed over to all i that is the empirical risk minimization of course, I can put a one by n in front but, that does not makes no difference. So, we have been doing. So, far by choosing and appropriate loss function we essentially minimize an empirical list correct this. So, what we are going to do as we already discuss this in last class we are going to use a special loss function that allows us to use the Kernel Trick.

(Refer Slide Time: 09:19)



ϵ -insensitive loss

- We employ ϵ -insensitive loss function:

$$L_{\epsilon}(y_i, g(X_i, W)) = \begin{cases} 0 & \text{If } |y_i - g(X_i, W)| < \epsilon \\ |y_i - g(X_i, W)| - \epsilon & \text{otherwise} \end{cases}$$

Here, ϵ is a parameter of the loss function.

- If prediction is within ϵ of true value, there is no loss.
- Using absolute value of error rather than square of error allows for better robustness.
- Also gives us optimization problem with the right structure.

What is a special loss function is the epsilon insensitive loss function we seen this earlier when we did the formulation of risk minimization empirical risk minimization and we discussed loss functions many different loss functions are listed and this is one of them and also briefly discussed it last class anyway. So, on an example X_i, y_i for the model $g(X, W)$ the loss will be 0 if $|y_i - g(X_i, W)| < \epsilon$ otherwise it is absolute value minus ϵ what is that mean if the prediction is within epsilon of the true value my prediction is $g(X_i, W)$ y_i is the true value on the example. So, if prediction is within epsilon of the true value than there is no loss epsilon of course, a parameter of the loss function that is why I put a epsilon as a subscript here.

So, given as specific value of epsilon if the prediction is within epsilon of the true value there is no loss otherwise the absolute value of the difference on over and above epsilon is the loss. Using absolute value of error rather than square of error earlier we just we are not using any epsilon we just using square of the error and minimizing it using absolute value rather than square of the error allows a better robustness. In the sense a few out layers do not give me do not influence the final fit to much but, more importantly the reason why you use this formula is that this particular loss function when we want to do empirical risk minimization of this class function that gives us optimization problem which has the structure for us the structure is one where I can employ the kernel function alright.

(Refer Slide Time: 11:14)

• We have chosen the model as:

$$g(X, W) = \Phi(X)^T W + b.$$

• Hence empirical risk minimization under the ϵ -insensitive loss function would minimize

$$\sum_{i=1}^n \max(|y_i - \Phi(X_i)^T W - b| - \epsilon, 0)$$

NPTEL

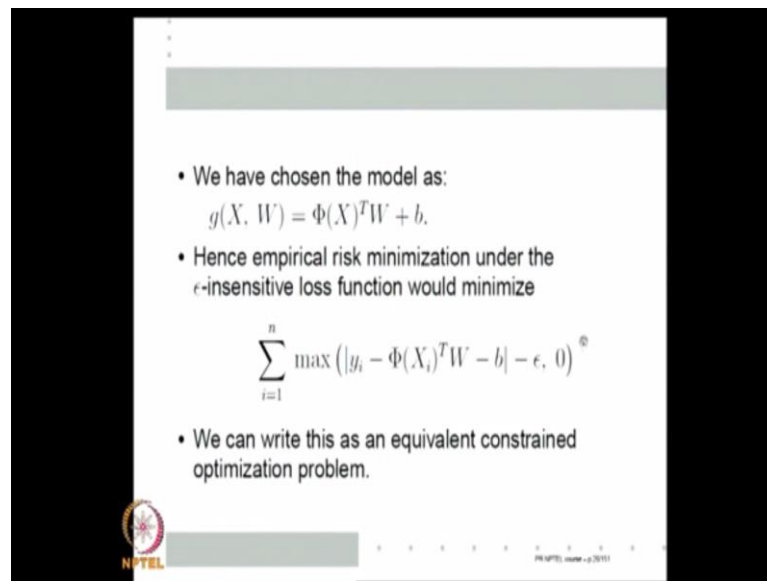
PR. NPTEL, course no. 75111

So, what is it I have to do? this is my model $g(X, W) = \Phi(X)^T W + b$ once again W and b both are included in this W and left hand side essentially I have to learn W on b . By minimizing the empirical risk. So, my empirical risk is for a for each example X_i, y_i my prediction is $\Phi(X_i)^T W + b$ the two prediction is y_i . So, I take the difference between them absolute value the difference minus epsilon if that quantity is greater than 0 that is the loss is that quantity less than 0 loss is 0. So, this is my empirical risk.

That is my epsilon insensitive loss function is the difference is less than epsilon the loss is 0 otherwise loss is this minus epsilon. So, loss is actually this minus epsilon comma 0

whichever is maximum.

(Refer Slide Time: 12:03)



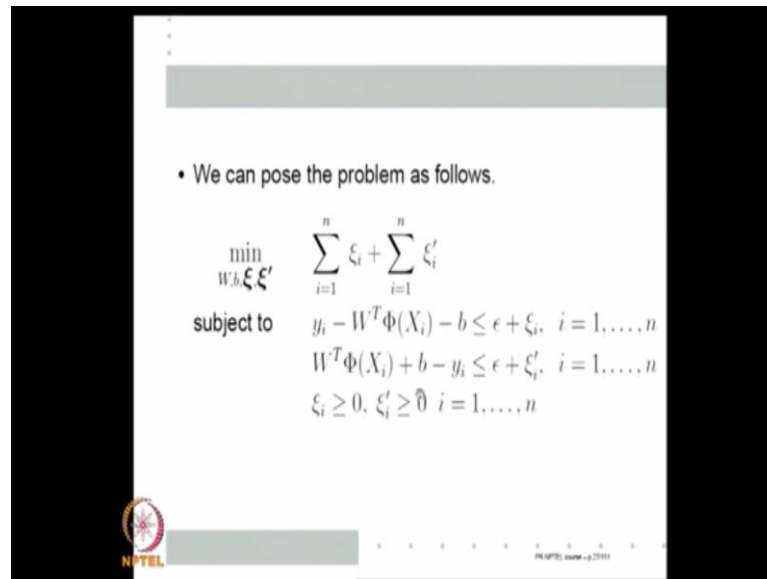
- We have chosen the model as:
$$g(X, W) = \Phi(X)^T W + b,$$
- Hence empirical risk minimization under the ϵ -insensitive loss function would minimize
$$\sum_{i=1}^n \max(|y_i - \Phi(X_i)^T W - b| - \epsilon, 0) \quad \text{⑧}$$
- We can write this as an equivalent constrained optimization problem.

NPTEL

PR NPTEL course 19/2019

So, that is what I have to maximize I can put one by n to make it actually empirical risk I have to put one by n but, that does not really make any difference and I am minimizing this. So, this what I want to minimize. What we are going to now do is we will write this unconstrained minimization into an equivalent constrained optimization problem because that will allow us to get into the structure we want. So, this is what I want to minimize I just remember this. So, basically given a W on b ideally what I want either y_i minus $\Phi(X_i)^T W + b$ less than epsilon I mean which or $\Phi(X_i)^T W + b$ minus y_i less than epsilon whichever is positive actually I want both of them to be less than epsilon because one of them will be negative.

(Refer Slide Time: 12:57)



• We can pose the problem as follows.

$$\min_{W, b, \xi, \xi'} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \xi'_i$$

subject to

$$y_i - W^T \Phi(X_i) - b \leq \epsilon + \xi_i, \quad i = 1, \dots, n$$

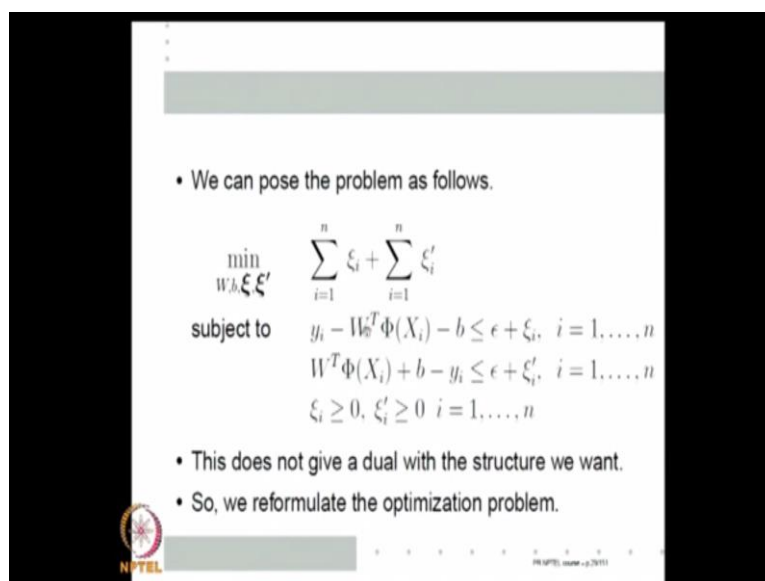
$$W^T \Phi(X_i) + b - y_i \leq \epsilon + \xi'_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad \xi'_i \geq 0 \quad i = 1, \dots, n$$

So, thus what have we write. So, the absolute value of this difference less than epsilon can be written as $y_i - W^T \Phi(X_i) - b$ I wanted to be less than epsilon I want $W^T \Phi(X_i) + b - y_i$ also to be less than epsilon. If one of them will be negative. So, only one of these constraints is really useful if that is less than epsilon I am done if it is not less than epsilon I am using a slack variable ξ_i . So, I need 2 level sets as likewise ξ_i and ξ'_i that tells me by how much the difference between $y_i - W^T \Phi(X_i) + b$ differs from epsilon. Note that, for any given i only one of ξ_i or ξ'_i is greater than 0 both cannot be greater than 0 because one of these left hand sides is negative it really does not matter which one so but, one of them is the operational 1 for any given i and that ξ'_i tells me that ξ_i or ξ'_i gives me the actual absolute value of $y_i - W^T \Phi(X_i) - b - \epsilon$.

So, essentially what is it saying is saying ξ_i one of ξ_i or ξ'_i will be greater than or equal to absolute value of $y_i - W^T \Phi(X_i) - b - \epsilon$ and is also greater than equal to 0. So, if I just add this ξ_i 's this ξ_i and ξ'_i is nothing but, max of this minus epsilon. So, I can always say such an unconstrained maximization problem as a constrained unconstrained minimization problem as a constrained minimization problem of this. So, of course, I am trying to find W, b, ξ_i and ξ'_i to minimize this while satisfying all this. This is nice problem but, which exactly is same as minimizing empirical risk as easy to see.

(Refer Slide Time: 14:57)



• We can pose the problem as follows.

$$\min_{W, b, \xi, \xi'} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \xi'_i$$

subject to

$$y_i - W^T \Phi(X_i) - b \leq \epsilon + \xi_i, \quad i = 1, \dots, n$$

$$W^T \Phi(X_i) + b - y_i \leq \epsilon + \xi'_i, \quad i = 1, \dots, n$$

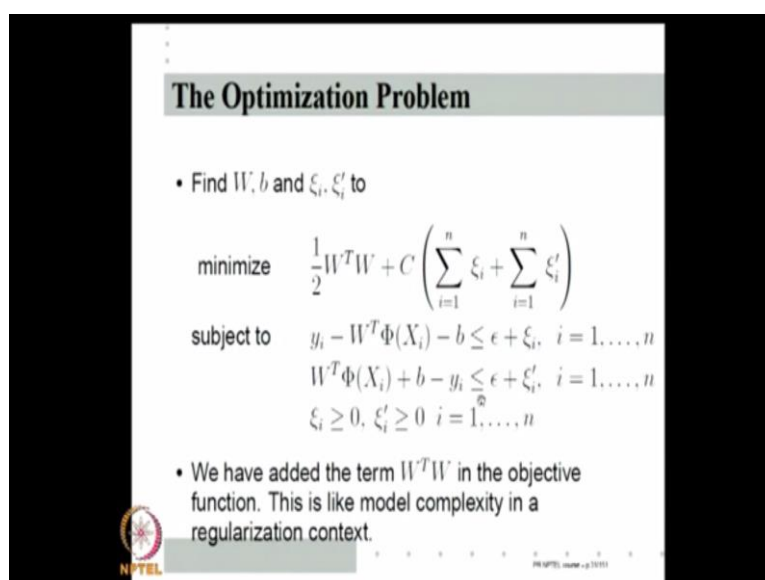
$$\xi_i \geq 0, \quad \xi'_i \geq 0 \quad i = 1, \dots, n$$

• This does not give a dual with the structure we want.

• So, we reformulate the optimization problem.

Minimizing this empirical risk is equivalent to solving this constrained optimization problem. So, the problem is nice, but I have we have a problem we have a problem in the sense this problem is nice but, it does not solve our purpose this does not give a dual with the structure we want. So, its slightly reformulate this what do we do keep the constrained same.

(Refer Slide Time: 15:15)



The Optimization Problem

• Find W, b and ξ_i, ξ'_i to

$$\text{minimize} \quad \frac{1}{2} W^T W + C \left(\sum_{i=1}^n \xi_i + \sum_{i=1}^n \xi'_i \right)$$

subject to

$$y_i - W^T \Phi(X_i) - b \leq \epsilon + \xi_i, \quad i = 1, \dots, n$$

$$W^T \Phi(X_i) + b - y_i \leq \epsilon + \xi'_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad \xi'_i \geq 0 \quad i = 1, \dots, n$$

• We have added the term $W^T W$ in the objective function. This is like model complexity in a regularization context.


But, add a half W transpose W just like. So, that it looks like you know what we solve for the classification why should I had this half W transpose W well we will come back to

that little later in today's class I will explain you why we want to add half $W^T W$ now let us say. So, if I had only this term the constraints and this term than is exactly like minimizing the empirical risk but, I am adding this. So, what we can think of this is with along with these constraints, this is the empirical risk. So, this in some senses like a model complexity.

So, this is the data error as you have already seen this term captures the empirical risk under the epsilon insensitive loss function. So, this is like some term which is modelled complexity. So, I can think of this as a regularized risk minimization we already seen in linear least squares regression that we want to we may want to put a regularization term. So, that we learn a smooth model instead of learning instead of completely putting all our attention in minimizing the data error that is minimizing the empirical risk and we only minimizing empirical risk we may not learn as nice a function as if we also pay a little attention to the complexity of the model there also we just mentioned that an often use model complexity term for linear models is the norm of W . So, the same norm of W we are using here as a regularization term as I promise to by before the end of today's class we will come back and ask why is this a good way to characterize the smoothness of the model where learning.

So, for now we simply say that earlier if I did not have that $W^T W$ as exactly equal to the empirical risk. So, in addition to minimizing empirical risk I add a model complexity term. So, that this problem now corresponds to a regularized empirical risk minimization. Now, we are in business this now is very easy to see except that instead of $\frac{1}{n} \sum x_i$ I have x_i and x_i' I essentially have similar kind of constrained and similar kind of problem structure as epsilon. So, I would expect similar kinds of things to happen.

(Refer Slide Time: 17:33)



- Like earlier, we can form the Lagrangian and then, using Kuhn-Tucker conditions, can get the optimal values of W^* and b .
- Given that this problem is similar to the earlier one, we would get W^* in terms of the optimal Lagrange multipliers as earlier.
- Essentially, the Lagrange multipliers corresponding to the inequality constraints on the errors would be the determining factors.
- We can use the same technique as earlier to formulate the dual to solve for the optimal Lagrange multipliers.

NPTEL

So, like earlier we can form the Lagrangian this is objective function this is constrained form the Lagrangian put in the Kuhn-Tucker conditions. That will give you the expressions for optimization values of W and b . Given this structural similarity we will expect the optimal values of W on b to be similar to the earlier expression. So, it means we can get W^* and b^* in terms of the optimal Lagrange multipliers what Lagrange multiplier essentially the Lagrange multipliers corresponding to the inequality constraints on the errors just like there also the Lagrange constraints and ξ_i 's are not important towards the constraints $\xi_i \geq 0$ Lagrange multiplier corresponding those constraints not important. Basically inequality constraints that on the errors in the data would be the determining factors and we can always use the same technique as earlier that is we formulate the dual to solve this.

We can once again this is a this is as you can see a quadratic cost function a convex cost function and linear constraints. So, I you can use the same technique as you use earlier for deriving the dual using the dual function.

(Refer Slide Time: 18:53)

The dual

- The dual of this problem is

$$\max_{\alpha, \alpha'} \sum_{i=1}^n y_i(\alpha_i - \alpha'_i) - \epsilon \sum_{i=1}^n (\alpha_i + \alpha'_i) - \frac{1}{2} \sum_{i,j} (\alpha_i - \alpha'_i)(\alpha_j - \alpha'_j) \Phi(X_i)^T \Phi(X_j)$$

subject to

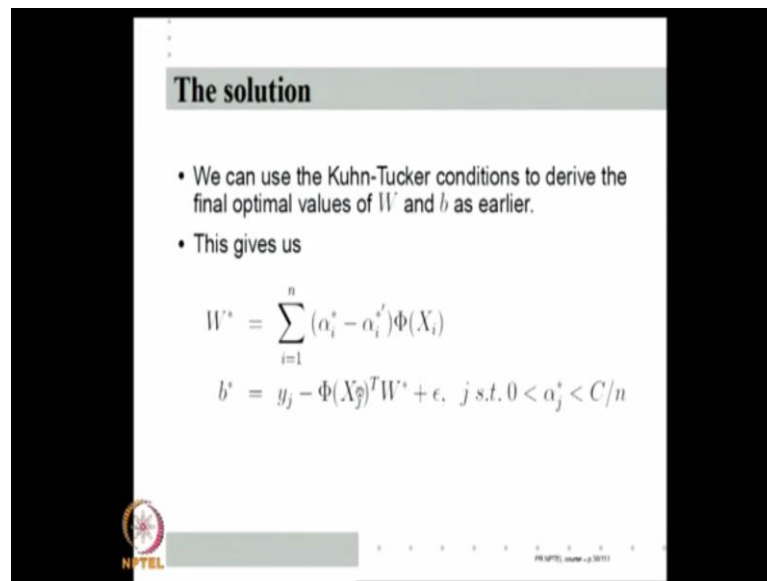
$$\sum_{i=1}^n (\alpha_i - \alpha'_i) = 0$$

$$0 \leq \alpha_i, \alpha'_i \leq C, \quad i = 1, \dots, n$$

And if you do all that the dual transform to be this we will just omit the details look at the derivations quite similar. So, now we have two Lagrange multipliers alpha and alpha prime. So, basically there 4 constraints 1,2,3,4. So, the constraints for this lets say this as the these are the xi i constraints these are the xi i prime constraints. So, Lagrange multiplier corresponds with xi constraints are alpha i Lagrange multiplier correspond with xi prime constraints are alpha I prime. These constraints are not important just like in the SVM case Lagrange multiplier correspond to these constraints would not enter into the dual. So, dual is in terms of alpha and alpha prime, I have we have alpha i i is equal to 1 to n and alpha i prime i is equal to 1 to n those 2 n variables are the Lagrange multiplier corresponding to these 2 n constraints.

So, that turns out to be the objective function and these half turn out to be constraints once again like earlier the all the Lagrange multipliers have to between 0 and C and there is one equality constraints and then the objective function the dual is also quadratic in both alpha i and alpha i prime and the training data themselves appear only as inner products phi X i transpose phi X j. So, all the structure that we had for the SVM dual still is there the details do not really matters to us the actual terms at different here compare to the SVM but, the overall structure is same it's a quadratic cost function one in a one linear equality constraint bound constraints and variables and the bound is between 0 and C where C is the same penalty constraints and the training data appear only as inner products in the objective function phi X i transpose phi X j.

(Refer Slide Time: 20:54)



The solution

- We can use the Kuhn-Tucker conditions to derive the final optimal values of W and b as earlier.
- This gives us

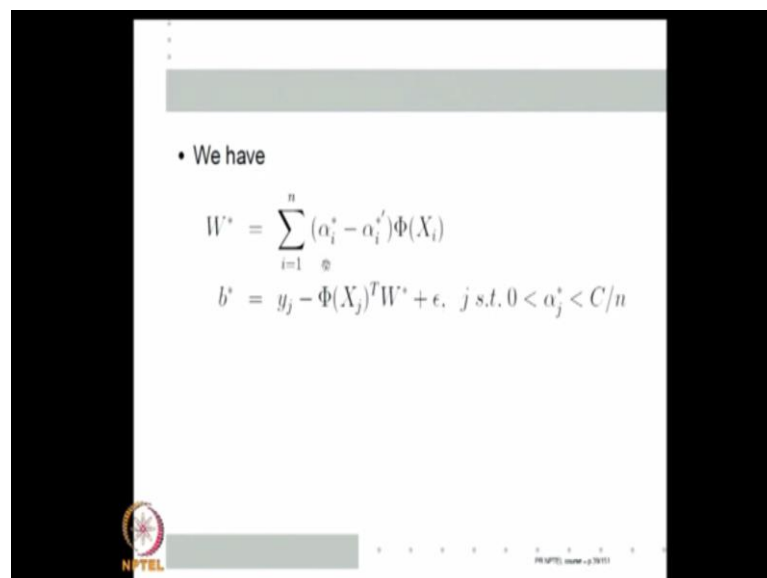
$$W^* = \sum_{i=1}^n (\alpha_i^* - \alpha_i^{\prime*}) \Phi(X_i)$$

$$b^* = y_j - \Phi(X_j)^T W^* + \epsilon, \quad j \text{ s.t. } 0 < \alpha_j^* < C/n$$

NPTEL

Using the Kuhn-Tucker conditions as earlier we can derive the optimal values of W and b and it trans out to be the same W star is sum all α_i star minus α_i star prime $\phi(X_i)$. So, as I said α_i 's are correspond with the first constraints α_i prime say correspond second constraints there. So, this happens to be the W star and b star.

(Refer Slide Time: 21:18)



- We have

$$W^* = \sum_{i=1}^n (\alpha_i^* - \alpha_i^{\prime*}) \Phi(X_i)$$

$$b^* = y_j - \Phi(X_j)^T W^* + \epsilon, \quad j \text{ s.t. } 0 < \alpha_j^* < C/n$$

NPTEL


Let us look at this expression.

First let us note certain thing about this primal problem. As you already seen, only 1 of x_i

and ξ_i is non-zero the other because on the LHS one of the terms is negative. So, for any given i only if at all only 1 of these constraints can be active. Only one of α_i and α_i' will be non-zero. α_i is strictly greater than 0 which means because Lagrange multipliers are α_i here and α_i' here the complementary slackness shows us that both α_i and α_i' cannot be simultaneously greater than 0 for them to be complementary. At least one of them should be equal to 0. Both inequalities should be satisfied by equality that is not possible. So, only one of α_i and α_i' would be greater than 0 at the optimality point and similarly, if y_i is negative, ξ_i is positive.

So, $y_i - W^T \Phi(X_i) - b$ if it is actually strictly less than ϵ then ξ_i will be 0. So, once again even among that particular set for each i only one of α_i and α_i' are greater than 0. As a set for each i in addition not all α_i 's are α_i' will be if I look at all the strictly greater than 0 Lagrange multipliers they may not be enough. The most probably will not be enough. They are like the support vectors in the classification case basically only when this inequality has to be strictly satisfied only that is why. So, that ξ_i is strictly greater than 0 only then the corresponding α_i will be positive.

(Refer Slide Time: 20:19)



• We have

$$W^* = \sum_{i=1}^n (\alpha_i^* - \alpha_i'^*) \Phi(X_i)$$

$$b^* = y_j - \Phi(X_j)^T W^* + \epsilon, \quad j \text{ s.t. } 0 < \alpha_j^* < C/n$$

- Note that we have $\alpha_i^* \alpha_i'^* = 0$. Also, $\alpha_i^*, \alpha_i'^*$ are zero for examples where error is less than ϵ .
- The final W^* is a linear combination of some of the examples – the support vectors.
- Note that the dual and the final solution are such that we can use the kernel trick.

NPTEL

So, which means we know by looking at the primal problem that only 1 of α_i and α_i' is non-zero and both of them will be 0 for example, where error is strictly less than ϵ . Now, when I think in my mind if I look at the W^* essentially W

star is nothing but, a linear combination of some of your data $\phi(X_i)$ and the coefficients are the Lagrange multipliers because in this term is either α_i^* or $-\alpha_i^*$. So, coefficients are simply Lagrange multipliers. So, the final W is a linear combination of some of the examples which are these examples the support vectors where ever the corresponding Lagrange multipliers strictly positive and the dual and the final solution both are such that the Kernel Trick is visible. The dual as we already seen the data $\phi(X_i)$ if we are only as inner products.

So, I can use the kernel there and the W^* is a linear combination of $\phi(X_i)$. So, once again I can use the Kernel Trick and the b^* comes as $\phi(X)^T W^* + b^*$. So, once again, I can use the Kernel Trick.

(Refer Slide Time: 24:40)

• Let $K(X, X') = \Phi(X)^T \Phi(X')$.

• The optimal model learnt is

$$g(X, W^*) = \sum_{i=1}^n (\alpha_i^* - \alpha_i') \phi(X_i)^T \phi(X) + b^*$$

$$= \sum_{i=1}^n (\alpha_i^* - \alpha_i') K(X_i, X) + b^*$$

• As earlier, b^* can also be written in terms of the Kernel function.

NPTEL

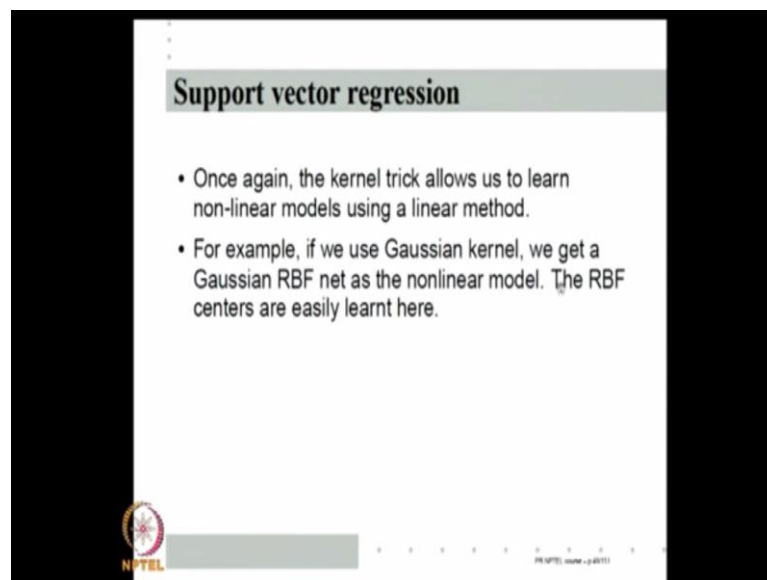
Specifically, suppose we have a kernel function then what is our optimal model $g(X, W^*)$ is all this is my expression for W . So, $W^T \phi(X) + b^*$. So, this can be canalized this can be done as $\sum K(X_i, X)$.

And b^* and b^* can also be canalized because b^* is $\phi(X)^T W^* + b^*$. W^* is linear in $\phi(X)$. So, this term also involves only $\phi(X_i)^T \phi(X_i)$. So, once again this can be canalized. So, this can be canalized and as earlier b^* can also be written in terms of the kernel function. So, the way we formulated the regression problem let us go back to our formulation.

This is how we formulate the regression problem. We put in this extra regularization context and otherwise it is like the SVM formulation this is actually minimization of empirical risk under epsilon insensitive loss function along with a regularization context. And for that problem that happens to be the dual dual has the very large structure.

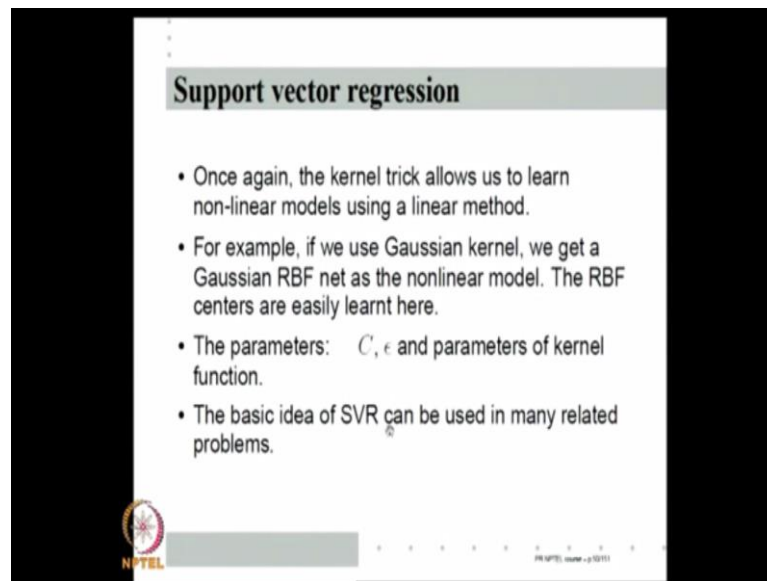
And with that dual I get my solution in the way. So, that the solution allows me the Kernel Trick specifically, I can write my final function like this where b^* can also be obtained through the kernels. So, essentially now I do not have to calculate ϕ 's, I only stay in the m dimensional space of the original space of the X i's and still learn nonlinear regression functions.

(Refer Slide Time: 26:42)



So, once again the Kernel Trick allows us to learn nonlinear models using a linear method for example, if we use Gaussian kernels we get a Gaussian RBF net. So, if this is Gaussian $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$ is exponential minus norm of x_i minus x_j whole square than this is nothing but, linear combination of Gaussian real basis functions. So, this is actually the output of an RBF network the stand RBF network Gaussian RBF network with a Gaussian real base function hidden nodes and linear output node.

(Refer Slide Time: 27:23)



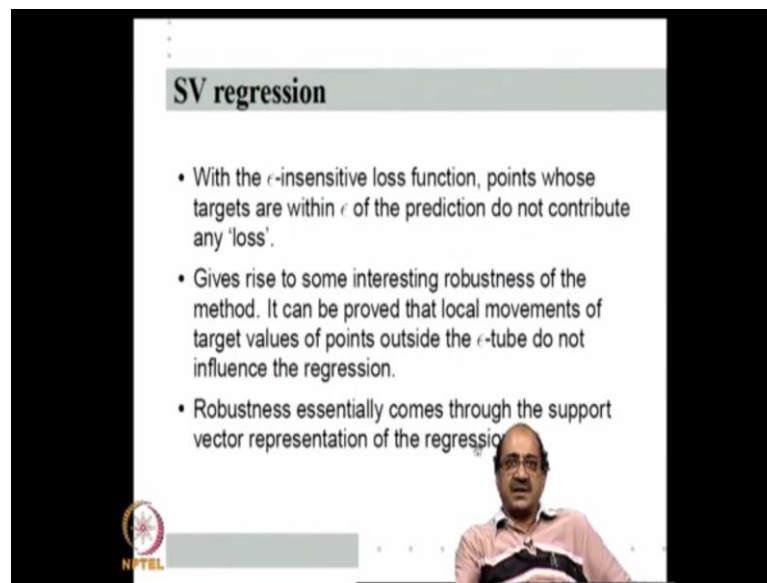
Support vector regression

- Once again, the kernel trick allows us to learn non-linear models using a linear method.
- For example, if we use Gaussian kernel, we get a Gaussian RBF net as the nonlinear model. The RBF centers are easily learnt here.
- The parameters: C , ϵ and parameters of kernel function.
- The basic idea of SVR can be used in many related problems.

NPTEL

So, the kernels automatically give you this. So, if you use a Gaussian kernel we get a Gaussian RBF net just like in the SVM case and just like in SVM case RBF centers are easily learnt here they happened to be the support vectors. There are parameters, of course, any algorithm as parameters, we have the penalty contrast C which is like the regularization constraints epsilon where choose the epsilon in the epsilon insensitive loss function and any parameter the kernel function for example, if I am using a Gaussian kernel I may use a sigma. So, these parameters have to be choosing properly, but, if I choose the parameters, I get fairly good performance and this basic idea of supported to regression is also used in many related problems.

(Refer Slide Time: 28:13)



SV regression

- With the ϵ -insensitive loss function, points whose targets are within ϵ of the prediction do not contribute any 'loss'.
- Gives rise to some interesting robustness of the method. It can be proved that local movements of target values of points outside the ϵ -tube do not influence the regression.
- Robustness essentially comes through the support vector representation of the regression.

NPTL

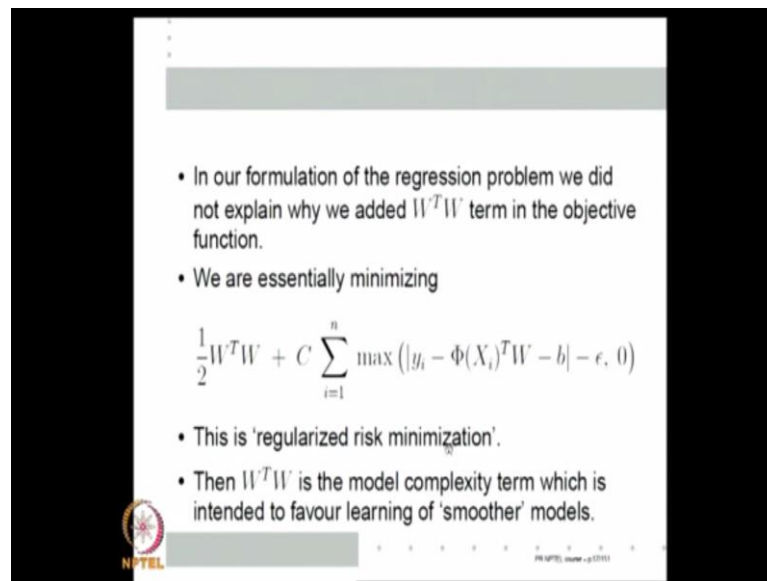
Basically, what happened was by using the epsilon insensitive loss function point whose targets are within epsilon of the prediction do not contribute to loss that as what as given rise to the nice structure in the problem and it also gives us some interesting like this the support vectors as we as we know in the general classification problem the support vectors are very good by product as we said they represent all the patterns which are closest to the mount debutante clauses. So, and once we solve a problem and get all the support vectors if am I training data set I remove all the things for the support vector i will still get the same SVM because those are the closest to the training class boundary.

So, if I can separate them we will anyway separate the rest of them. So, the SVM on the entire training data will be same as SVM reruns on the support vectors. So, essentially same way the support vectors represent the regression function two for example, the substant because the epsilon loss function will can show that local movement of target values of target values of points outside the epsilon tube that is the epsilon tube is like the $W^T X + b$ is equal to plus one or minus one through two parallel hyper planes in our classifier here given our predictor function $W^T \phi(X)$ any prediction that within epsilon of this does not contribute any loss.

So, I can put an epsilon tube around this prediction function and we can show that if you move targets outside this that does not much influence a regression. So, robustness

comes through this support vector representation of the regression.

(Refer Slide Time: 30:13)

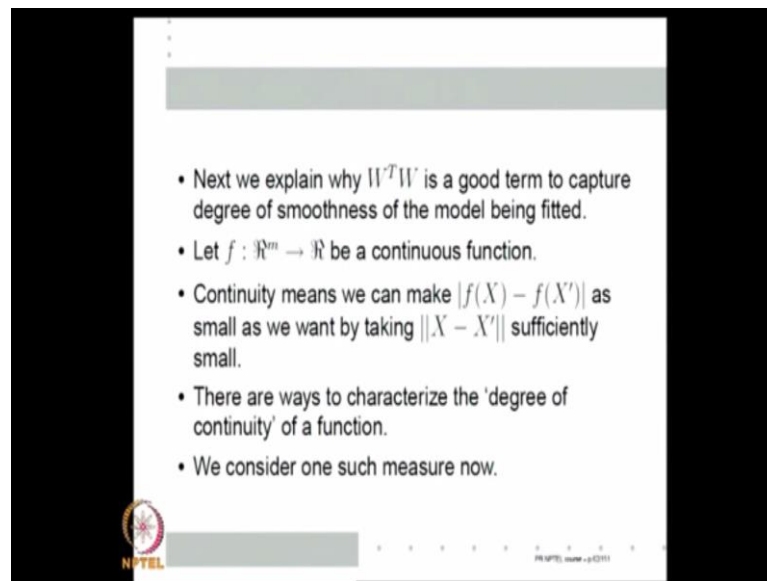


- In our formulation of the regression problem we did not explain why we added $W^T W$ term in the objective function.
- We are essentially minimizing
$$\frac{1}{2} W^T W + C \sum_{i=1}^n \max(|y_i - \Phi(X_i)^T W - b| - \epsilon, 0)$$
- This is 'regularized risk minimization'.
- Then $W^T W$ is the model complexity term which is intended to favour learning of 'smoother' models.

Now, let us get back to our regularization term as I promise you I said I will tell you why it is. All to add $W^T W$, we chosen an epsilon insensitive loss function we want to minimize empirical insensitive loss function and the rest of the formulation is fine except the term $W^T W$ and we said it is like a model complexity term. So, basically this is my empirical risk $\frac{1}{n} \sum_{i=1}^n \max(|y_i - \Phi(X_i)^T W - b| - \epsilon, 0)$.

This is what is captured by $\sum_{i=1}^n x_i^2$ plus $\sum_{i=1}^n x_i^2$ prime subject to those constraints in addition I added this. So, this is essentially what I am minimizing. So, this is the empirical risk this the model complexity and that is what I minimize. So, this is regularized a risk minimization the only thing we have to show is why is this a good regularization term. So, $W^T W$ is a model complexity term which is intended to favor learning of smooth models why do we use regularization why do we use a model complex term just because I am minimizing my data error is not enough as we discuss when we discuss statistical learning theory we need to learn simpler models the simpler the model the better it is. So, essentially we want more smoother models. So, $W^T W$ term is a complexity term which intended to favor smoother model.

(Refer Slide Time: 31:40)

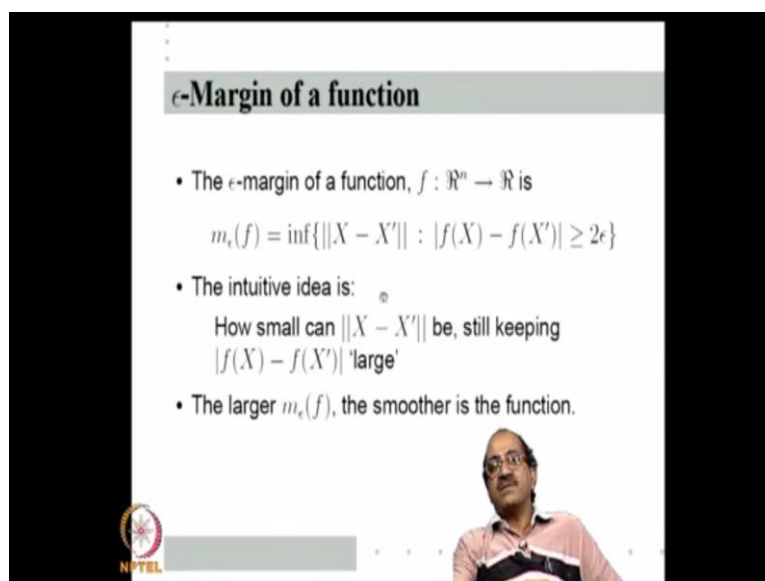


- Next we explain why $W^T W$ is a good term to capture degree of smoothness of the model being fitted.
- Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a continuous function.
- Continuity means we can make $|f(X) - f(X')|$ as small as we want by taking $\|X - X'\|$ sufficiently small.
- There are ways to characterize the 'degree of continuity' of a function.
- We consider one such measure now.

So, we have to ask why is this a good term to capture degree of smoothness of a function being fitted smoothness is being nicely continuous there many continuities a continuous function is kind of a continuity that binds the property a function either continuous or not continuous but, all continuous functions are not same or a say $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is a continuous function what does continuity mean I can make $f(X) - f(X')$ the difference between $f(X)$ and $f(X')$ as small as I want by taking the norm of the difference $X - X'$ sufficiently small but, if f is even if f is continuous if it is varying rapidly than to make the difference between $f(X)$ and $f(X')$ small by the same epsilon I may have to make the difference between the X and X' very very small.

So, for a given difference between $f(X)$ and $f(X')$ how small should I have to make $X - X'$ is a kind of one kind of measure to say how fast or slowly varying the function is and which is one measure of how smooth the function is. So, there are ways to characterize what can be called degree of continuity a degree of smoothness a function were going to look at just one such measure we just look at one such measure how to characterize the degree of continuity function.

(Refer Slide Time: 33:06)



ϵ -Margin of a function

- The ϵ -margin of a function, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is
$$m_\epsilon(f) = \inf \{ \|X - X'\| : |f(X) - f(X')| \geq 2\epsilon \}$$
- The intuitive idea is:
How small can $\|X - X'\|$ be, still keeping $|f(X) - f(X')|$ 'large'?
- The larger $m_\epsilon(f)$, the smoother is the function.

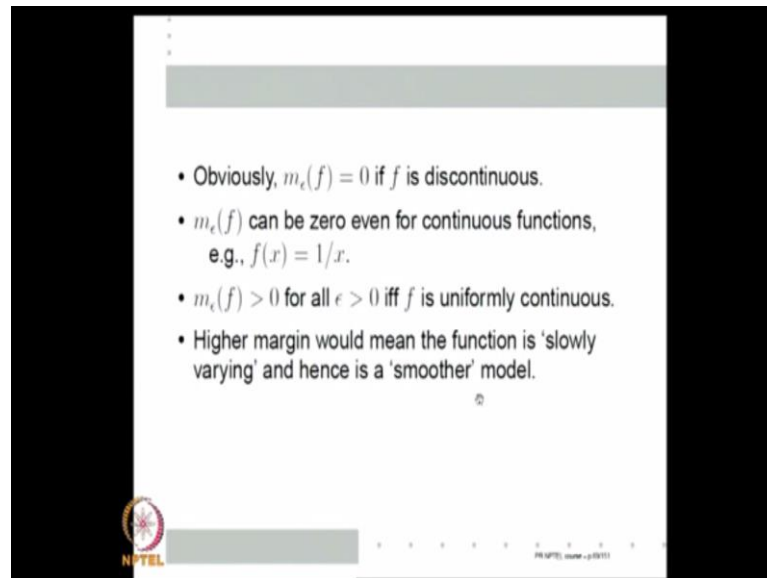
The measure we look at is what is known as an epsilon margin of a function given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and some epsilon greater than 0 we for that epsilon the epsilon margin of the function we generate is m_ϵ of f is infimum over all the numbers norm X minus X' such that $|f(X) - f(X')| \geq 2\epsilon$ let say given an epsilon I want to keep $|f(X) - f(X')|$ at least two epsilon where $|f(X) - f(X')|$ has to be large. So, do say $|f(X) - f(X')|$ has to be greater than 2 epsilon I am asking of all X minus X' that will allow $|f(X) - f(X')|$ to be greater than 2 epsilon what is this smallest of that. So, how small X and X' and how close X and X' be still keeping $|f(X) - f(X')|$ large.

So, if $m_\epsilon(f)$ is small that means X and X' even very small difference between X and X' can make $|f(X) - f(X')|$ large than is not very good on the other hand if m_ϵ is large than its good the larger $m_\epsilon(f)$ the smoothly the function $m_\epsilon(f)$ the epsilon margin being large means what the smallest of norm X minus X' which ensures $|f(X) - f(X')| \geq 2\epsilon$ is still large. So, even if I want to make $|f(X) - f(X')|$ more than 2 epsilon I can still keep I do not have to keep my X and X' very close. So, that is what the basic idea here is. So, the larger the epsilon margin the smoother is the function.

Now, it is easy to see that for suppose f is discontinuous if f is discontinuous than epsilon margin will be 0 because around the discontinuity at the point of discontinuity, no matter

how small I make X minus X prime this can still be greater than 2ϵ because $f(X)$ minus $f(X')$ is strictly greater than 0. So, they will be an ϵ such that this is always greater than 2ϵ no matter how small I make X minus X' .

(Refer Slide Time: 35:26)



So, this is easy to see that the epsilon margin is 0 the function discontinuous but, interesting the function the epsilon margin can be 0 even if the function continuous if you look at a function one-xth obviously is continuous but, as x comes closer to 0 it is varying fast and faster because its rate of change is one by the order one by x square. So, x becomes closer to 0 its varying fast and faster. So, its not a really smooth function because its varies very very fast. So, they can be a continuous function like this for which also epsilon margin is 0. One can show that for every epsilon greater than 0 if the epsilon margin has to be greater than 0 then f is to be uniformly continuous do not worry if you do not know what uniform continuous we really do not need it.

Essentially, higher the margin the smoother the function is that the function is slowly varying and hence it is smoother that is basically what this is. We have defined this as the epsilon margin and is easy to see that higher the epsilon margin the smoother the function.

(Refer Slide Time: 36:39)

SVR and margin

- Consider regression with linear models. Then,

$$|f(X) - f(X')| = |W^T(X - X')|.$$
- For all X, X' with $|W^T(X - X')| \geq 2\epsilon$, $\|X - X'\|$ would be smallest if $|W^T(X - X')| = 2\epsilon$ and $(X - X')$ is parallel to W .
That is, $X - X' = \pm \frac{2\epsilon W}{W^T W}$.
- Thus, $m_\epsilon(f) = \frac{2\epsilon}{\|W\|}$.
- Thus in our optimization problem in SVR, minimizing $W^T W$ promotes learning of smoother models.

So, now let us calculate the epsilon margin for linear models. So, we have a linear model than $f(X) - f(X')$ because $f(X)$ is $W^T X + b$ and $f(X')$ is $W^T X' + b$ then $f(X) - f(X')$ will be $W^T(X - X')$. So, to find epsilon margin whatever to do given X and X' such that $W^T(X - X') \geq 2\epsilon$ find all X and X' which satisfy $W^T(X - X') \geq 2\epsilon$ among all such pairs calculate $\|X - X'\|$ and ask for which pair is $\|X - X'\|$ smallest that is what is infimum in the definition?

So, we are asking when would $\|X - X'\|$ be smallest if we have to make $W^T(X - X') \geq 2\epsilon$ firstly the concern is greater than 2ϵ because I want $\|X - X'\|$ to be smallest and this inner product is nothing but, $\|W\| \|X - X'\| \cos \theta$ into some cost of angle. So, do say its better to take it to be 2ϵ . So, $W^T(X - X') = 2\epsilon$ would be satisfied by X and X' which have smaller norms than $W^T(X - X') \geq 2\epsilon$ and among all X and X' that will satisfy $W^T(X - X') = 2\epsilon$.

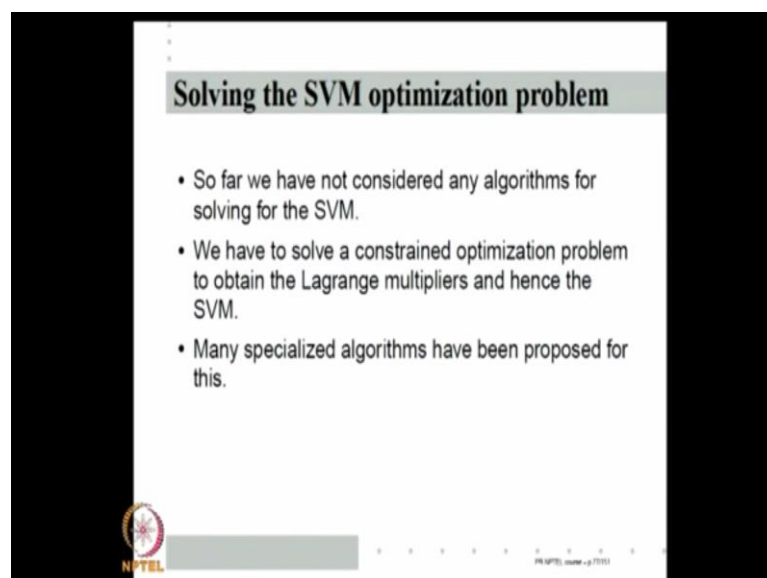
The one that we have the smallest norm if that pair such that the vector $X - X'$ is parallel to W , here we know this a norm $\|W\| \|X - X'\| \cos \theta$ into some factor which is less than 1 that factor becomes one of the two vectors are parallel. So, 2 vectors

are not parallel because that factor is less than one corresponding a norm has to be larger. So, that is still 2ϵ . So, the smallest the pair X and X' with the smallest norm $X - X'$ which still satisfies this is such the $W^T (X - X')$ is equal to 2ϵ and the vector $X - X'$ is parallel to W what is that mean $X - X'$ should be parallel to W . So, it should be k times W and $W^T (X - X')$ should be 2ϵ .

So, if is k times W $k^T W^T W$ should be equal to 2ϵ there is an absolute value. So, k can be only plus minus 2ϵ by $W^T W$. So, $X - X'$ is either plus 2ϵ by $W^T W$ into W or minus of that. So, my epsilon margin is nothing but, the norm of this smallest $X - X'$ if I take norm of this that will be 2ϵ by $W^T W$ into root of $W^T W$. So, that will give me 2ϵ by norm W . So, for linear models the epsilon margin is nothing but, 2ϵ by norm W this is what we call the margin of the hyper plane the classifier.

So, essentially when I added a term to minimize norm W I am promoting learning of the smoother. So, for all linear models because the epsilon margin of the function is inversely proportional to norm W or norm W square is a good regulation term to add we are done it regularized least squares also earlier without explaining here is the reason why norm W square is a good regularization. So, this completes our discussion on regression.

(Refer Slide Time: 40:18)



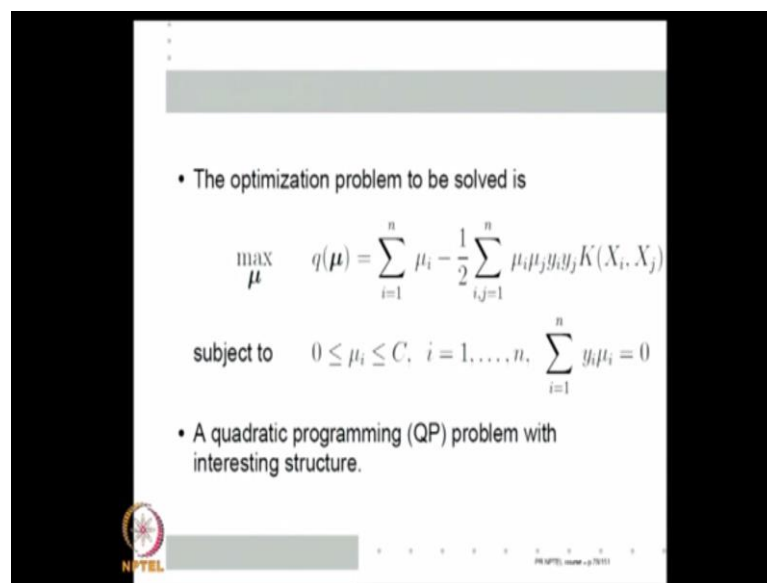
Solving the SVM optimization problem

- So far we have not considered any algorithms for solving for the SVM.
- We have to solve a constrained optimization problem to obtain the Lagrange multipliers and hence the SVM.
- Many specialized algorithms have been proposed for this.

NPTEL

So, far we looked at support vector machines for classification regression but, than both cases we just said solve the dual we never said how do I solve the dual dual is a constrained optimization problem were not really given any algorithm for solving the dual we have to solve a constrained optimization problem to obtained the lagrange multipliers and once you get the Lagrange multipliers we can get the SVM. So, there is one optimization problem namely the dual which is a quadratic programming problem which needs to be solve we are not looked at any algorithm. So, far for this dual is many some interesting structure. So, that many specialized algorithm proposed for this.

(Refer Slide Time: 41:00)



• The optimization problem to be solved is

$$\max_{\mu} \quad q(\mu) = \sum_{i=1}^n \mu_i - \frac{1}{2} \sum_{i,j=1}^n \mu_i \mu_j y_i y_j K(X_i, X_j)$$

subject to $0 \leq \mu_i \leq C, \quad i = 1, \dots, n, \quad \sum_{i=1}^n y_i \mu_i = 0$

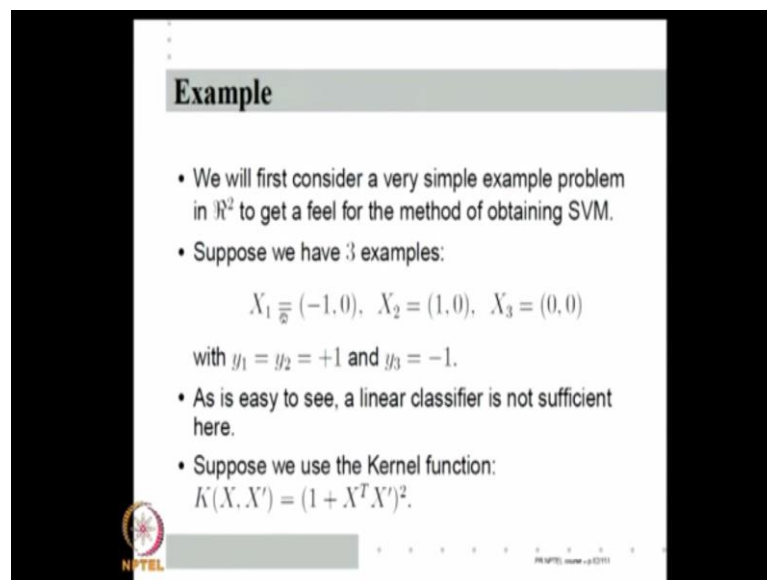
• A quadratic programming (QP) problem with interesting structure.

So, this is the dual. Let us remember this is the dual we want to maximize over mu some objective function which I called q mu which is equal to summation i is equal to one to n mu i minus half mu i mu j y i y j K X i X j this like a quadratic form of the K X i X j's. So, if I got X one to X n as the data vectors and I make a matrix whose i j'th element is K X a X j mu i mu j K X i X j is nothing but, the quadratic form of the matrix with respect to the vectors mu the vector mu and we have to also of course, multiply this matrix by y i y j the i j'th element has to multiply by either plus one minus one remember y i y i are plus one minus one. So, the product is either plus or minus depending on whether both X i and X j in the same class or different classes.

So, this essentially a quadratic form. So, this a quadratic cost function subject to one linear inequality constraints and bound constraints and variables. So, quadratic

programming problem with very interesting structure and hence there are many special algorithms. So, we would not consider any algorithm details I will just mention where this specialized algorithms comes from but, before going there we will solve one simple problem just to get a feel for how one solves for an SVM.

(Refer Slide Time: 42:17)

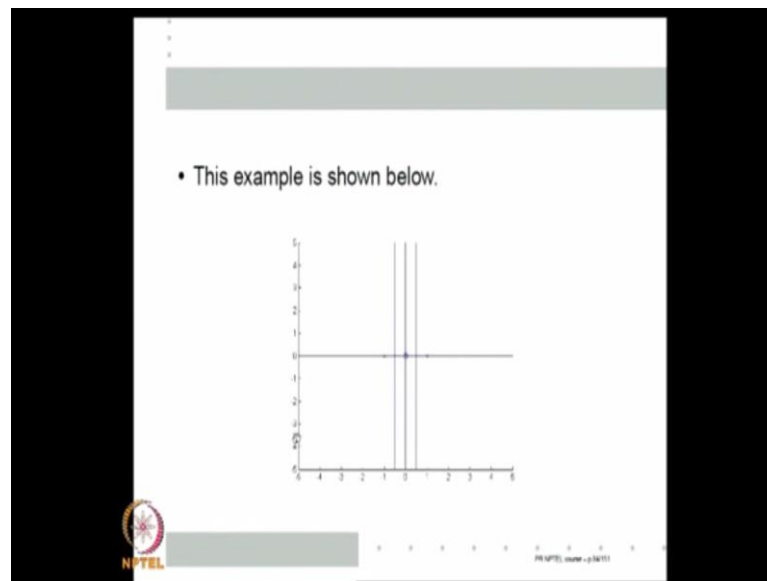


Example

- We will first consider a very simple example problem in \mathbb{R}^2 to get a feel for the method of obtaining SVM.
- Suppose we have 3 examples:
$$X_1 = (-1, 0), \quad X_2 = (1, 0), \quad X_3 = (0, 0)$$
with $y_1 = y_2 = +1$ and $y_3 = -1$.
- As is easy to see, a linear classifier is not sufficient here.
- Suppose we use the Kernel function:
$$K(X, X') = (1 + X^T X')^2.$$

So, we will first consider a very simple example in \mathbb{R}^2 let us say only 3 examples in \mathbb{R}^2 minus 1,0,1,0,0,0 minus 1,0 and 1,0 are in class plus 1,0,0 is in class minus matter of fact they are all of them are y component 0's is actually one dimensional problem. So, we call it 2 dimensional problem. So, is a very simple thing for all that is simple a linear classifier still not sufficient.

(Refer Slide Time: 42:48)



Because, see this my problem. So, this is 1 class the 2 dots or 2 stars are at the same class and the center circle is a different class. So, this is not linearly separable there in a line as you already seen 3 points which are collinear cannot be shattered by a hyperplane is the middle point is one class so, I have chosen that so, the middle point is 1 class the outer 2 other 2 classes. So, for all that is a very simple three point thing it is still not solvable by linear classifier.

So, let us say we chose a kernel function of this kind which essentially give us second degree polynomial in the original space so, this is my thing for example, I know one classifier is this I just put now the origin at that 0,0,3 in just to once again see. So, basically if I have some line parallel to y axis on either side. So, the thing is in between the lines is one class outside the lines other class so, it is like a pair of lines kind of a classifier. So, essentially if I square the x coordinate than this 2 will become same. So, I can actually get a threshold on the x coordinate square as a classifier. So, its a quadratic classifier but, of course, we know this but, my SVM has to automatically fine instead of I just chose a kernel function and go ahead and solve for the SVM.

(Refer Slide Time: 44:20)

• Recall, the examples are

$$X_1 = (-1, 0), \quad X_2 = (1, 0), \quad X_3 = (0, 0)$$

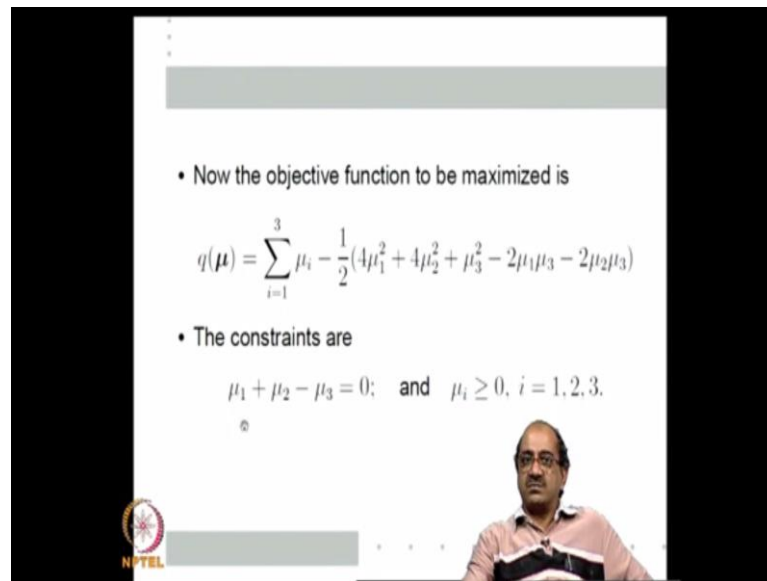
• The objective function involves $K(X_i, X_j)$. These are given in a matrix below.

$$[(1 + X_i^T X_j)^2] = \begin{bmatrix} 4 & 0 & 1 \\ 0 & 4 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

NPTEL

So, how do I solve for the SVM and it know the to write the dual as $K X_i X_j$ n a. So, let us first calculate $K X_i X_j$ $K X_i X_j$ is 1 plus X_i transpose X_j whole square we can write that as a matrix notation these are the X_i 's. So, 1 plus X_1 transpose X_1 whole square X_1 transpose X_1 is 1. So, 1 plus 1 is 2 square is 4 X_1 transpose X_2 is 0 sorry is minus 1. So, minus 1 plus 1 is 0. So, K 1 plus X_1 transpose X_2 is 0 X_1 transpose X_3 is 1,0. So, this becomes 1 and similarly, for X_2 transpose X_2 is once again 1. So, the Kernel function is 4 and X_2 transpose X_3 ,0. So, the Kernel function is 1 and finally, X_j transpose X is 0. So, this Kernel function is 1 by this obviously this is a symmetric matrix.

(Refer Slide Time: 45:24)



• Now the objective function to be maximized is

$$q(\mu) = \sum_{i=1}^3 \mu_i - \frac{1}{2}(4\mu_1^2 + 4\mu_2^2 + \mu_3^2 - 2\mu_1\mu_3 - 2\mu_2\mu_3)$$

• The constraints are

$$\mu_1 + \mu_2 - \mu_3 = 0; \quad \text{and} \quad \mu_i \geq 0, \quad i = 1, 2, 3.$$

The slide also features an NPTEL logo in the bottom left corner.

So, what will be mu dual now? My dual is q mu of the objective function is i is equal to 1 3 mu i plus the quadratic form on this. So, 4 mu 1 square plus 4 mu 2 square plus mu 3 square than I have all the class terms mu 1 mu 2 has coefficient 0 of course, a way mu 1 mu 3 is 2 mu 2 mu 3 is also 2 but, both mu 1 mu 3 and mu 2 mu 3 will be minus terms because 1 and 3 are a different classes and 2 and 3 also in different classes. So, it will be 4 mu 1 square 4 mu 2 square mu 3 square minus 2 mu 1 mu 3 minus 2 mu 2 mu 3 this is my objective function. What am my constraints? mu i y i is equal to 0. So, that is mu 1 plus mu 2 minus mu 3 mu 1 and mu 2 are in class plus x y 1 and y 2 is plus 1 y 3 is minus 1. So, mu 1 plus mu 2 minus mu 3 is equal to 0 and mu's are positive this is my problem to be solve. So, let us form.

So, this is the dual but, does not really matter to us this is just a constraint optimization problem with this at the objective function these are the constraints. So, I can once again use Kuhn-Tucker conditions. So, I have to form the Lagrangian. Lagrangian is this q mu plus Lagrange multipliers when constraints allow one mag Lagrange multiplier for this constraint. Let us call that lambda. I have some 3 Lagrange multipliers I have to write these constraints as minus mu i less than equal to 0.

(Refer Slide Time: 46:54)

- The lagrangian for this problem is
$$L(\mu, \lambda, \alpha) = q(\mu) + \lambda(\mu_1 + \mu_2 - \mu_3) - \sum_{i=1}^3 \alpha_i \mu_i$$
- Using Kuhn-Tucker conditions, we have $\frac{\partial L}{\partial \mu_i} = 0$ and $\mu_1 + \mu_2 - \mu_3 = 0$.
- This gives us four equations; we have 7 unknowns. We use complementary slackness conditions on α_i .
- We have $\alpha_i \mu_i = 0$. Essentially, we need to guess which $\mu_i > 0$.

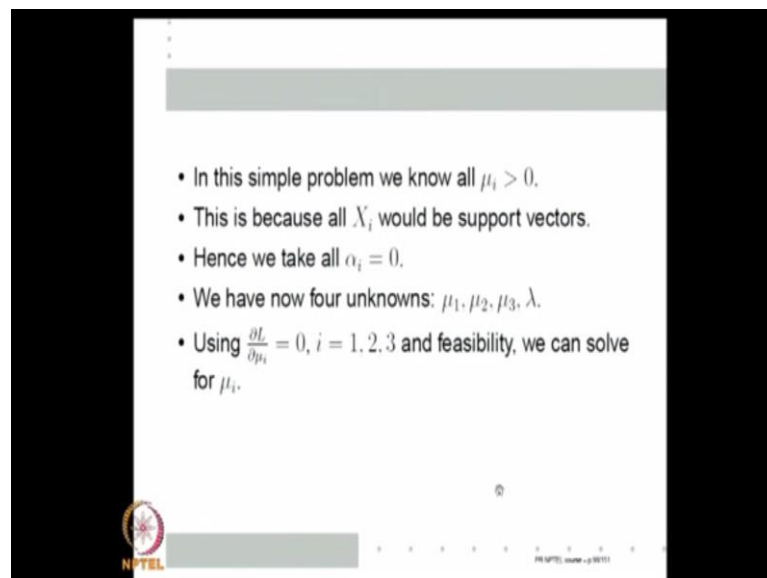
So, ultimately I get that as my Lagrangian $q(\mu) + \lambda(\mu_1 + \mu_2 - \mu_3) - \sum \alpha_i \mu_i$ where α_i 's are the 3 Lagrangian multipliers for these 3 constraints that is my Lagrangian. Now, to solve for solve this using Kuhn-Tucker condition. what are the Kuhn-Tucker condition? I have my optimization variables are μ_1, μ_2, μ_3 . So, I can take partial derivatives of L with respect to μ_1, μ_2, μ_3 is equal to 0 that gives me 3 equations I know μ 's have to be feasible.

So, they have to satisfy all the constraints. So, $\mu_1 + \mu_2 - \mu_3$ should be equal to 0 that's a constraint. So, that I get 4 so, this gives us 4 equations how many unknown I have $\mu_1, \mu_2, \mu_3, \lambda$ and 3 α 's I got 7 unknowns. So, you need 3 more equations have somewhere where do we get the extra equations I have to use the inequality constraints and we can use the complementary slackness on α_i 's. What does complementary slackness α_i ? I say α_i 's correspond to the constraints $\mu_i \geq 0$. So, the complementary slackness tells us that $\alpha_i \mu_i$ is equal to 0 which means if μ_i is strictly greater than 0 then α_i is 0 otherwise μ_i equal to 0. So, μ_i is equal to 0 I do not need to know μ_i and if μ_i strictly greater than 0 α_i is equal to 0 I do not need to know α_i .

So, because that 3 complementary slackness is essentially they give me the remaining three equations essentially in the sense they will tell me which α_i 's are μ_i 's are 0 but, I really do not know which are 0 I have to actually in a general constrained

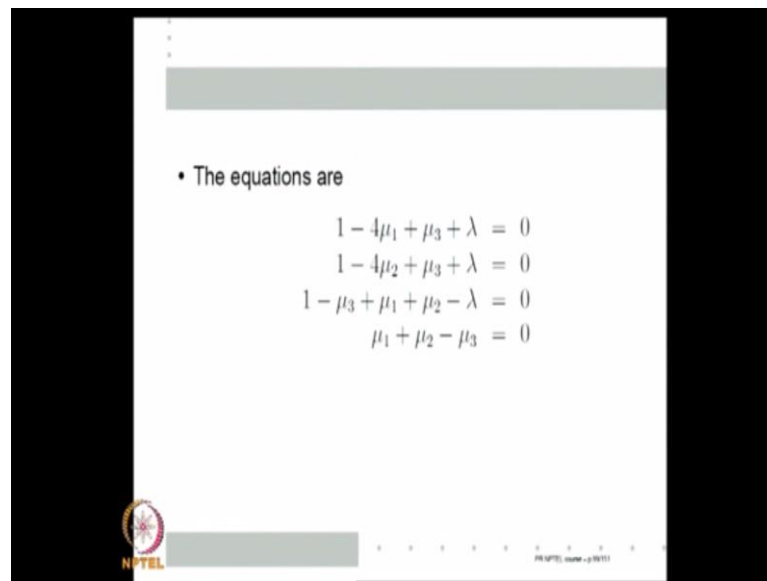
optimization problem like this I have to make a guess work and which μ_i 's are greater than 0. So, they can take the corresponding α_i to be 0 and solve this what 1 normally does if you are solving a constrained optimization problem in pen and paper is to make this various guesses like this and then see which guess gives rise to a consistent solution in SVM's in general in this particular type problem but, in SVM's in general this is very simple to make because we know that μ_i greater than 0 has a physical interpretation a μ_i is greater than 0 the corresponding X_i is a support vector that means is closest to the separating bound. So, I can actually geometrically think of which excess have to at least in \mathbb{R}^2 I can do this in this problem of course, we know all the 3 are support vectors.

(Refer Slide Time: 49:27)



So, we know that all the μ_i will be greater than 0. So, we can take all α_i to be 0 that means we have only 4 unknowns and we got 4 equations. So, we can solve for it.

(Refer Slide Time: 49:42)



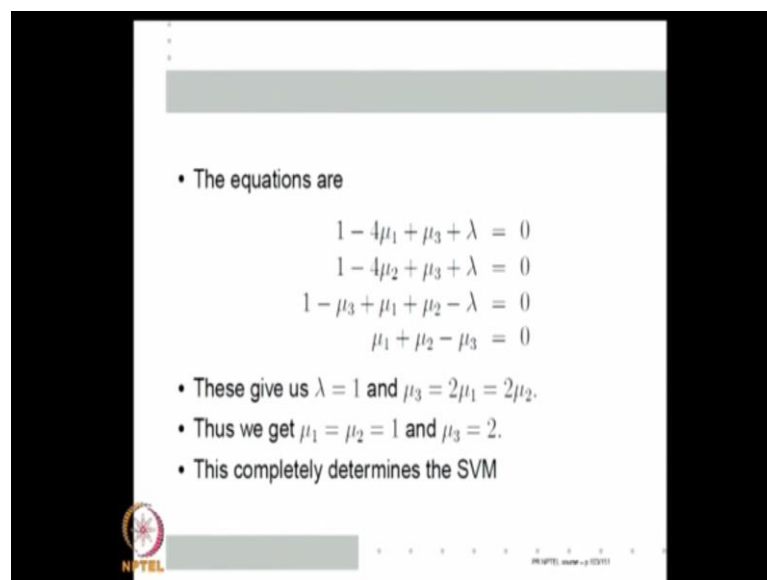
• The equations are

$$\begin{aligned}1 - 4\mu_1 + \mu_3 + \lambda &= 0 \\1 - 4\mu_2 + \mu_3 + \lambda &= 0 \\1 - \mu_3 + \mu_1 + \mu_2 - \lambda &= 0 \\\mu_1 + \mu_2 - \mu_3 &= 0\end{aligned}$$

NPTEL

So, if I put this 4 equations down that will happen to be the 4 equations it is easy enough to see. Essentially, I am differentiating this if I differentiate this could to see when differentiate $q \mu_i$ get 1 from this because in y and a μ term from this the square and then the μ_2 or μ_3 term from the cross terms. And then, from here I will get a λ term α is 0.

(Refer Slide Time: 50:23)



• The equations are

$$\begin{aligned}1 - 4\mu_1 + \mu_3 + \lambda &= 0 \\1 - 4\mu_2 + \mu_3 + \lambda &= 0 \\1 - \mu_3 + \mu_1 + \mu_2 - \lambda &= 0 \\\mu_1 + \mu_2 - \mu_3 &= 0\end{aligned}$$

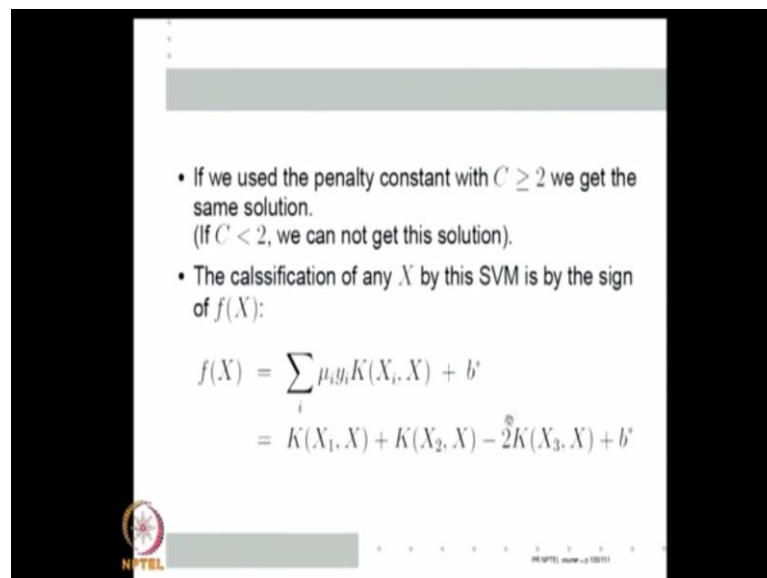
- These give us $\lambda = 1$ and $\mu_3 = 2\mu_1 = 2\mu_2$.
- Thus we get $\mu_1 = \mu_2 = 1$ and $\mu_3 = 2$.
- This completely determines the SVM

NPTEL

So, roughly one can see that that the general structure of the equations it that just happened because you can differentiate and see now 1 and this particular case is chosen.

So, the equations are very simple to solve we already can see that $\mu_1 + \mu_2 - \mu_3 = 0$ I got $\mu_1 + \mu_2 - \mu_3$ here if I put that equal to 0. I get λ equal to 1 and looking at these 2 equations. if I subtract 1 from the other I know μ_1 is equal to μ_2 and from here I know μ_3 is equal to $\mu_1 + \mu_2$. So, I know μ_3 is equal to $2\mu_1$ is equal to $2\mu_2$. So, μ_1 is equal to μ_2 and μ_3 is twice of them once I put that in I its absolutely straightforward to see that the solution is μ_1 is equal to μ_2 is equal to 1 and μ_3 is equal to 2. This completely does in the SVM that is all the SVM is I have all the non0 Lagrange multipliers and the corresponding support vectors.

(Refer Slide Time: 51:15)



- If we used the penalty constant with $C \geq 2$ we get the same solution.
(If $C < 2$, we can not get this solution).
- The classification of any X by this SVM is by the sign of $f(X)$:

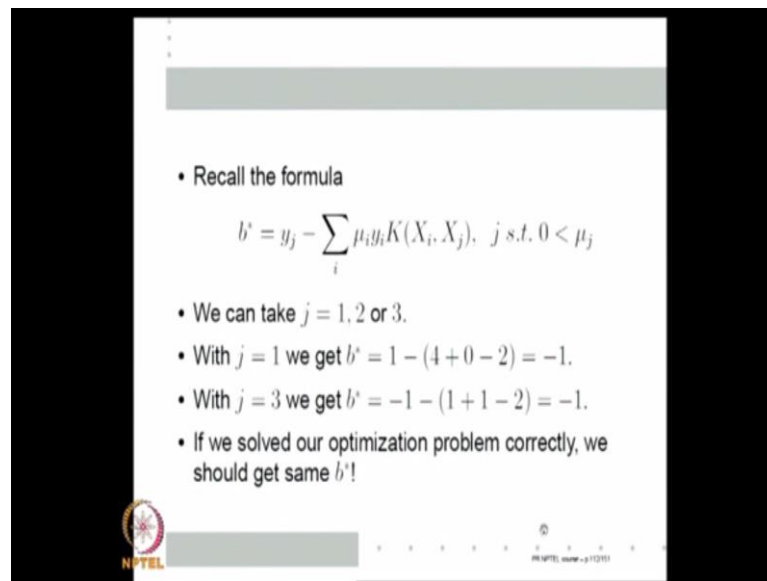
$$f(X) = \sum_i \mu_i y_i K(X_i, X) + b^*$$

$$= K(X_1, X) + K(X_2, X) - 2K(X_3, X) + b^*$$

Just one thing I am not uses the penalty constant in this formulation because the high problem but, if I use the penalty constant with greater than 2. I would get the same solution because using the penalty constant simply means that in the I have an extra constraint that says that all μ 's have to between 0 and 2 even with unconstrained my μ 's are utmost 2. So, if C is greater than or equal to 2 the penalty constant would not make a difference the other hand by mistake we change should be less than 2 we will not get the solution we may get a bad solution will lets not go there. So, the classification by this SVM for any new vector x is determines the sign of $f(X)$ where $\mu_i y_i K(X_i, X) + b^*$ lets first calculate b^* to understand what this is. So, this is what $\mu_i y_i K(X_i, X)$ is $K(X_1, X) \mu_1$ is 1 y_1 is 1 $K(X_2, X) \mu_2$ is one-second is 1 minus 2 $K(X_3, X)$ because μ_3 is 2 and y_3 is minus 1 plus b^* i because i know X_1, X_2 and X_3 . If you give me

X_i can calculate all this.

(Refer Slide Time: 52:23)



- Recall the formula


$$b^* = y_j - \sum_i \mu_i y_i K(X_i, X_j), \quad j \text{ s.t. } 0 < \mu_j$$

- We can take $j = 1, 2$ or 3 .
- With $j = 1$ we get $b^* = 1 - (4 + 0 - 2) = -1$.
- With $j = 3$ we get $b^* = -1 - (1 + 1 - 2) = -1$.
- If we solved our optimization problem correctly, we should get same b^* !


So, let us first calculate b^* you remember b^* is given by this y_j minus $\mu_i y_i K(X_i, X_j)$ for any j 's such that μ_j strictly greater than 0 in this case all three μ_j 's strictly greater than 0. So, if I put j is equal to 1 what do I get b^* is y_1 y_1 is 1 minus that 3 terms here first 2 terms will be positive second term negative because $y_1 y_2$ is positive y_3 is negative μ_1 is 1 $K(X_1, X_1)$ is 4 that is $4 - 0 + 1 = 3$ and the other one is 1 I get a 2 from μ_3 and minus from. So, that gives me minus 1.

I do not have to choose j is equal to 1 any j will do if j is equal to 3 y_3 is minus 1 with respect to that I once again put that in this I get minus from it obviously basically if we solved our optimization problem correctly no matter which j such that μ_j star μ_j is greater than 0 I chose I should get the same b^* matter of fact this is one way I can check whether we have solve the optimization problem correctly but, when we do it numerically this may not work. So, we may take b^* to be the average of all b^* obtained with different j 's the minor issue.

(Refer Slide Time: 53:37)




- We have $X_1 = (-1, 0)$, $X_2 = (1, 0)$, $X_3 = (0, 0)$ and $K(X, X') = (1 + X^T X')^2$.
- Hence, taking $X = (x_1, x_2)^T$, we have

$$\begin{aligned}
 f(X) &= K(X_1, X) + K(X_2, X) - 2K(X_3, X) + b^* \\
 &= (1 - x_1)^2 + (1 + x_1)^2 - 2(1) - 1 \\
 &= 2x_1^2 - 1
 \end{aligned}$$


So, we have X_1, X_2, X_3 is this and we know $K(X, X')$ is $1 + X^T X'$ whole square. So, if I take X to be X_1, X_2 I know my function is K of X_1, X plus K of X_2, X minus $2K$ of X_3, X plus b^* . What will be X_1, X is this. So, this is $1 - X_1$ whole square. So, $1 - X_1$ whole square X_2 is this. So, it will be $1 + X_1$ whole square X_3 is this. So, X_3 transpose X is always 0. So, 2 into 1 b^* is minus 1. So, $1 - X_1$ whole square plus $1 + X_1$ whole square will give me $2X_1$ square plus 2 that plus 2 will cancel this minus 2 I get minus 1. So, my $f(X)$ is $2X_1$ square minus 1.


(Refer Slide Time: 54:31)



- Hence this SVM will assign class +1 to $X = (x_1, x_2)^T$ if

$$2x_1^2 \geq 1 \quad \text{or} \quad |x_1| \geq \frac{1}{\sqrt{2}}$$

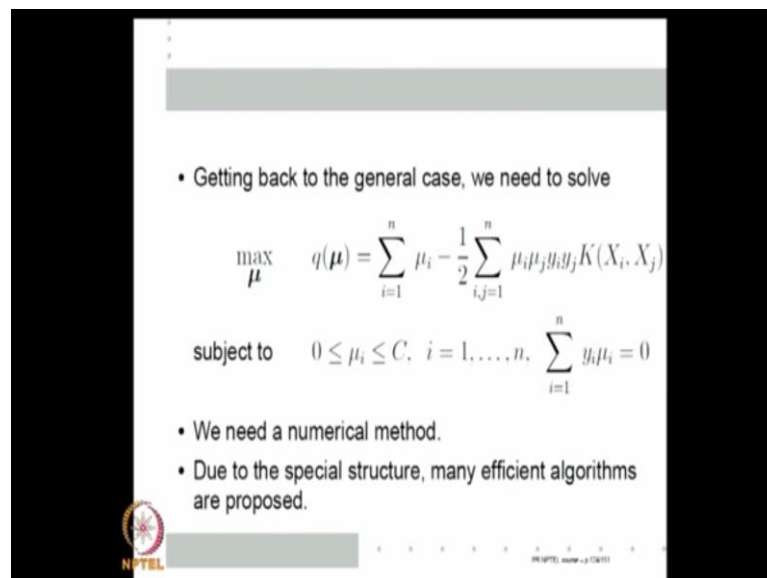
- Why not $|x_1| \geq (1/2)$?
- We are maximizing margin of the hyperplane in ' x^2 '-space.
- The final SVM is intuitively very reasonable and we solve essentially the same problem whether we are seeking a linear classifier or a nonlinear classifier.



Which mean, this SVM will assign class plus 1 to X any given X if $X^T X$ is greater than 1 or $X^T X$ is greater than $\sqrt{2}$ this exactly what we as we said the square of the first thing should be greater than some threshold that is the SVM that is magically come over we did nothing we just chose a Kernel function. Now, you may say it should not the classifier be mode $X^T X$ greater than half. So, that i i put it exactly between the 0 and minus 1 or 0 1 plus 1 that is what we are asking by maximizing margin well, we shall understand that we are maximizing because we are using a kernel function where maximizing margin in the quadratic space in the $X^T X$ space.

So, in the $X^T X$ space I wanted between 0 1 plus 1 that is at half. So, in the $X^T X$ space will be 1 by root 2 does we maximizing margin in the transforms space. So, the final SVM is intuitively very reasonable and we solve exactly the same problem whether we are seeking a linear classifier or a nonlinear classifier. So, this is a good example, so we understand how to solve SVM's.

(Refer Slide Time: 55:47)



- Getting back to the general case, we need to solve

$$\max_{\mu} \quad q(\mu) = \sum_{i=1}^n \mu_i - \frac{1}{2} \sum_{i,j=1}^n \mu_i \mu_j y_i y_j K(X_i, X_j)$$

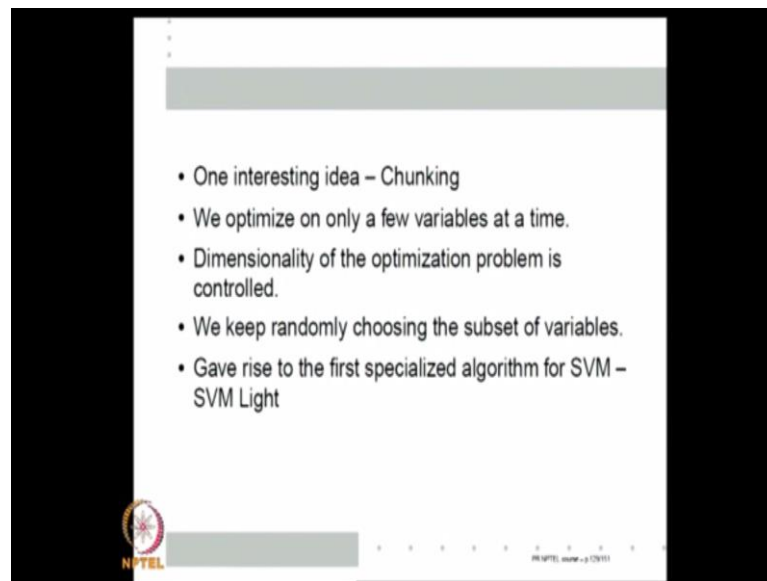
subject to $0 \leq \mu_i \leq C, \quad i = 1, \dots, n, \quad \sum_{i=1}^n y_i \mu_i = 0$

- We need a numerical method.
- Due to the special structure, many efficient algorithms are proposed.

NPTEL

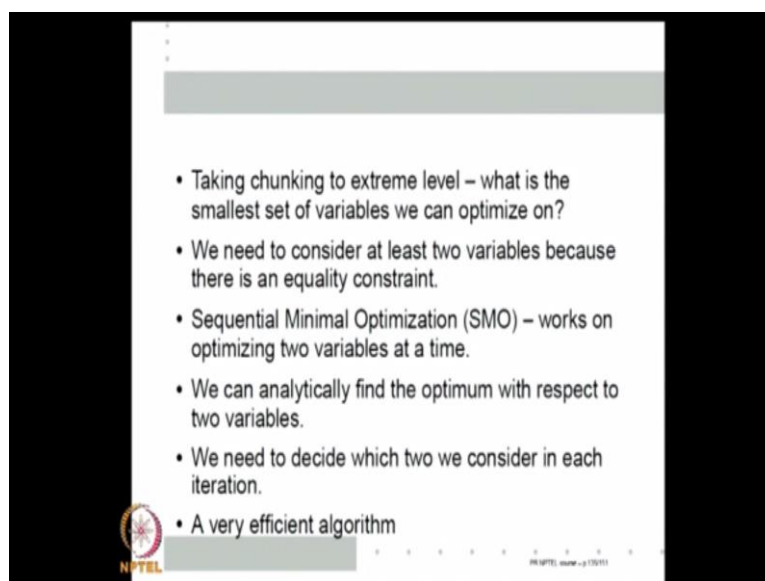
Now, let us move on and look at the general case this is what you want to solve and we need a numerical method as I said due to the special structure there are many efficient algorithms.

(Refer Slide Time: 56:00)



We essentially look at only one structure here what is called Chunking. Chunking is simply that when you have a optimization of many variables we do optimization a few variables at a time something that all of you might remember is what is called coordinate rise minimization if you have function of $X_1 X_2 \dots X_n$. We first keep X_2 to X_n fix and minimize with respect X_1 as one dimensional problem than minimize with respect to X_2 than minimize with respect to X_3 and cyclically keep doing it. Lets has simple idea of Chunking we optimize on only a few variables at a time which allows the dimensionality of each optimization problem to be small and at each iteration we randomly choose a subset of variables of course, we have to use some proper heuristic to decide how to choose as long as all variables keep coming into the set being optimized again and again essentially chunking works. This basic idea gave rise to the first specialized and very efficient algorithm for solving SVM which is called SVM light.

(Refer Slide Time: 56:59)



And then, I can take the I want to take Chunking to the extreme level SVM light takes 5 or 10 or 15 variables at a time but, I want to take chunking to the extreme level I should ask what is the smallest chunk that I can optimize on because SVM has one equality constraint I can optimize at most on 2 variables at a time I cannot do change one variable keeping all the others fix because if I met a feasible point and change only one variable I go to infeasible point because there is a equality constraint.

So, there is a special algorithm call sequential minimal optimization which works on this principle of optimizing two variables at a time basically this very interesting algorithm because if I keep rest of them fix And I will say i will optimize only μ_1 μ_2 than by analytically manipulating the expression I can calculate the optimum with respect with respect to μ_1 μ_2 which means I do not have to do any numerical optimization at all I spend all my time deciding which two variable which 2 Lagrange multipliers I should do optimization next time. So, using a variety of heuristic this algorithm does a very good job of solving the SVM problem is a very efficient algorithm. So, I will just re talk about s this SMO and if and a couple of other issues in solving the SVM problem in the next class.

Thank you.