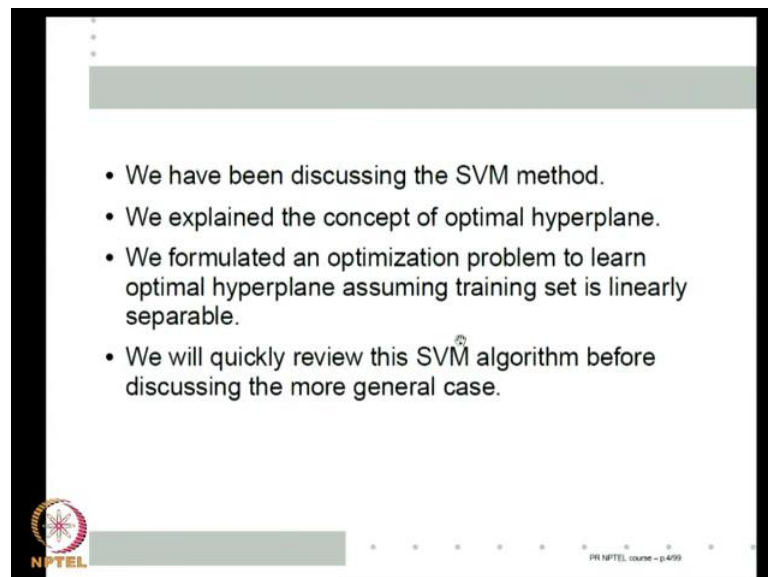


Pattern Recognition
Prof. P. S. Sastry
Department of Electronics and Communication Engineering
Indian Institute of Science, Bangalore

Lecture - 33
SVM formulation with slack variables; non-linear SVM classifiers

Hello and welcome to this next lecture on pattern recognition, we have started discussing the support vector machines from last class. As you already discussed this is an approach where the idea is that we effectively transform the feature vectors into high dimensional space and learn a linear classifier there, because you already know how to learn a linear classifier and learning linear classifier is in general more efficient. We want to learn a linear classifier in a high dimensional space, which in the original space will be a non-linear classifier that is the basic idea. We said that to for the idea to work, we actually do not learn any linear classifier, but learn what we called an optimal separating hyperplane. So, to briefly review we have been discussing the SVM method.

(Refer Slide Time: 01:08)



The one major new concept is the optimal hyperplane, right? We formulated the optimization problem of how one learns optimal hyperplane. We seen that learn for learning optimal hyperplane, we essentially have to solve a quadratic optimization problem with linear constraints. Basically, the optimal hyperplane is one which maximizes the separation and we looked at the formulation only under the assumption that the training set is linearly separable that is the case we have been considering so far.

So, what we will do in this class is we quickly review the, we start with the optimization formulation again quickly go through the dual and then move on and see how we can do the same formulation in general even when the training set is not linearly separable

(Refer Slide Time: 02:12)

The optimization problem for SVM

- The optimal hyperplane is a solution of the following constrained optimization problem.
- Find $W \in \mathbb{R}^m, b \in \mathbb{R}$ to

$$\begin{aligned} &\text{minimize} && \frac{1}{2}W^T W \\ &\text{subject to} && 1 - y_i(W^T X_i + b) \leq 0, \quad i = 1, \dots, n \end{aligned}$$

- Quadratic cost function and linear (inequality) constraints.
- Kuhn-Tucker conditions are necessary and sufficient. Every local minimum is global minimum.

NPTEL PR NPTEL course - p1599

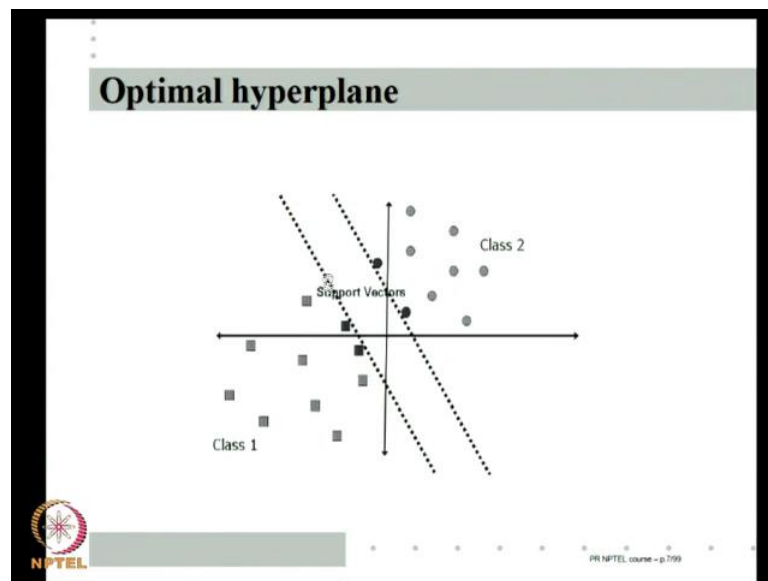
So, the optimal hyperplane is a solution to this optimization problem, what is the optimization? It is a constraint optimization problem namely to minimize transpose W subject to 1 minus y i into W transpose X i plus b less than or equal to 0. Optimization variables are W and b, find W and b to minimize half W transpose W subject to this constraints. What does this constraint imply? As we have already seen if W b, W comma b represents a separating hyperplane, then y i into W transpose X i plus b will be equal to 1, we will be greater than or equal to 1.

So, essentially y i into W transpose X i plus b greater than or equal to 1, which is same as 1 minus y i into W transpose X i plus b less than or equal to 0, tells you that this W b is a separating hyperplane and the separating hyperplane is such that W transpose X i plus b is equal to minus 1 and W transpose X i plus b is equal to plus 1. The two parallel hyperplanes that we considered last class no pattern is in between, so a separating hyperplane which, for which the margin of separation is given by 1 by norm w, that is any W b that satisfies this constraint will be a separating hyperplane, for which the closest pattern is 1 by norm W away.

Hence among all such W b if I minimize this, then I am finding a separating hyperplane, which has the maximum margin of separation meaning the closest training pattern is far away from for this hyperplane compared to any other separating hyperplane. So, that is why it is a constrained optimization problem among all separating hyperplanes that is a all W which satisfy this constraint, find the one which has the maximum margin of separation this is the optimization problem.

So, the optimal hyperplane is a solution to this problem and this is a, this this is the cost function, so which is quadratic, convex quadratic, all constraints are linear. So, this is the most efficient constrained optimization problem one can solve, for this problem because this is convex and constraints are linear, Kuhn Tucker conditions are (()) and sufficient every local minimum is a global minimum. We have seen all the background of constraint optimization last class, so let us quickly write down the Kuhn-Tucker conditions for this problem.

(Refer Slide Time: 04:40)



Before that let me give you the as I have been telling you that these constraints essentially mean that you know if the hyperplane is this at the centre, then on either side i have got these two parallel hyperplanes and this is the margin of separation. So, what the optimal hyperplane is doing is that finding that hyperplane which has the highest margin of separation, that is what my optimization problem. I do not want a separating hyperplane like this, but I want a separating hyperplane like that, a separating hyperplane

for which the closest pattern is more distance away and this. So, this is a non optimal hyperplane whereas, that is a optimal hyperplane. So, what is what is my Lagrangian for this problem? So, Lagrangian is my cost function plus sum of all the constraints you have to Lagrangian multiplies into constraints.

(Refer Slide Time: 05:38)

The slide contains the following text and equations:

- The Lagrangian is given by

$$L(W, b, \mu) = \frac{1}{2}W^T W + \sum_{i=1}^n \mu_i [1 - y_i(W^T X_i + b)]$$
- The Kuhn-Tucker conditions give

$$\nabla_W L = 0 \Rightarrow W^* = \sum_{i=1}^n \mu_i^* y_i X_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \mu_i^* y_i = 0$$

$$1 - y_i(X_i^T W^* + b^*) \leq 0, \quad \forall i$$

$$\mu_i^* \geq 0, \quad \& \quad \mu_i^* [1 - y_i(X_i^T W^* + b^*)] = 0, \quad \forall i$$

The slide also features the NPTEL logo in the bottom left corner and the text '© NPTEL course - p 1299' in the bottom right corner.

So, my Lagrangian will be half W transpose W plus sum over i mu i into 1 minus y I W transpose X i plus b, mu i are the Lagrangian multipliers. So, it is a function of W b the optimization variables and all the mu i which are the Lagrangian multipliers. So, what are my Kuhn Tucker conditions? The gradient of L with respect to the optimization variable should be 0, so if I equate gradient with respect to W to 0, if I take W k I get W here and from here there is one linear W term, so i get mu i y i X i. So, essentially at the optimal w, W star is summation mu i y i X i for all the optimal variables we put star. So, mu i, star is the optimal values of the Lagrangian multipliers, W star is the optimal value of W, this the normal to the separating hyperplane.

So, from gradient of L with respect to W is 0, gives me that W star is this, gradient of L with respect to b 0, gives me there is only one b term, but that is linear. So, I get i is equal to 1 to n mu i y i into b, so differentiate it I get summation i is equal to 1 to n mu i y i is equal to 0. This is the condition that the gradient of the Lagrangian with respect to optimization variables is 0, then we need feasibility, so all constraints should be satisfied. So, 1 minus y i into X i transpose W star plus b star is less than or equal to 0, for all i and

then I have complementary slackness, so I first I want all the Lagrangian multipliers to be positive or non negative and μ_i^* into the constraint.

As I have already mentioned last class our optimization problem is such that there is one constraint for each example, essentially because I want $W^* b$ to be a separating hyperplane. For each $X_i y_i$, each example $X_i y_i$ it should be on the right side of the W hyperplane, so each constraint is straight to one example, each example has one constraint, so each constraint is a corresponding X_i . In the, in my Lagrangian each constraint will have its own Lagrange multiplier, so my each Lagrange multiplier is tied to a particular example.

So, I can talk of Lagrange multiplier for that example, because each Lagrange multiplier is for a constraint and each constraint is for an example. Now, this complementary slackness gives us something very interesting as we seen last class. Let us say S is the set of all indices for which these corresponding Lagrange multipliers are strictly positive, you know Lagrangian multipliers have to be greater than or equal to 0, some of them might be 0, some of them might be strictly positive, S may be the index such that Lagrange multipliers are strictly positive.

(Refer Slide Time: 08:13)

The slide contains the following text and equations:

- Let $S = \{i \mid \mu_i^* > 0\}$.
- By complementary slackness condition,

$$i \in S \Rightarrow y_i(X_i^T W^* + b^*) = 1$$
 Implies X_i is closest to separating hyperplane.
- $\{X_i \mid i \in S\}$ are called Support vectors. We have

$$W^* = \sum_i \mu_i^* y_i X_i = \sum_{i \in S} \mu_i^* y_i X_i$$
- Optimal W is a linear combination of Support vectors.
- Support vectors constitute a very useful output of the method.

The slide also features the NPTEL logo in the bottom left corner and the text "© NPTEL course - p.15/59" in the bottom right corner.

Then the complementary slackness says if this is strictly, this strictly greater than 0 then this has to be 0, because the product has to be 0. So, if i belongs to S then y_i into X_i transpose $W^* + b^*$ is equal to 1, what is that this mean? This means that the the

X_i is on the hyperplane $W^* \text{ transpose } X_i, X + b$ is equal to plus 1 or minus 1 depending on whether y is plus 1 or minus. So, is on one of the two parallel hyperplanes, which means that X_i is a pattern that is closest to the separating hyperplane. So, the X_i corresponding to the non zero Lagrangian multipliers or the training patterns which are closest to the separating hyperplane, right? The complementary slackness says if μ_i^* is greater than 0 then this has to be 0, right?

So, if I think of S as the set of all Lagrange multipliers, in this of all Lagrange multipliers which are strictly positive, then for all those i I know this is satisfied which means all the corresponding X_i are closest to the separating hyperplane. We will call such X_i corresponding to positive Lagrangian multipliers as support vectors, these are training patterns they are closest to the separating hyperplane. These are the my optimal W^* I already know is summation or $\sum \mu_i^* y_i X_i$, so if this summation need be taken only over those μ_i^* which are strictly greater than 0. So, essentially my W^* is a linear combination of these support vectors, if I define this as the support vectors.

Also another thing that should remembered is that for my optimization problem, if μ_i^* is equal to 0 then the constraint is not active, so which means essentially these X_i, X_i corresponding to the support vectors are the critical X_i . Because if I take away all the other examples and leave only these I still get the same separating hyperplane, the optimization problem will not change. That is why W^* is a linear combination of support vectors, we have already seen while discussing perceptrons that any separating hyperplane will be a linear combination of X_i , but for the optimal separating hyperplane is a linear combination of some of the X_i not all, only the what we call the support vectors and the linear combination involves the corresponding Lagrange multipliers.

So, optimal W is a linear combination of support vectors and in this sense support vectors constitute a very useful output of the method, right? For this problem so these are the support vectors, the pattern that are darken they are closest to the supporting. So, as you can see geometrically if I throw away all the remaining data and just keep the support vectors right for that pattern classification problem also this will be the best hyperplane. So, in some sense support vectors are the most critical examples of this pattern set of this example set, thus as we shall see later on much more clearly, support

vectors constitute a very useful extra output of the method in addition to getting the W star like this.

(Refer Slide Time: 11:57)

The SVM solution

- The optimal hyperplane – W^*, b^* given by:

$$W^* = \sum_i \mu_i^* y_i X_i = \sum_{i \in S} \mu_i^* y_i X_i$$

$$b^* = y_j - X_j^T W^*, \quad j \text{ s.t. } \mu_j^* > 0$$
 (Note that $\mu_j^* > 0 \Rightarrow y_j(X_j^T W^* + b^*) = 1$)
- Thus, W^*, b^* are determined by $\mu_i^*, i = 1, \dots, n$.
- We can use the dual of the optimization problem to get μ_i^* .

So, let us sum up the SVM solution, SVM tries to learn the optimal hyperplane, optimal hyperplane is defined to be one that maximizes the separation, is a separating hyperplane which has the maximum separation. The optimal hyperplane is given by W star transpose X plus b is equal to 0, then the W star and b star are given by this, W star is summation over i , μ_i star $y_i X_i$ and b star how do I know b star? From the complementary slackness, whenever μ_j star is greater than 0 I know y_j into X_j transpose W star plus b star is equal to 1. Because my complimentary slackness tells me that whenever the μ_i is Lagrange multiplier is 0 this is equal to 1.

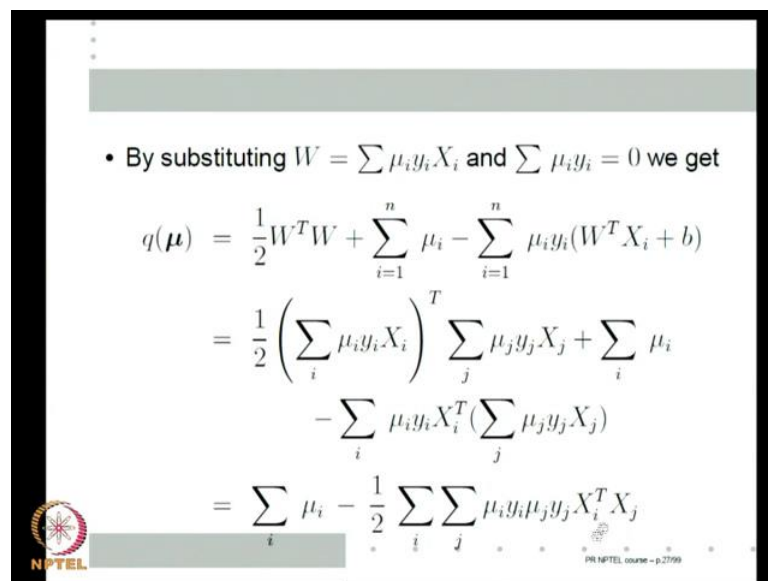
When this is 1 b multiply both side y_j , I remember y_j square is equal to 1 because y_j is plus 1 or minus 1, so that gives me b star is y_j minus X_j transpose W star right. So, once I have all the μ_i stars I can calculate W star and second calculate W star and have all the μ_i stars I can take any j at that μ_j star is greater than 0 and b star will be this. Theoretically does not matter which j you take you should get the same value of b star, so that W star and b star are determined by your μ_i stars your Lagrangian multipliers.

So, if I can get the optimal Lagrangian multipliers for this problem I am done and because I essentially want to get Lagrangian multipliers one way of doing solving this problem, one way of getting the optimal hyperplane is to find μ_i star by solving the

am maximizing $q(\mu)$ so a μ at which $q(\mu)$ takes value minus infinity is not a place where maximum will be attained. Hence we need to maximize q only over those μ which satisfy $\sum \mu_i y_i = 0$. I am not interested in finding value of $q(\mu)$ for those μ for which $\sum \mu_i y_i$ is not equal to 0 for those μ anyway I know is minus infinity.

So, we need to maximize q only over those μ where $\sum \mu_i y_i = 0$, so once I take $\sum \mu_i y_i = 0$ the b term drops out. Now, I need to only find the infimum of this with respect to W , that means I have to differentiate with respect to W and find the gradient with respect to W equal to 0 which we already done earlier. So, I know infimum with respect to W is attained at $W = \sum \mu_i y_i X_i$ because this X differentiate with respect to W and equate with 0. So, how do I get $q(\mu)$ for a given μ now? For that μ using that $W = \sum \mu_i y_i X_i$ and now I substitute in this expression $W = \sum \mu_i y_i X_i$ and impose $\sum \mu_i y_i = 0$ then I get $q(\mu)$.

(Refer Slide Time: 16:27)



• By substituting $W = \sum \mu_i y_i X_i$ and $\sum \mu_i y_i = 0$ we get

$$\begin{aligned}
 q(\mu) &= \frac{1}{2} W^T W + \sum_{i=1}^n \mu_i - \sum_{i=1}^n \mu_i y_i (W^T X_i + b) \\
 &= \frac{1}{2} \left(\sum_i \mu_i y_i X_i \right)^T \sum_j \mu_j y_j X_j + \sum_i \mu_i \\
 &\quad - \sum_i \mu_i y_i X_i^T \left(\sum_j \mu_j y_j X_j \right) \\
 &= \sum_i \mu_i - \frac{1}{2} \sum_i \sum_j \mu_i y_i \mu_j y_j X_i^T X_j
 \end{aligned}$$

NPTEL
PR NPTEL course - p.2799

Let us calculate $q(\mu)$, this is what we have to do, we have to substitute $W = \sum \mu_i y_i X_i$ and my $q(\mu)$ is this, now I am taking an infimum, because I know how what to do for infimum or just to do the substitutions. My $q(\mu)$ is $W^T W$ μ_i into 1 minus this I have written as $\sum \mu_i$ and the second term because this term $\sum \mu_i y_i W^T X_i + b$. So, this $\sum \mu_i y_i b$ will go to 0 because I have to put $\sum \mu_i y_i = 0$.

equal to 0, so I have to put $\sum_i \mu_i y_i X_i$. So, this $\sum_i \mu_i y_i X_i$ whole transpose, this $\sum_i \mu_i y_i X_i$ once again this i is a dummy index. So, I have put this $\sum_j \mu_j y_j X_j$ plus $\sum_i \mu_i y_i X_i$ minus $\sum_i \mu_i y_i X_i$ has gone, $\sum_i \mu_i y_i X_i$ I write this as $X_i^T \sum_j \mu_j y_j X_j$.

So, both these summations are same I have this summation of i and j here $\sum_i \mu_i y_i X_i^T X_j$ here also $\sum_j \mu_j y_j X_i^T X_j$ this is minus 1 this is half. So, I will get a minus half. So, $q(\mu)$ is $\sum_i \mu_i y_i X_i^T X_j$ minus half $\sum_{i,j} \mu_i \mu_j y_i y_j X_i^T X_j$ that is my $q(\mu)$.

(Refer Slide Time: 17:55)

• Thus, the dual problem is:

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

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

- Quadratic cost function and linear constraints
- Training data vectors appear only as innerproduct
- Optimization is over \mathbb{R}^n irrespective of the dimension of X_i .

NPTEL © NPTEL course - p.2999

So, my dual now turns out to be maximize $q(\mu)$ which is given by this as we have just now seen subject to $\mu_i \geq 0$ and we have to impose the extra condition $\sum_i y_i \mu_i = 0$, right? So, this is my dual problem this is also a constrained optimization problem, the optimization variables are μ , right? The cost function is this because my variables are μ this is a quadratic function; this term is linear in μ ; this term is quadratic in μ . So, this is a quadratic function, constraints are linear $\mu_i \geq 0$ $\sum_i y_i \mu_i = 0$.

So, it is once again a constrained optimization problem with quadratic cost function and linear constraints. The training data vectors $X_i X_j$ unlike in the primal way they appear inside the constraint they appear only in the objective function, there also only as an inner product $X_i^T X_j$, right? The optimization variables here are μ how many are

there n, right? As many examples 1 mu per example so there are n examples, so there are n Lagrange multipliers and so the dimensionally the optimization problem is n irrespective of the dimension of x, no matter what the dimension of X is dimension of this optimality is optimization problem is equal to the number of examples. We will just remember these things for later use.

(Refer Slide Time: 19:19)

Optimal hyperplane

- The optimal hyperplane is a solution of

$$\min_{W,b} \quad \frac{1}{2}W^TW$$

subject to $y_i(W^TX_i + b) \geq 1, \quad i = 1, \dots, n$

- Instead of solving this primal problem, we solve the dual and obtain optimal Lagrange multipliers, μ_i^* .

NPTEL logo and footer text: PR NPTEL course - p.3299

So, sum up the optimal hyperplane is a solution of this problem, this optimization problem we call this the primal problem. Instead of solving this primal problem we solve the dual and obtained optimal Lagrange multiplier.

(Refer Slide Time: 19:32)


• The dual problem is:

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

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

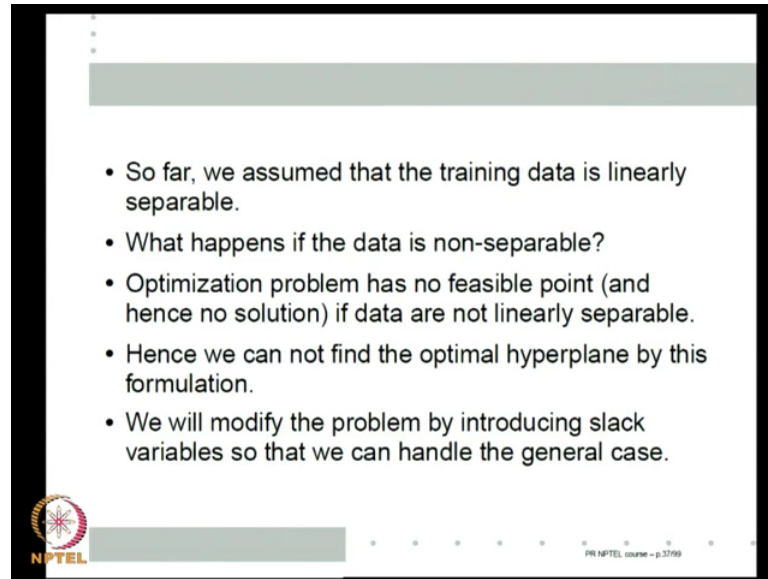
• Then the final solution is:

$$W^* = \sum \mu_i^* y_i X_i, \quad b^* = y_j - X_j^T W^*, \quad j \text{ such that } \mu_j > 0$$

 © NPTEL course - p.32/99

This is the dual problem, we solve this dual problem so when we solve this we get all the μ_i and once we get all the μ_i , this is my final solution. What do you mean by W^* b^* is my solution? Now, you give me any X then I calculate $W^* \text{ transpose } X + b^*$ if this is positive I put it in class plus 1, if it is negative I put it in class minus 1. So, $W^* b^*$ gives me my optimal hyperplane, so from now on give me any new feature vector X I will say I will calculate $W^* \text{ transpose } X + b^*$ and sign of that will give me the class. As I said we will we are considering only two class problem, so this is the full SVM method for linearly separable case. So, we define the optimal hyperplane to be this, then that is the dual and that is my final solution.

(Refer Slide Time: 20:23)



- So far, we assumed that the training data is linearly separable.
- What happens if the data is non-separable?
- Optimization problem has no feasible point (and hence no solution) if data are not linearly separable.
- Hence we can not find the optimal hyperplane by this formulation.
- We will modify the problem by introducing slack variables so that we can handle the general case.

Now, in all this we have assumed that the training data is linearly separable, what happens in training data if it is not separable? What do you mean by training data not separable? If the training data is not separable, what it means is there will not exist a W and b to satisfy these constraints. These constraints are satisfied by a W and b only if W and b are a separating hyperplane, if there does not exist a separating hyperplane because training data is not linearly separable, then there will not be a W and b that satisfies the constraint. That means there is no feasible, so feasible point for this optimization problem.

There is no feasible point for the optimization problem obviously there is no optimal point. So, if the data is not separable the optimization problem is no feasible point and hence no optimal, so no solution, right? Which means in general if the data is not linearly separable we cannot find the optimal hyperplane by this formulation, so we need to change this formulation, so that we can. Basically, why are we not able to find a solution? Because my constraints may not be satisfied, whenever constraints are not fully satisfied the standard technique in constraint optimization we used to use what are called slack variables. So, that we will keep our constraints not hard, but soft so to say and then we we we manage to get a feasible point, so that is what we will do next.


(Refer Slide Time: 21:42)

Using slack variables

- When data are not linearly separable, we can try:

$$\begin{aligned} &\text{minimize} && \frac{1}{2}W^T W + C \sum_{i=1}^n \xi_i \\ &\text{subject to} && y_i(W^T X_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ &&& \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned}$$

- Opt. variables: W, b, ξ_i .

 © NPTEL course - p.30/99

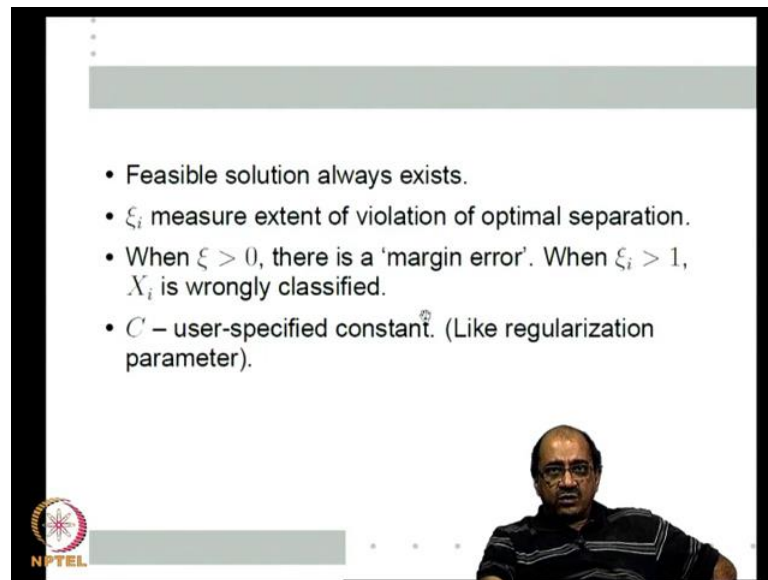
So, what we do is for general case, normally what is that we want y_i into W transpose X_i plus b greater than or equal to 1, if that is satisfied forget about the ξ_i right now, we actually wanted y_i into W transpose X_i plus b greater than or equal to 1 for separation, but suppose there is no W, b that can satisfy this. Then I say if you cannot get it greater than 1 get it at least greater than $1 - \xi_i$, right? Where ξ_i some positive number, so essentially ξ_i measures a penalty of how much you are unable to satisfy your required separability constraint.

If I did not have this $1 - \xi_i$ y_i into W transpose X_i plus b greater than 1 is my separability constraint, I am unable to get W, b so that for each i it is definitely greater than or equal to 1, I said at least get it greater than or equal to $1 - \xi_i$ for some slack ξ_i . Now, I do not want too much slack if I have too much slack obviously you know you will take W to be 0, right? Because you want to minimize this you take W to be 0 and put everything into this slack, so I do not want too much slack, so the ξ_i are the slack variables I will penalise ξ_i . So, if ξ_i are large this term will be large, I am minimizing not only half W transpose w , but also summation i is equal to 1 to n ξ_i , right?

So, my optimization variables now are W, b and ξ_i , find W, b and ξ_i , ξ_i is $\xi_1, \xi_2, \dots, \xi_n$. So, that I added the ξ_i here, but I do not want too much slack, so I added a penalty here. So, the first thing about this is that a feasible solution always is this you

give me any $W b$, I put a corresponding y_i and X_i here for each of the i . If it is greater than 1 I am fine I will take ψ_i to be 0, if it is not greater than 1 whatever is this slack I need, I take that into ψ_i , right? So, any $W b$ is feasible now, because I can always find ψ_i to satisfy this, right?

(Refer Slide Time: 23:28)



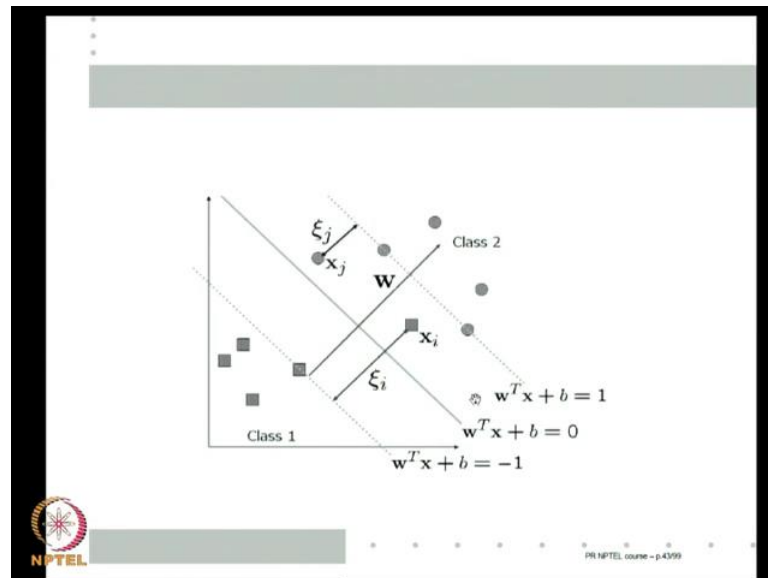
So, there is no problem about finding the feasible point, feasible solution always exist any $W b$ is feasible. ψ_i measures the extent of violation of the optimal separation that also we seen. So, basically if ψ_i is greater than 0 then there is a margin error, if ψ_i is greater than 1 X_i is wrongly classified you see. If ψ_i is greater than 0 what it means, $W y_i$ into $W \text{ transpose } X_i$ plus b is is not greater than or equal to $()$ it might be only greater than or equal to 0.5, even if it is greater than or equal to 0.5 the sign of $W \text{ transpose } X_i$ plus b is same as sign of y , right?

So, I am not getting the optimal separation because we wanted so much separation, we do not, did not want any training pattern between $W \text{ transpose } X$ plus b is equal to plus 1 and minus 1. That much separation I am not getting if ψ_i is greater than 0, but if ψ_i is greater than 1, that means this sign is changing which means I have not even classifying it correctly, right? That is that is how ψ_i measure the X under violation.

Now, finally, what we are minimizing is my margin plus this some of these slack variables they tells you how much violation of separation I am doing and C is a kind of a exchange rate constant. How much of increased margin of separation am I willing to take

at the expense of you know making so many errors in classification to say. So, this tells me some measure of the error in optimal separation, so this C tells you how much of error I am willing to trade with how much of increase in margin, C is used as a different parameter much like regularisation constant later on we will see that it is actually regularisation constant, C is user specified.

(Refer Slide Time: 25:34)



So, we can see this formulation what we are saying is this is W transpose X plus b is equal to 0, this is W transpose X plus b is equal to 1 W transpose X plus b is equal to minus 1. So, I take this to be my w , W on b then if a class two pattern came here, so for all the patterns which are beyond this anyway my ψ will be 0. If a class two pattern came here between these two this ψ_j would be greater than 0, but less than 1 still on the right side. But I am actually counting this as my margin right because I have taken this W on b , I am counting this as my margin with respect to this much margin, this is one pattern that is not giving me that much margin, right?

On this side, on the other hand this pattern is actually on the wrong side of the hyperplane, right? If I want to take this as the separating, so here this ψ I would be greater than 1 right, but essentially what it means is given this given this pattern is here and this pattern here, normally I would not have been able to find any separating hyperplane. What I decided to do is call these two patterns errors then I get a separating hyperplane with such a nice merge, that is what this formulation is. Separate bit a penalty

by misclassifying a few patterns and then separate the rest of it; that is what I am trying to do.

(Refer Slide Time: 26:56)

• The optimization problem now is

$$\min_{W, b, \xi} \quad \frac{1}{2} W^T W + C \sum_{i=1}^n \xi_i$$

subject to

$$1 - \xi_i - y_i(W^T X_i + b) \leq 0, \quad i = 1, \dots, n$$

$$-\xi_i \leq 0, \quad i = 1, \dots, n$$

So, that is my optimization problem, right? Minimize over W , b and ψ_i , so I put that as the all the ψ_i as the bold ψ_i , half W transpose W plus C times is sum of slack for I have written this in now my standard optimization thing in optimization all my constraints have to be written as something less than or equal to 0. So, I wrote it as instead of writing y_i into W transpose X_i plus b greater than or equal to $1 - \psi_i$, I wrote as $1 - \psi_i - y_i$ into W transpose X_i plus b less than or equal to 0.

Similarly, instead of writing ψ_i greater than or equal to 0 I wrote minus ψ_i less than or equal to 0, this is my optimization problem now. Now, what is my Lagrangian? Once again my Lagrangian will be a function of the optimization variable W , b , ψ_i and there will be 1 Lagrange multiplier per constraint so there will be n Lagrange multipliers for these n constraints, call them μ_1 to μ_n . There will be another n Lagrange multiplier for these n constraints, call them λ_1 to a λ_n , then my Lagrangian will be a function of W , b , all the ψ_i and then μ_i and λ_i , right?

Lagrange will be my cost function plus sum over μ_i into this constraint plus sum over λ_i into this constraint, right? So, my Lagrange will be W , b , ψ_i , μ , λ this is my cost function plus μ_i times my $1 - \psi_i - y_i$ into W transpose X_i plus b plus λ_i times minus ψ_i , so I wrote it as minus $\lambda_i \psi_i$.

(Refer Slide Time: 28:37)

• The Lagrangian now is

$$L(W, b, \xi, \mu, \lambda) = \frac{1}{2} W^T W + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (W^T X_i + b)) - \sum_{i=1}^n \lambda_i \xi_i$$

• μ_i are the lagrange multipliers for the separability constraints as earlier.

• λ_i are the lagrange multipliers for the constraints $-\xi_i \leq 0$.

NPTEL © NPTEL course - p4509

So, μ_i are the Lagrange multipliers for this separability constraints, as earlier for each examples. λ_i is the Lagrange multipliers for the constraints $\xi_i \leq 0$. So, this is my Lagrangian, so what will be my Kuhn-Tucker conditions? I have to get derivative with respect to W is equal to 0, derivative with respect to b is equal to 0, gain with respect to ξ_i that means partial derivative with respect to each of ξ_i is equal to 0, that is the derivative constraints.

Then I need feasibility, all constraints should be satisfied and μ_i should be greater than or equal to 0, λ_i should be greater than or equal to 0 and there are complimentary slackness. So, let us write all our Kuhn Tucker conditions, first is gradient of L with respect to W is equal to 0, so as W the function still the same it has not changed. So, I get one W from here and I get $\mu_i y_i X_i$ here right? So, gradient of L is equal to 0 still gives me the same old thing W^* is summation over i $\mu_i y_i X_i$.

(Refer Slide Time: 29:38)

The Kuhn-Tucker conditions give us

- $\nabla_W L = 0 \Rightarrow W^* = \sum_{i=1}^{\ell} \mu_i^* y_i X_i$
- $\frac{\partial L}{\partial b} = 0 \Rightarrow \sum \mu_i^* y_i = 0$
- $\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \mu_i^* + \lambda_i^* = C, \forall i$
- $1 - \xi_i - y_i(W^T X_i + b) \leq 0; \xi_i \geq 0; \forall i$
- $\mu_i \geq 0; \lambda_i \geq 0, \forall i$
- $\mu_i(1 - \xi_i - y_i(W^T X_i + b)) = 0; \lambda_i \xi_i = 0, \forall i$

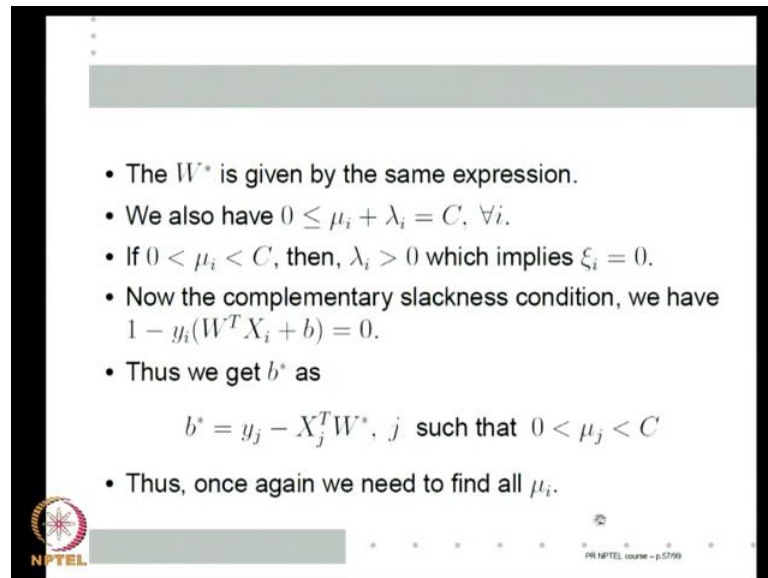
NPTEL

© NPTEL course - p.5/59

Similarly, with respect to b I get the same old thing, what will it give me if I differentiate with respect to each of these ψ_i ? There is one ψ_i here, so we get a C term plus C here I get a minus μ_i from here, minus λ_i from here. So, c minus μ_i minus λ_i is equal to 0 for each i , which is same as μ_i plus λ_i is equal to C for each i . So, this this the derivative with respect to remaining optimization variables will tell me μ_i^* plus λ_i^* is equal to C for every i , then I have the feasibility constraints all my ψ_i have to be positive.

These are my other constraint $1 - \psi_i - y_i W^T X_i + b \leq 0$ the separability constraints. All my Lagrange multipliers have to be non zero, non negative and my complement is slackness, μ_i into this constraint should be equal to 0. Similarly, $\lambda_i \psi_i$ should be equal to 0, these are my full complement of Kuhn-Tucker conditions for this optimisation problem.

(Refer Slide Time: 30:52)



• The W^* is given by the same expression.

• We also have $0 \leq \mu_i + \lambda_i = C, \forall i$.

• If $0 < \mu_i < C$, then, $\lambda_i > 0$ which implies $\xi_i = 0$.

• Now the complementary slackness condition, we have $1 - y_i(W^T X_i + b) = 0$.

• Thus we get b^* as

$$b^* = y_j - X_j^T W^*, j \text{ such that } 0 < \mu_j < C$$

• Thus, once again we need to find all μ_i .

NPTEL

PH NPTEL course - p.5709

Let us see what this means, W^* is still given by the same expression as you have already seen, we also have the extra condition $\mu_i + \lambda_i = C$ for every i and because both μ_i and λ_i are non negative $\mu_i + \lambda_i$ will also be non negative, so we have this for every i . What does this mean? If μ_i is strictly between 0 and C , right? My complementary slackness, remember my complementary slackness implies if μ_i is greater than 0, then this term is 0, this constraint becomes equal to 0. Similarly, if λ_i is greater than 0 then ψ_i becomes equal to 0, right?

Now, because I have this, if I consider all μ_i which are not only strictly greater than 0, but are strictly less than C , then $\mu_i > 0$, the first complementary slackness tells me that the constraint is 0, because λ_i is greater than 0. Why is λ_i a greater than 0? Because μ_i is strictly less than C and I need $\mu_i + \lambda_i = C$, λ_i has to be strictly greater than 0, which implies $\psi_i = 0$, right? So, if I consider any μ_i such that μ_i is strictly between 0 and C , then λ_i is greater than 0, so $\psi_i = 0$ and $\mu_i > 0$, so this term is equal to 0.

In this term I know $\psi_i = 0$, I get $1 - y_i(W^T X_i + b) = 0$. Now, the complementary slackness conditions give you $1 - y_i(W^T X_i + b) = 0$. So, earlier for any $\mu_i > 0$ this is true, now for all μ_i which are strictly between 0 and C this is true. Now, this obviously immediately gives me like earlier my b^* , right? Only thing is instead of this formula

being true for any μ that is strictly greater than 0, it is only for those μ which are strictly between 0 and C.

Now, I can think of my vector that determine b to be one that are not stuck at 0 or C. See now that, see earlier I had only μ_i greater than 0, now I have μ_i plus λ_i is equal to C and λ_i is also greater than or equal to 0, which means μ_i is now strictly between 0 and C anyway I need to have $0 \leq \mu_i \leq C$. Similarly, $0 \leq \lambda_i \leq C$. So, unlike earlier problems μ also has a strict upper bound, if μ_i is strictly between their bounds not equal to either of the bounds, then that μ_i , that determines these X_i which are closest to the separating boundary and all.

Earlier, the X_i that are closest to the separating hyperplane are those whose corresponding Lagrange multipliers are strictly positive. Now, with this slack variable formulation, I have to look at those X_j for the for whom the corresponding Lagrange multipliers are not stuck at either end, they are neither stuck at 0, nor stuck at C. Those μ correspond to X_j which satisfies this which means those are the X_j that are closest to the separating hyperplane now. Apart from that once again if you I can (()) a star and b^* so I am done. So, once again I can use the dual to find the μ .

(Refer Slide Time: 34:18)

• We can derive the dual as before. The dual function is

$$q(\mu, \lambda) = \inf_{W, b, \xi} L(W, b, \xi, \mu, \lambda)$$

where the lagrangian is given by

$$L(W, b, \xi, \mu, \lambda) = \frac{1}{2} W^T W + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \mu_i (1 - \xi_i - y_i (W^T X_i + b)) - \sum_{i=1}^n \lambda_i \xi_i$$

NPTEL
© NPTEL course - p.59/99

So, what will be the multipliers, now there are the μ and the λ s, so q is a function of μ 1 λ , it is obtained as infimum over the all the optimization variables of the

primal problem namely $W b \psi$ of the Lagrangian, Lagrangian function $W b \psi$ and the two sets of Lagrange multipliers $\mu_1 \lambda$, so where the Lagrangian as we have already seen is given by this. So, basically for a given a given a particular μ_1 to μ_1 and λ_1 to λ_n , if I want the value of q at that $\mu_1 \lambda$, I need to find infimum of this expression at that $\mu_1 \lambda$ with respect to $W b$ and all the ψ_i , right?

Let us first consider finding infimum with respect to ψ_i , what is, what is the ψ_i dependence of this term, for each i I have a ψ_i into $C - \mu_i - \lambda_i$ term. So, for a given λ and μ if $C - \lambda_i - \mu_i$ is not 0, then I can take that ψ_i to be as negative as I want, when I am taking the infimum or as positive as I want and hence infimum will be minus infinity.

(Refer Slide Time: 35:41)

- In the lagrangian we have the term $\sum_i (C - \mu_i - \lambda_i)\xi_i$.
- Since we take infimum w.r.t. ξ_i , we need to impose $C = \mu_i + \lambda_i, \forall i$.
- When we impose this, all the terms containing λ_i or ξ_i drop out and hence now the q function would be same as earlier.
- We only need to ensure (in the dual) that $\lambda_i \geq 0$ and $C = \mu_i + \lambda_i, \forall i$.
- This is easily done by ensuring $0 \leq \mu_i \leq C$.

Because we have this term, because you are taking infimum with respect to ψ_i , unless we impose C is equal to μ_i plus λ_i infimum will be minus infinity and that we are not interested. So, once again we only work, we only need to look at $\mu_1 \lambda$, so that C is equal to μ plus λ . When we, when we consider μ and λ so that for each i μ_i plus λ_i is equal to C , then all the ψ_i terms will drop out, for each ψ_i i have got a $c - \mu_i - \lambda_i$ and similarly, all the λ_i terms drop out, right?

The moment I impose the c is equal to μ_i plus λ_i , then all terms containing λ_i or ψ_i or dropout of the expression. Once the drop out of the expression what is left is the old expression and I need to find infimum with respect to μ and we already that is that old q function we know, right? So, if you want to find the new q function, all you have to do is that we get the old q function only anyway, because we have to anyway impose this. But in addition to the old constraints I have to also impose the constraints that λ_i is greater than equal to 0 and C is equal to μ_i plus λ_i . Because I have to impose C is equal to μ_i plus λ_i , my q function will not contain any λ_i terms.

So, q is still only function of μ , so I do not have to actually maximize to our λ_i , as long as I can somehow ensure that these two conditions are satisfied. These two are easily satisfied by imposing this condition, and once I impose this condition on μ , I can take $c - \mu_i$ to be λ_i , because λ_i need not have to does not acquire the cost function, λ_i does not affect the cost function as long as I satisfy the constraints on λ_i I am done. So, let us say I solve it only for μ with μ as this constraint, then $C - \mu_i$ I take to be λ_i , then I have to satisfy λ_i greater than or equal to 0, as well as μ_i plus λ_i is equal to C .

(Refer Slide Time: 37:42)

The dual

- The dual problem now is:

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

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

- The only difference - upper bound also on μ_i .

NPTEL

© NPTEL course - p 0599

So, which means my dual now will be exactly be same as the old dual, the cost function. All the old constraint will be there, earlier I had only μ_i greater than or equal to 0, now

I have μ_i less than equal to C also. So, the only difference is in upper bound, right? This very nice see in the primary problem, the problem is totally different I get a plus C into summation ψ_i in the cost function, each of my constraints is changed I have got n extra constraints.

So, if you wanted to solve the primal lots of extra work needs to be done, but if you want to solve the dual whether I am thinking that the patterns or linear separable all I am willing to use slack variables, so that I have always have a feasible solution and always get a solution. The only difference is whether or not I put an extra upper bound which is anyway used as different constant. So, essentially solved this optimisation problem the dual by by just choosing some convenient upper bound and you know that is all right. The only difference is in upper bound I am doing, so that is for another reason for preferring to solve the dual, right?

(Refer Slide Time: 38:45)

• The primal problem is

$$\text{minimize } \frac{1}{2}W^T W + C \sum_{i=1}^n \xi_i$$

subject to $y_i(W^T X_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n$
 $\xi_i \geq 0, \quad i = 1, \dots, n$

• As $C \rightarrow \infty$, we get back the old problem.

NPTEL

PR NPTEL course - p.0509

Of course it is very easy to see, this is my primal problem. So, if C tends to infinity, if C very a large that even a little bit of positive ψ_i , this term becomes so large that that will not be the where the minimum is attained. So, ultimately my minimum will be attained only for those ψ_i for that place where ψ_i are all 0, if ψ_i are all 0 it is like the old problem. So, if C tends to infinity we get the old problem, if C tends to infinity I get the old dual obviously, right?

(Refer Slide Time: 39:19)


• The dual problem is:

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

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

• We solve dual and the final optimal hyperplane is

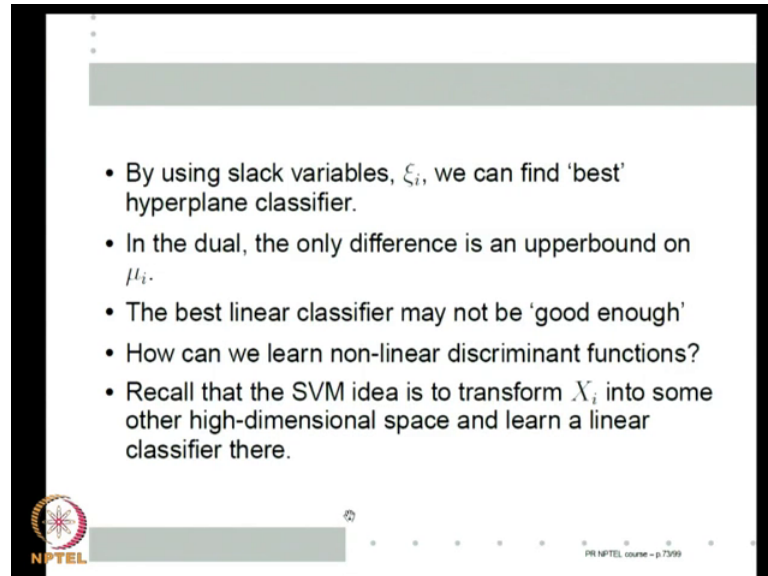
$$W^* = \sum \mu_i^* y_i X_i,$$
$$b^* = y_j - X_j^T W^*, \quad j \text{ such that } 0 < \mu_j < C.$$

 © IN NPTEL course - p 09/09

So, this is my dual, so we solve the dual and once again this is given by my, this is my optimal hyperplane. The only difference is earlier the b^* is for any j of that μ_j greater than 0 and now it is μ_j strictly between 0 and C . So, this extra one upper bound on C is all the difference in the optimization problem between solving the problem assuming training data is linearly separable or you know taking general data and finding the best possible linear separation. Is best possible in the sense I have just decided to make a few errors on some patterns and add the others in one specific way, by adding these slack variables into the optimization problem.

So, if this is a formulation with somehow I am minimising some measure of the error in classification while maintaining a maximum margin, right? When in the linearly separable case the optimization problem is very nice to see, because it maximizes margin and the separability constraint. In this case because of this C , I am essentially adding apples and oranges I am taking some part of the margin, I am some margin plus some constant into the others are making I am minimising that. But anyway this is what one does has to do in general, because in general we have no way of knowing whether the training data is linearly separable or not.

(Refer Slide Time: 40:42)



The slide contains the following text:

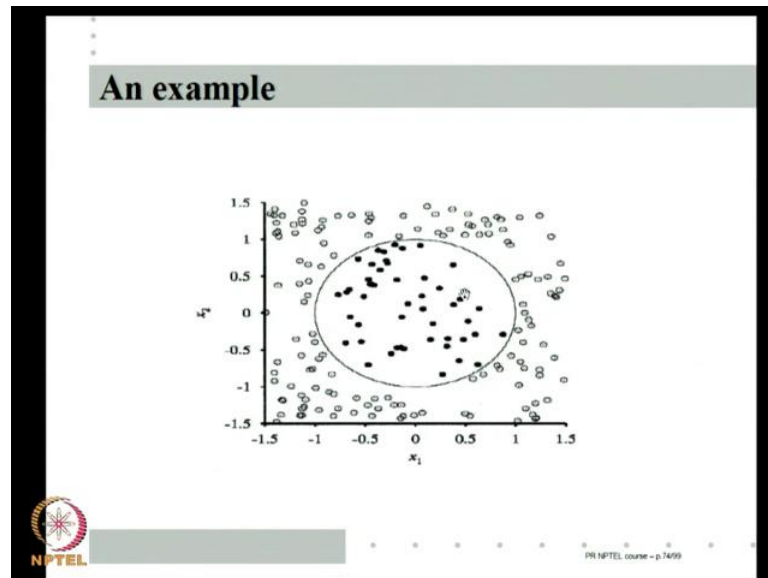
- By using slack variables, ξ_i , we can find 'best' hyperplane classifier.
- In the dual, the only difference is an upperbound on μ_i .
- The best linear classifier may not be 'good enough'
- How can we learn non-linear discriminant functions?
- Recall that the SVM idea is to transform X_i into some other high-dimensional space and learn a linear classifier there.

At the bottom left of the slide is the NPTEL logo. At the bottom right, it says "© NPTEL course - p 73/99".

So, by using the slack variables ψ_i , we can find the best hyperplane classifier. So, this is the optimal separating hyperplane very good, so we can find, instead of finding any linear classifier now, we have a concept of some optimal linear classifier and we can find it by solving a nice quadratic optimization problem as a quadratic cost function and linear constraints. The only difference is when we are solving the dual as we have seen the only difference is an upper bound and μ , so I just throw in a constant c as an upper bound on μ and I always solve only the dual, so it is very nice.

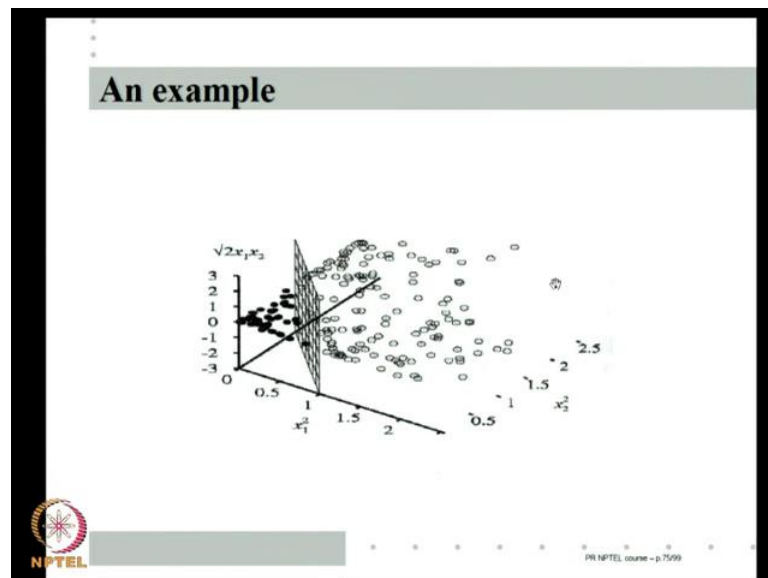
Of course, as we have already seen the best linear classifier may not be good enough, the idea is we want to learn nonlinear discriminant functions and recall that the whole idea of SVM is to transform ψ into some other high dimensional space and then learn a linear classifier. The idea is ψ will be transformed to z and then in that z space you learn a linear classifier, right?

(Refer Slide Time: 41:38)



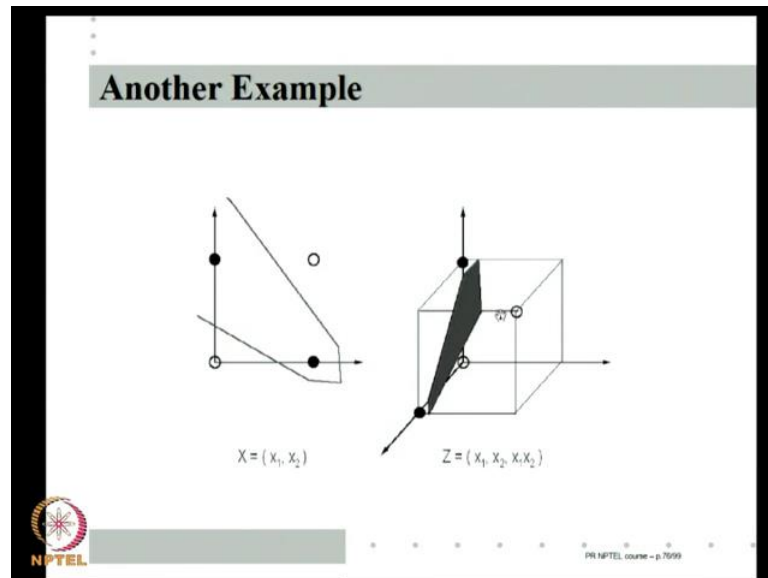
So, as you just to recall this is the idea, so if this is the original feature space and i needed a quadratic discriminant function to separate the two classes.

(Refer Slide Time: 41:48)



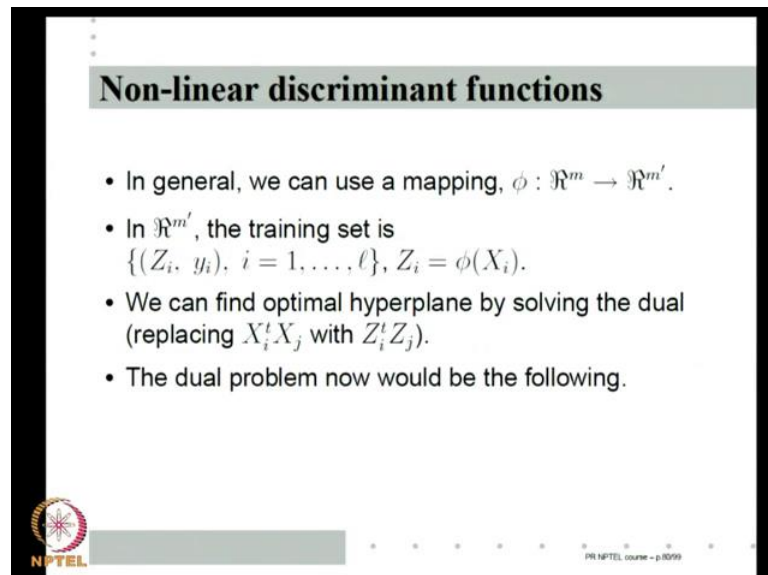
If I map the patterns into high dimensional space, I would be able to separate with the hyperplane. Essentially this is in the $X_1 X_2$ space, whereas this is in the $X_1^2 X_2^2$ space. So, a hyperplane can separate the classes now, whereas I needed a second ordered discriminant function, quadratic discriminant function in the original space.

(Refer Slide Time: 42:14)



It is another example, if I have an X r problem, right? To separate the two classes I need a nonlinear classifier a single layer line cannot separate them. But I can always map them to a higher dimensional space, which is $X_1 X_2$ and $X_1 X_2$ and in that space a hyperplane can separate them, so that is the basic idea.

(Refer Slide Time: 42:33)



So, in general we used some mapping ϕ from m dimensional space to m' dimensional space. In the new space my training set now will be $z_i y_i$ where z_i is $\phi(X_i)$. So, we will be in the next, in the rest of this class and also most of our discussions of

same we will use z and $\phi(X)$ interchangeably, $\phi(X_i)$ will be z_i that is what the new transform, this is the vector corresponding to the training vector X_i in the transform space. So, once we use this ϕ to transform all the training patterns from m dimensional space to m' dimensional space, my new training set has become z_i, y_i , so z_i is equal to $\phi(X_i)$.

So, we can find the optimal hyperplane by solving the dual, what was our dual? If I wanted to find the best linear classifier, this is what I have to solve where X_i, X_j is the training patterns, so nowhere else the training patterns have appeared. So, if my new training my (z_i, y_i) transformed them to z_i I just replace X_i, X_j by z_i, z_j that is all or $\phi(X_i), \phi(X_j)$. So, we can find the optimal hyperplane by solving the dual just simply replacing X_i, X_j with z_i, z_j .

(Refer Slide Time: 43:51)

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

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

- This is an optimization problem over \mathbb{R}^n (with quadratic cost function & linear constraints) irrespective of ϕ and m' .
- But computationally expensive?

© NPTEL course - p.0309

The dual problem will be what now? I just replace X_i, X_j with z_i, z_j z_i, z_j is same as $\phi(X_i), \phi(X_j)$ this. Am I done? The good thing about it is there is an optimization problem over n the number of examples, it has nothing to do with the function ϕ or m' . Nothing to do meaning the the dimensionality of this optimization problem does not depend on the function ϕ or the m' space. So, originally I may have hundred dimensional space, now $\phi(X)$ may be ten thousand dimensional space, but in both spaces I will be solving may be a thousand dimensional optimization problem because I have thousand examples.

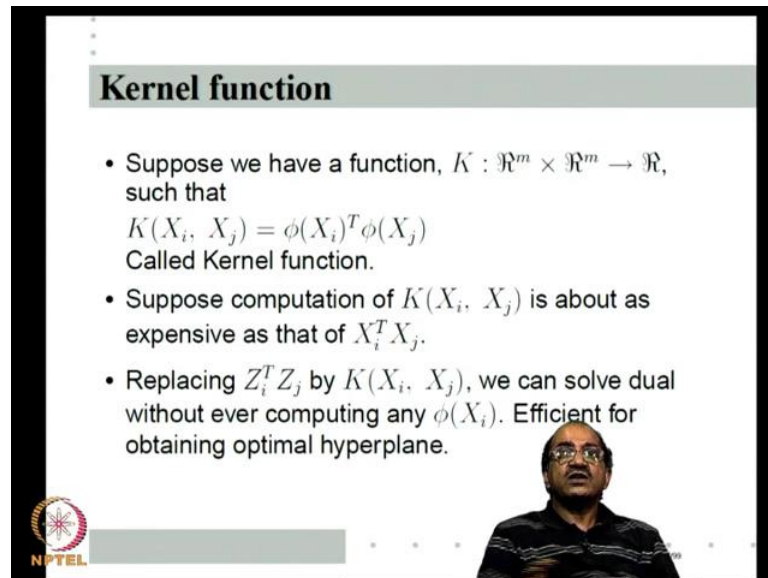
So, the optimization problem remains only of $r(\cdot)$ at least earlier we said namely solved the problem in the z space, I have to solve a problem of the dimension of the z space that is m prime, but now I am not solving that, because I am solving the dual. The dimensionality of my optimization problem remains n the number of examples, it does not become m prime. So, even if m prime is ten thousand or one million my dimensional optimization problem does not, does not increase it will stay at n . This is the good thing, but does it mean that I have really got rid of computational cost; of course, during learning I have to solve this problem to find μ_i given all my training examples, so I have to calculate $\phi^T X_i X_j \phi$.

Of course, I can pre compute it because these do not depend on the optimization variables here, so I can find all the n squares or the n square number $\phi^T X_i X_j \phi$ one time and then solve use those constants here and solve this optimization problem. Of course, still it can be expensive, because I have to first transform each of the X through ϕ , ϕ might be a difficult transformation and in that high dimensional space may be ten thousand dimensional space or one million dimensional space I have to calculate the inner products.

But I did it do it only once and then solve this μ obtained, is that all? No, what do I do with the μ I have to use my W stars, so the W star will live in the ϕ space so W star is the same dimension as ϕ . So, how to store W star? The one million dimensional W star and every time now you give me a new X I still have to do $W^T X$, right? I have not, so every time I use the classifier I suffered the computational cost of going into high dimensional space and finding inner products there.

Of course, in addition even though the optimization problem remains after dimension n , I do not know because effectively I am learning a hyperplane a very high dimensional space, are thousand examples enough. The second problem we will solve much later, but this SVM method as efficiently tackled that particular let us see the reasons for it later. But right now let us concentrate on the computational problem, how do I get rid of having to do this $\phi^T X_i X_j \phi$?

(Refer Slide Time: 46:52)



Kernel function

- Suppose we have a function, $K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$, such that
$$K(X_i, X_j) = \phi(X_i)^T \phi(X_j)$$

Called Kernel function.
- Suppose computation of $K(X_i, X_j)$ is about as expensive as that of $X_i^T X_j$.
- Replacing $Z_i^T Z_j$ by $K(X_i, X_j)$, we can solve dual without ever computing any $\phi(X_i)$. Efficient for obtaining optimal hyperplane.

NPTEL

The basic idea is the following, suppose we have a function which takes (()) of m dimensional vectors and gives me a real number, such that k of X_i, X_j is $\phi(X_i)^T \phi(X_j)$, we call such a k as Kernel function. What do you mean by I have such a function? Of course, I have such a function once I know ϕ I can define $k(X_i, X_j)$ to be $\phi(X_i)^T \phi(X_j)$ or basic idea is that such a function is less expensive than calculating $\phi(X_i)^T \phi(X_j)$, right?

We may not have time to look at Kernel function in this class, but not to make a big secrecy of its. For example, if I take $k(X_i, X_j)$ to be $1 + X_i^T X_j$ whole square $(1 + X_i^T X_j)^2$ needs one extra multiplication compared to $X_i^T X_j$. Suppose, $X_i^T X_j$ is 100, X_i is 100 dimensional, $X_i^T X_j$ needs hundred multiplications, one addition we do not care $1 + X_i^T X_j$ whole square we will need one more multiplication, right?.

So, we are looking at Kernel function like that we will (()) square as a Kernel function would be $\phi(X_i)^T \phi(X_j)$ for a ϕ which effectively learns the quadratic discriminant function in the original space, right? So what we considered right in the beginning two class ago about mapping X_1, X_2, \dots, X_n square $X_1^2, X_2^2, \dots, X_n^2$ and then finding a linear classifier there, where we said we need order n^2 dimensional z . Finding an inner product in that space would be like $1 + X_i^T X_j$ whole square, we will see that later.

So, in that sense instead of doing ten thousand multiplications to find $\phi^T X_i$ and $\phi^T X_j$ and also many multiplication to actually find $\phi^T X_i$ and $\phi^T X_j$, we will just do one extra multiplication and $X_i^T X_j$ and all. So, we are looking for functions which are not particularly more expensive than doing $X_i^T X_j$. Let us say such a function exists, that let us just assume for now such a function exists, then what can that function give us? Then replacing $\phi^T X_i$ and $\phi^T X_j$ by $k X_i^T X_j$ we can solve the dual, this is my dual $\phi^T X_i$ and $\phi^T X_j$ is $k X_i^T X_j$, so I will just list a $\phi^T X_i$ terms for $\phi^T X_j$ I have put $k X_i^T X_j$ here.

So, $k X_i^T X_j$ is much less computationally expensive than $\phi^T X_i$ and $\phi^T X_j$ which means I have solved at least the computational problem of calculating $\phi^T X_i$ and $\phi^T X_j$ while learning the μ . So, replacing $\phi^T X_i$ and $\phi^T X_j$ by $k X_i^T X_j$ we can solve the dual without ever computing any $\phi^T X_i$, so for obtaining the optimal hyperplane, right? So, one problem is solved now, so I can efficiently solve the dual see that is the whole idea why we went to the dual.

(Refer Slide Time: 50:33)

• The primal problem is

$$\text{minimize } \frac{1}{2} W^T W + C \sum_{i=1}^n \xi_i$$

subject to $y_i(W^T X_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n$
 $\xi_i \geq 0, \quad i = 1, \dots, n$

• As $C \rightarrow \infty$, we get back the old problem.

NPTEL

PR NPTEL course - p0509

If I want to do the same thing in the primal, this X_i will become $\phi^T X_i$, right? So, I cannot use the Kernel function here if and this is not even linear constraints anymore because I have a $\phi^T X_i$ there, right? I mean occurs will be linear constraint in the sense $X_i^T \phi$ are constants linear in W and b even though it is a linear constraint, this $\phi^T X_i$ because it comes into the constraints, right? Every time I do any, I will use any

numerical algorithm in every iteration I have to calculate the corresponding $\phi(X_i)$ in the primal.

(Refer Slide Time: 51:07)

• The dual problem is:

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

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

• We solve dual and the final optimal hyperplane is

$$W^* = \sum \mu_i^* y_i X_i,$$

$$b^* = y_j - X_j^T W^*, \quad j \text{ such that } 0 < \mu_j < C.$$

NPTEL PR NPTEL course - p.0909

On the other hand in the dual, because it comes only with $X_i^T X_j$ and nowhere else I am, I am able to substitute this by the Kernel right.

(Refer Slide Time: 51:30)

Kernel function

- Suppose we have a function, $K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$, such that

$$K(X_i, X_j) = \phi(X_i)^T \phi(X_j)$$
 Called Kernel function.
- Suppose computation of $K(X_i, X_j)$ is about as expensive as that of $X_i^T X_j$.
- Replacing $X_i^T X_j$ by $K(X_i, X_j)$, we can solve dual without ever computing any $\phi(X_i)$. Efficient for obtaining optimal hyperplane.
- What about storing W^* ? Computing $\phi(X)^T W^*$ for new patterns?

NPTEL PR NPTEL course - p.0909

So, what the kernel has given me is that by replacing this $\phi(X_i)^T \phi(X_j)$ by $K(X_i, X_j)$ I have been able to efficiently solve the dual problem without actually calculating $\phi(X_i)$ and obtain all my hyperplanes, all my Lagrange multipliers. So, far so good, but

does it mean that I have solved the problem of inefficiency in inherent with this idea, one of the inefficiencies inherent with the idea is every time we may use new pattern to classify. I have to transform it in this high dimensional space and then classify it with the corresponding name.

So, even though for obtaining my mu I do not have to calculate phi X i what do I do with my mu? Ultimately, obtain my mu so that I can calculate W star and I calculate W star, because next time you give me an X I have to calculate phi X transpose W star to classify it. So, W star stays in my million dimensional space so I have to store W star and I have to calculate phi X every time you give me a mu X and do this, what do I do about this?

(Refer Slide Time: 52:34)

Kernel function based classifier

- Let μ_i^* be soln of Dual. Then $W^* = \sum \mu_i^* y_i \phi(X_i)$.
- Then we have

$$b^* = y_j - \phi(X_j)^T W^* = (y_j - \sum_i \mu_i^* y_i \phi(X_i)^T \phi(X_j))$$
- Given a new pattern X we only need to compute

$$f(X) = Z^T W^* + b^*$$

$$= \sum_i \mu_i^* y_i \phi(X_i)^T \phi(X) + b^*$$

$$= \sum_i \mu_i^* y_i K(X_i, X) + (y_j - \sum_i \mu_i^* y_i K(X_i, X_j))$$

Well as it turns out I do not have to do that also, it is very very interesting thing. Let us say mu i star are the solutions of the dual, then we know that W star is summation mu i star y i phi i X i. So, for example what will be b star we seen b star, will be y j minus phi X i transpose W star, now W star is given by this, so it will be y j minus summation over i mu i star y i phi X i transpose phi X j. Wherever phi X i transpose phi X j comes I know I do not have to do any phi X phi j computation, I can convert it into Kernel function and k psi X j is much less expensive than phi X i transpose phi X j.

So, for example to compute b star I do not need to calculate phi X at all, if I know all the mu stars, then I can just put k X i X j here and calculate b star of course b star is not the

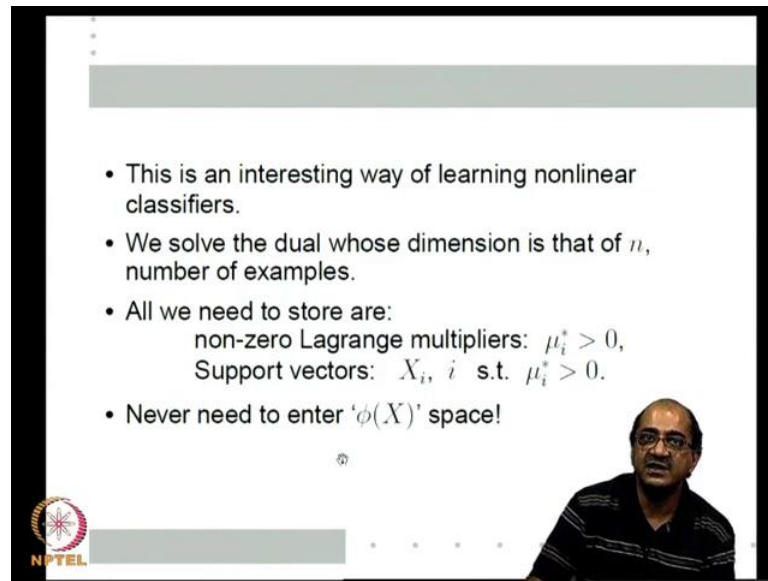
only thing I need, I need much more than b^* . What you have to do give a new pattern, what you have to compute? I have to compute $f^T X$ right ultimately my classification is depend on the sign of this $f^T x$, what is this $f^T x$? $z^T W^* + b^*$, what is $z^T \phi(x)$? So, $\phi^T X^T W^* + b^*$, that is what I have to calculate and find its sign.

What will be that be, so W^* I write it as $W^T \phi(x)$, W^* is $\mu_i y_i \phi(X_i)$, so this is summation over i $\mu_i y_i \phi(X_i)^T z$ is $\phi^T X + b^*$, now I know how to write b^* . So, this is summation over i $\mu_i y_i \phi(X_i)^T \phi(X_j)$ can be written as $k(X_i, X_j) + y_j \mu_i y_i k(X_i, X_j)$. This is for a particular j which have for which the corresponding μ_j^* is between 0 and c . As you can see while finding the μ I do not have to calculate any $\phi(x)$, after obtaining the μ you give me any μ^* and I want to calculate its actual $\phi^T X^T W^* + b^*$.

Once again I do not have to calculate $\phi(X)$ at all I can do with the Kernel functions I just have to calculate this, right? This is a really interesting way of learning non linear classifiers. Now, no matter what $\phi(X)$ is, no matter what how high dimensional m' is I am solving n , a constraint optimization problem and that to a quadratic cost function linear constraints constrained optimization problem. Whose dimension is simply a number of examples and to solve this I do not have to actually calculate my $\phi(x)$, ϕ is only in the background is only effectively, because I Kernelise my inner product.

So, in by while solving the dual, where $\phi(X_i)^T \phi(X_j)$ would have come I would have put $k(X_i, X_j)$, I solve the dual get all my μ , once I get my μ I have to effectively calculate this, that I can calculate like this using my μ and the Kernel function. So, I never actually calculate W^* , never actually calculate $\phi(X)$ for any x , right?

(Refer Slide Time: 55:54)

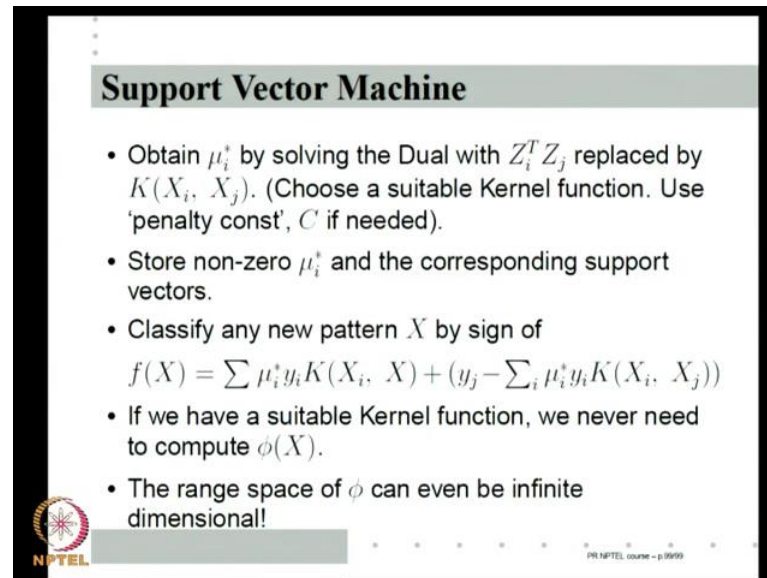


- This is an interesting way of learning nonlinear classifiers.
- We solve the dual whose dimension is that of n , number of examples.
- All we need to store are:
non-zero Lagrange multipliers: $\mu_i^* > 0$,
Support vectors: X_i, i s.t. $\mu_i^* > 0$.
- Never need to enter ' $\phi(X)$ ' space!

It is an interesting way of learning non-linear classifiers, we solved the dual whose dimension is that of n , the number of examples and all we need to store or all the non-zero Lagrange multiplier is μ_i^* and the corresponding support vectors as X_i , right? All X_i corresponding to $\mu_i^* > 0$ I have to solve. If I solve them, if I store them you give me any μ, X I calculate this and then I am done. I can classify here X by doing this, I never need to enter the $\phi(X)$ space, even though I am actually using some transformation inside ϕ , which maps my m dimensional space to some other high dimensional m' dimensional space, I am finding a linear classifier there, I never calculate $\phi(X)$ for any x , I never ever enter that m' dimensional space, right?

The only connection between that space and me is the Kernel function, somehow to find the right Kernel function. We will see next class how we can find Kernel functions, but provided that I can find the kernel function, $k(X_i, X_j)$ is equal to $\phi(X_i)^T \phi(X_j)$ such a way that $k(X_i, X_j)$ is much simpler to calculate than $\phi(X_i)^T \phi(X_j)$ then I am done. So, no matter what are the transformations ϕ I am interested in, if I can properly find a Kernel function to represent it then I am done, right?

(Refer Slide Time: 57:18)



Support Vector Machine

- Obtain μ_i^* by solving the Dual with $Z_i^T Z_j$ replaced by $K(X_i, X_j)$. (Choose a suitable Kernel function. Use 'penalty const', C if needed).
- Store non-zero μ_i^* and the corresponding support vectors.
- Classify any new pattern X by sign of
$$f(X) = \sum \mu_i^* y_i K(X_i, X) + (y_j - \sum_i \mu_i^* y_i K(X_i, X_j))$$
- If we have a suitable Kernel function, we never need to compute $\phi(X)$.
- The range space of ϕ can even be infinite dimensional!

NPTEL © NPTEL course - p19999

What is the final support vector machine idea? Is that I obtain μ_i^* by solving duals where that $X_i^T X_j$ is replaced by $Z_i^T Z_j$, which is replaced by $K(X_i, X_j)$ that is essentially it chooses a suitable Kernel function and use a penalty constant C if needed. Using a penalty constant is simply putting an upper bound in my dual optimization problem, so I will just throw in a C and choose some suitable Kernel we will see what are suitable next class and I solve the dual.

Once I solve the dual, I store all the non zero μ_i^* and their corresponding support vectors. Now, that is my support vector machine this represents my classifier, how do I represent my classifier? Once I solve them, I am now ready to work, you give me any new pattern x , this is what I calculate and I classify it by the sign of this. So, if you have a suitable Kernel function we never need to calculate $\phi(X)$ for any x , what does this mean?

For all I care for the range space of ϕ can be infinite dimension, I never go there. So, this is a very, very interesting idea of effectively finding a linear classifier in a high dimensional space. So, we next class we will look at the Kernel functions, what kernel functions are suitable and how different kernel functions will allow you to find different non linear classifiers.

Thank you.