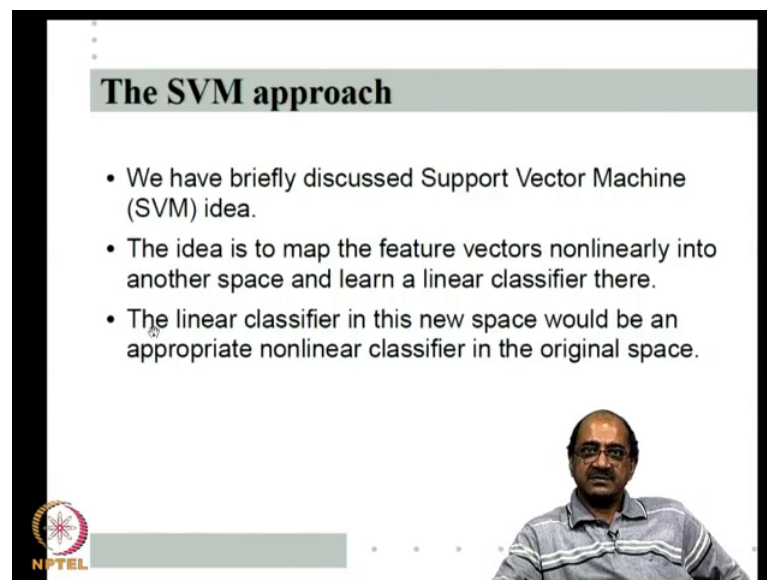


Pattern Recognition
Prof. P. S. Sastry
Department of Electronics and Communication Engineering
Indian Institute of Science, Bangalore

Lecture - 31
Support Vector Machines – Introduction,
obtaining the optimal hyperplane

Hello and welcome to this next lecture on pattern recognition course. We have been looking at essentially neural network model for the last few classes. We looked at feed forward networks with sigmoidal activation function as well as real based function networks. Learning algorithms for using them, both of them are useful models for learning non linear classifiers and regression functions, right? Then end of last class we said that we will move on to a next approach for learning non linear classifiers, where instead of wanting a general parameterized class of non linear classifiers. We are looking at learning a linear classifier in a transformed space, right? That is the so called support vector approach. So, that is what we start with today the support vector machine approach, right?

(Refer Slide Time: 01:08)



The SVM approach

- We have briefly discussed Support Vector Machine (SVM) idea.
- The idea is to map the feature vectors nonlinearly into another space and learn a linear classifier there.
- The linear classifier in this new space would be an appropriate nonlinear classifier in the original space.

NPTEL

The slide features a small video inset of Prof. P. S. Sastry in the bottom right corner. The NPTEL logo is visible in the bottom left corner of the slide.

We have briefly discussed it last class, so the basic idea is the following. You from the original feature space, you non linearly map the feature vectors into some other space. Normally a higher dimensional space, so not necessarily, so you map into a different space. In the new space learn a linear classifier. Now, this works because the linear

classifier in the new space could be an appropriate non linear classifier in the original space.

If you do the mapping properly if we choose the right mapping, then the linear classifier in the new space would be a nice non linear classifier in the original space, right? As an old Bob Dylan song refrain use to have one man sailing is another man's floor, so what is linear depends on which space you are in? So, what is linear in one space could be non-linear in another space.

(Refer Slide Time: 02:14)

- Recall the simple example we saw earlier.
- Let $X = [x_1 \ x_2]$ and let $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^5$ given by

$$Z = \phi(X) = [1 \ x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1x_2]$$
- Now,

$$g(X) = a_0 + a_1x_1 + a_2x_2 + a_3x_1^2 + a_4x_2^2 + a_5x_1x_2$$
 is a quadratic discriminant function in \mathbb{R}^2 ; but

$$g(Z) = a_0 + a_1z_1 + a_2z_2 + a_3z_3 + a_4z_4 + a_5z_5$$
 is a linear discriminant function in the ' $\phi(X)$ ' space.

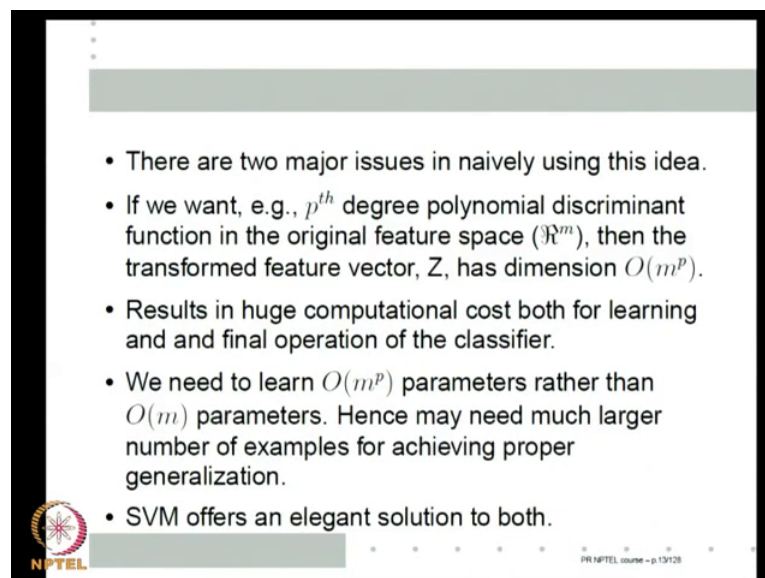
So, we shown an example of this last class. Let us recall it it is a two dimensional. Example, suppose those will feature space is two dimensional and let us say x_1, x_2 are the two feature components. Now, let us define a function a non linear function that maps \mathbb{R}^2 to \mathbb{R}^5 ; that means you are transforming the two dimensional feature space into a five dimensional feature space. So, it is a function phi that maps \mathbb{R}^2 to \mathbb{R}^5 given by given any x , the image of x phi x which we called by z is given by $1 \ x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1x_2$. $1 \ x_2$, it is a five component vector, right? Because it is an \mathbb{R}^5 , I am sorry its actually six component vector, it should have been \mathbb{R}^6 . I am sorry but, it is not \mathbb{R}^5 , but it is \mathbb{R}^6 . Now, the main idea is that if I look at a function of x , the two dimensional vector x which is a $0, a_1 \ x_1, a_2 \ x_2, a_3 \ x_1^2, a_4 \ x_2^2, a_5 \ x_1x_2$.

Then this is a quadratic discriminant function in the original two dimensional space, because it has quadratic terms in the feature components. On the other hand, if I take a

function of the vector z , let us call the components of z as you know this is $z_0, z_1, z_2, z_3, z_4, z_5$, right? So, essentially because of this one it is like an augmented feature vector in \mathbb{R}^5 , the rest of the feature vectors is \mathbb{R}^5 , but I am doing an augmentation, so we will call this as z_0 .

Then z_1, z_2, z_3, z_4, z_5 , then if I read a $g(z)$ as $a_0 + a_1 z_1 + a_2 z_2 + a_3 z_3 + a_4 z_4 + a_5 z_5$, this is a linear discriminant function in the $\phi(x)$ space, right? So, what is linear in terms of the components of z would be non linear in terms of the components of x . This is the basic idea. So, if I can find the right transformation, can use techniques for learning linear models for learning non linear classifiers in the original space.

(Refer Slide Time: 04:23)



- There are two major issues in naively using this idea.
- If we want, e.g., p^{th} degree polynomial discriminant function in the original feature space (\mathbb{R}^m), then the transformed feature vector, Z , has dimension $O(m^p)$.
- Results in huge computational cost both for learning and final operation of the classifier.
- We need to learn $O(m^p)$ parameters rather than $O(m)$ parameters. Hence may need much larger number of examples for achieving proper generalization.
- SVM offers an elegant solution to both.

NPTEL

PR NPTEL course - p13/28

Of course, there are as we mentioned if I just naively try this idea, there are two major problems, right? By naively try I mean simply think of a function ϕ transform all x 's into z 's and then on use any linear method that will act. What is the problem? As we can see from the previous example, when I wanted quadratic, I need all quadratic terms in the in the components of z , right? So, I need $x_1^2, x_2^2, x_1 x_2$ and so on. So, essentially if I need a p^{th} degree polynomial at the discriminant function in the original space.

Then in the in the transformed vector I should have terms that contain all p^{th} degree terms in components of x , which means they'll be order m to the power p terms in the

vector z . So, z would have very huge dimension. So, even in the quadratic case I have taken 2, but if I have got let us say a 100 dimensional vector which is not very uncommon pattern recognition, then I will have of the order of 100 square components in z ; that is 10,000. On the other hand if I want a cubic discriminant function, discriminant function up to third degree polynomial, then I would not have 10 power 4, but, I would have 10 power 6 components in z a billion components in z instead of being a 100 components in z .

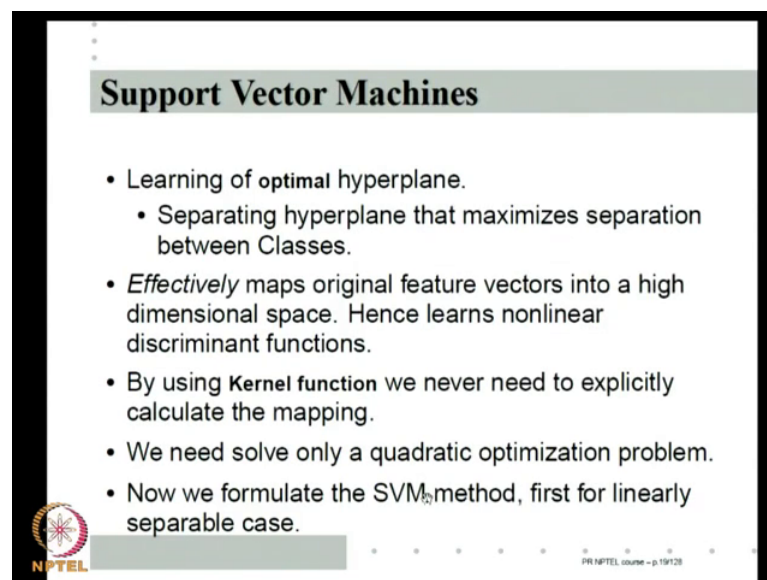
So, the the dimension of z can be very huge even for second and third degree polynomial discriminant functions that we want to learn. They are not particularly non linear just quadratic or cubic functions in the original feature space. The transform space can have huge dimension and this can results in huge computational cost. Let us say what do I have to do when to classify? I have to calculate $w^T x$, if x is a feature vector for any linear classifier I essentially calculate something like $w^T x$ in the original space. It will be in only 100 multiplications, if it is a 100 diameter space. But if you are learning a cubic classifier a a cubic discriminant function it will be a million component z , so it will mean a million multiplication.

Every time I do a simple $w^T x$ operation it is a million multiplications, right? So, learning the classifier as well as the final operation of the classifier can be very costly. Now, during learning I might have to store, let us say a 1000 100 dimensional vectors, instead of that now to store a 1000 million dimensional vector, right? That is that is a that is a you know a 10,000 fold increase in the memory needed. So, it it it incurs large computational cost not only that to learn the linear classifier in this new space. I have to learn order m power p parameters, w will be the same order as the feature vector.

So, if I am learning a linear classifier in the jet space the w in that linear classifier should have the same dimension as the vector z . So, that means I have to learn order m power p parameters. If I am learning a p th degree polynomial in the original space. If I am learning a linear in the original space, I will learn only order m because instead of learning a 100 parameters for a hyperplane in in the 100 dimensional feature space, I am learning a million parameters, because now my feature space has become a million dimensional. Now, to learn 100 parameters based on v c dimension hyperplanes, v c dimensional hyperplanes is of the order of the dimension of the space is actually dimension of the space plus 1.

So, in the 100 dimensional space the family of hyperplanes should have v c dimension of 100, so maybe I need 1000 examples here. Because this space has dimension 1 million I may need 10 million examples. 10 million examples are huge, we cannot get them is even getting 1000 examples may be difficult, but 10 million examples is much more difficult, right? Now, this is also a a a a an extra problem, so basically both in terms of computation and in terms of being able to learn well with the available examples. This this new idea will not work. So, the nice thing about what we call support vector machines is that they offer an elegant solution for both these problems, so let us look at what these are?

(Refer Slide Time: 09:12)



Support Vector Machines

- Learning of optimal hyperplane.
 - Separating hyperplane that maximizes separation between Classes.
- *Effectively* maps original feature vectors into a high dimensional space. Hence learns nonlinear discriminant functions.
- By using Kernel function we never need to explicitly calculate the mapping.
- We need solve only a quadratic optimization problem.
- Now we formulate the SVM method, first for linearly separable case.

NPTEL logo and footer text: PPR NPTEL course - p.19/128

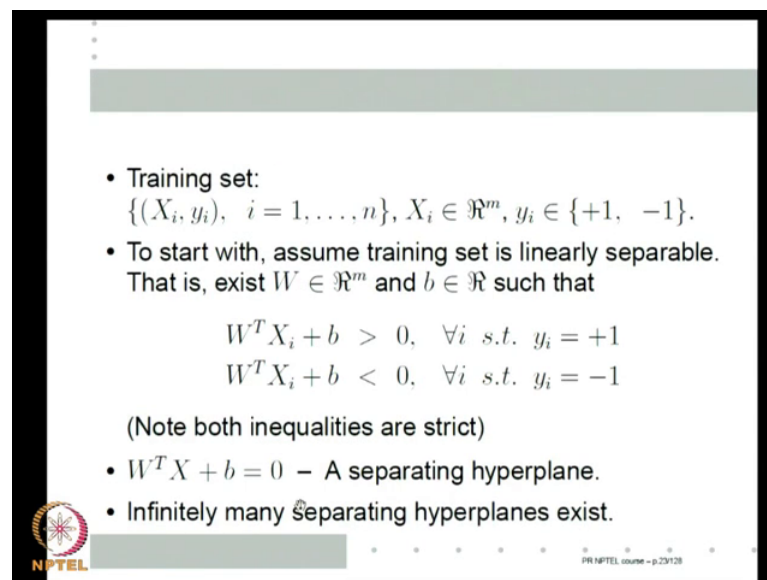
They essentially learn not any hyperplane, but what we term optimal hyperplane? An optimal hyperplane is a separating hyperplane that maximizes separation between classes, we will formally define it in a couple of minutes. Now, learning optimal hyperplane is essentially what allows SVM's to learn well with fewer examples even in a very large dimensional space.

Then they only effectively map the original feature vectors into high dimensional space. Of course, they learn a non linear discriminant function by learning a linear discriminant or or linear classifier in a high dimensional space, but this is done using what is known as a Kernel function, whereby we never actually explicitly do this mapping. We never calculate ϕ of x , we never work in the range space of ϕ , we shall see why, how it is

feasible? By using Kernel functions, we never actually explicitly calculate this mapping. So, we only effectively find a linear classifier in the high dimensional space.

So, ultimately as it turns out for the sake of solving a quadratic optimization problem we will look at the dimension of the quadratic optimization. It turns out to be dimension of the number of examples rather than the dimension of the feature vector. But we solve only a quadratic optimization problem and are able to do this mapping into high dimensional space and learning a linear classifier there. So, with that introduction, let us move on to formulating the SVM method. First we look at formulating it for the linearly separable case and then we look at the more general case, okay?

(Refer Slide Time: 10:55)



• Training set:
 $\{(X_i, y_i), i = 1, \dots, n\}, X_i \in \mathbb{R}^m, y_i \in \{+1, -1\}$.


• To start with, assume training set is linearly separable. That is, exist $W \in \mathbb{R}^m$ and $b \in \mathbb{R}$ such that

$$W^T X_i + b > 0, \quad \forall i \text{ s.t. } y_i = +1$$
$$W^T X_i + b < 0, \quad \forall i \text{ s.t. } y_i = -1$$

(Note both inequalities are strict)

• $W^T X + b = 0$ – A separating hyperplane.

• Infinitely many separating hyperplanes exist.

 NPTEL PR NPTEL course - p.25128

So, this is the training set given to us we will revert back to our old training set notation. Remember when we are considering all linear models we are using x's for feature vector x or and y's for targets or classes and small n as the number of examples. So, earlier when we considered neural networks because y is happen to be outputs of the network nodes, we used d i as the desired output in the training set. Now, that we do not have that problem we will revert back to the originals notation we have.

So, x i's are the inputs and y i's are the targets of the class labels. So, we are given n examples x i are in m dimensional input space. So, we will take m to be the dimension of the feature vector space and we choose class labels for here to be plus 1 minus 1 earlier when we considered linear classifiers, we mostly chosen 0 and 1 as the two class labels.

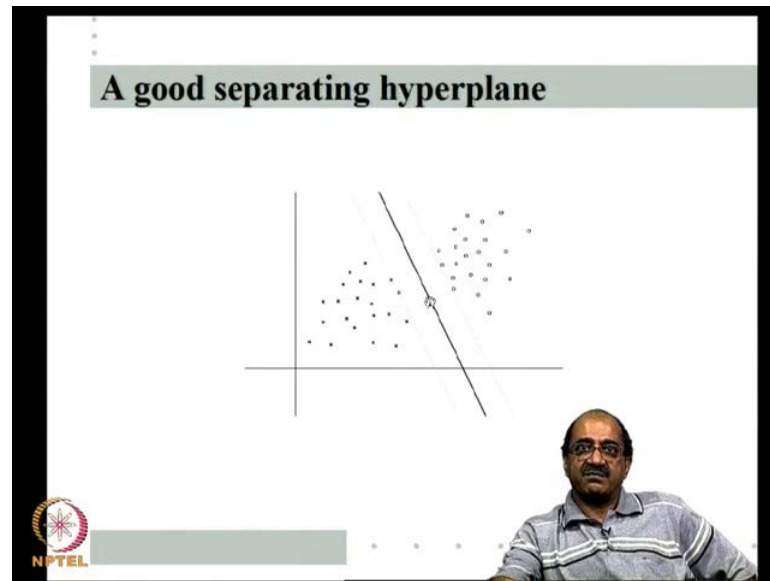
But we will take plus 1 and minus 1, we are only going to look at two class problem. Multiclass will be solved much like in the linear models by using you know one of the many techniques as we discussed there.

So, we will consider only the two class problem, the two class will be taken as plus 1 minus 1. To start with we're assuming that the training set is linearly separable. That means there is a weight vector w belonging to \mathbb{R}^m and a constant b belonging to \mathbb{R} . Of course, there are such that whenever y_i is plus 1; that means x_i belongs to the class plus one $w^T x_i + b$ is greater than 0. Whenever y_i is minus 1, we have $w^T x_i + b$ less than 0, so for all class one patterns $w^T x_i + b$ greater than 0. All class minus 1 patterns $w^T x_i + b$ less than 0.

Because I have only finitely many patterns, if there are separable. They will be exactly separable, right? So, no pattern need be on the separating hyperplane, so we will take both inequalities to be strict. Now, of course, we know we have already seen when we considered perceptron that such a thing is called a separating hyperplane, such as $w^T x + b = 0$ is the equation of a separating hyperplane as we seen when we discussed perceptron. If there is one separating hyperplane, if the if the training data is linearly separable, then there are infinitely many separating hyperplanes, right?

So, the question comes which is a good separating hyperplane? Earlier say for example, when we did perceptron all we are bothered is the about finding any separating hyperplane; that is all we actually proved about perceptron that it can find a separating hyperplane. But because we know that exist infinitely many separating hyperplanes, can we say that some hyperplanes have better than others?

(Refer Slide Time: 13:37)

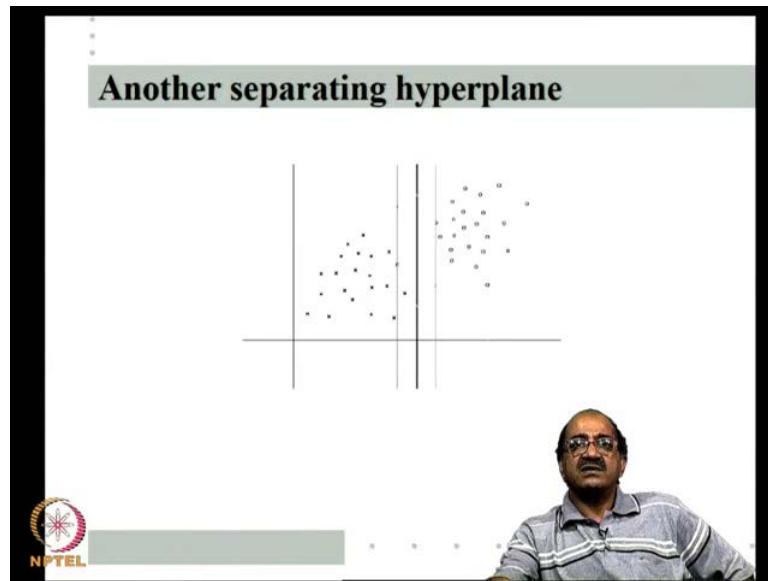


So, here is one way of looking at it look at that dark line in the center for now forget about the two faint lines on either side. Let us say for this problem we would like to say that this is a good separating hyperplane, whereas this is not so good hyperplane. Why do we say, why do why why do you want to say that? Essentially looking at this spread the classes seem to be spread like this, so this line seem to be exactly you know bisect the main line along which classes are spread, right?

But that is intuitive that is more difficult to capture, but another way of seeing it is you can see the nearest patterns on either side of the hyperplane are so far apart for this hyperplane. Whereas, for this hyperplane the nearest patterns are much closer to the hyperplane, so this hyperplane is making more commitment about unseen parts of the space compare to this hyperplane. This hyperplane is keeping itself as forever from both side patterns as possible. Hence, getting lot of what I may call marginal separation, what we will actually ultimately call margin of separation, right?

So, on either side of the hyperplanes that I am finding the the patterns are sufficiently far away. That gives me some intuitive feel that this might be a better hyperplane. So, this is essentially what is one intuitive way of thinking good hyperplane as it turns out even more formal is a good separating hyperplane. If the closest pattern to the hyperplane on an either side is far away for one hyperplane, then that hyperplane is better than another hyperplane for which it is closer.

(Refer Slide Time: 15:20)



As here it is too close, so this this is not as good as separating hyperplane as that one.

(Refer Slide Time: 15:30)

- Recall that we assume training set is linearly separable and hence
$$W^T X_i + b > 0, \quad \forall i \text{ s.t. } y_i = +1$$
$$W^T X_i + b < 0, \quad \forall i \text{ s.t. } y_i = -1$$
- Since the training set is finite, $\exists \epsilon > 0$ s.t.
$$W^T X_i + b \geq \epsilon, \quad \forall i \text{ s.t. } y_i = +1$$
$$W^T X_i + b \leq -\epsilon, \quad \forall i \text{ s.t. } y_i = -1$$

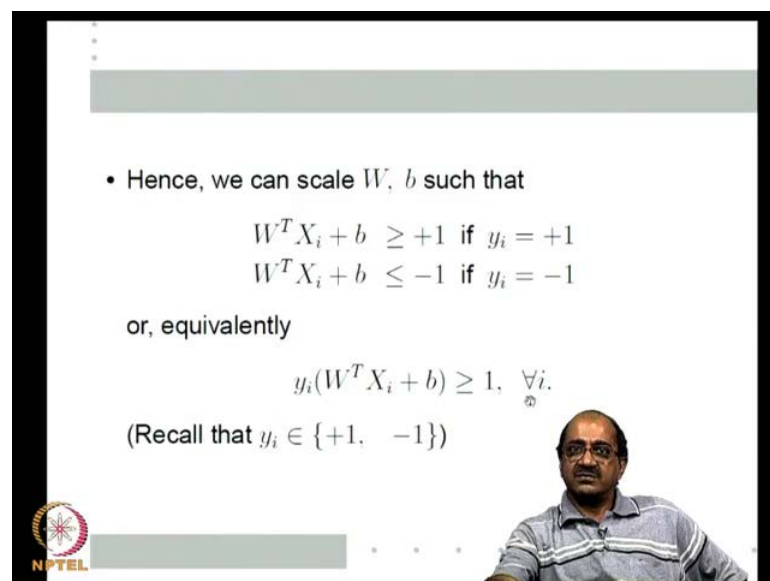
So, let us now try to formalize this notion that the closest pattern should be far away. Recall that the training set is linearly separable and hence we have this W, b . Now, because the training set is finite and for each of the class plus one pattern is strictly greater than 0, each of the class minus 1 pattern is strictly less than 0. So, if I take minimum of this over all class one patterns, it will be some epsilon 1, which is greater than 0 strictly. If you take maximum of this over all class minus one patterns, this some

minus epsilon 2, where epsilon 2 is positive.

If I take smaller of the epsilon 1 epsilon 2 is epsilons, I can always satisfy $w^T x_i + b$ is greater than or equal to epsilon. If x_i is class plus 1 less than or equal to minus epsilon, if it is minus 1. Because this they are only finitely many patterns and if they all if $w^T x_i + b$ is positive for all of them, there must be a smallest number for it. We have will calling that epsilon 1. Similarly the largest number for this we are calling that minus epsilon 2 and then we are taking smaller of epsilon 1 or epsilon 2.

So, that we can satisfy this, once we can satisfy this, now I can now divide the center equation by epsilon. Because it is a positive number, right? Dividing the the inequality by epsilon simply means, I am scaling W on b . Now, if $W^T x_i + b$ is equal to 0 a separating hyperplane, then $k w^T x_i + k b$ is also a separating hyperplane. So, $k W$ plus $k b$ will also be a separating hyperplane.

(Refer Slide Time: 17:11)



• Hence, we can scale W, b such that

$$W^T X_i + b \geq +1 \text{ if } y_i = +1$$
$$W^T X_i + b \leq -1 \text{ if } y_i = -1$$

or, equivalently

$$y_i(W^T X_i + b) \geq 1, \forall i.$$

(Recall that $y_i \in \{+1, -1\}$)

So, I just I can I can scale the W on b to satisfy $W^T x_i + b$ is greater than plus 1, if x_i is in class plus 1 is less than minus 1. If x_i is in class minus 1, what is this mean? Now, we can rewrite this as because y_i which is chosen to be plus 1 minus 1. I can write both these inequalities together compactly as y_i into $w^T x_i + b$ greater than equal to 1.

Actually to be able to write like this and hence simplify our algebra is why we have

taken y_i to be plus 1 minus 1 in this case. So, we can always scale any separating hyperplane, so that it will satisfy y_i into $w^T x_i + b$ is greater than equal to 1. What does it mean?

(Refer Slide Time: 17:53)

• When the training set is separable, any separating hyperplane, W, b , can be scaled to satisfy

$$y_i(W^T X_i + b) \geq 1, \forall i.$$

• Then there are no training patterns between the two parallel hyperplanes

$$W^T X + b = +1$$

and

$$W^T X + b = -1$$

NPTEL

When training set is separable, any separating hyperplane W, b can be scaled. That is we can always take a W, b to satisfy y_i into $w^T x_i + b$ greater than or equal to 1. What is that mean? If x_i is class plus 1, then $w^T x_i + b$ is greater than plus 1. If it is class minus 1, $w^T x_i + b$ is less than minus 1, which means for no training pattern $w^T x_i + b$ is between plus 1 and minus 1, which is same as saying there is no training patterns between the two parallel hyperplanes.

$w^T x + b$ is equal to plus 1 $w^T x + b$ is equal to minus 2. $w^T x + b$ is equal to 0 is my hyperplane. But for all class one pattern $w^T x + b$ is greater than 1 greater than or equal to 1. All class minus one patterns $w^T x + b$ is less than or equal to minus 1, which means if I take the two parallel hyperplanes $w^T x + b$ is equal to plus 1 and $w^T x + b$ is equal to minus 1, there will be no training pattern between them.

(Refer Slide Time: 18:55)

Optimal hyperplane

- Distance between these two hyperplanes is: $\frac{2}{\|W\|}$.
Called margin of the separating hyperplane.
- Hence distance between the hyperplane and the closest pattern is $\frac{1}{\|W\|}$.
- Intuitively, more the margin, better is the chance of correct classification of new patterns.
- **Optimal Hyperplane** – separating hyperplane with maximum margin.

NPTEL logo and footer text: PR NPTEL course - p.37128

Now, the distance between these two hyperplanes will be 2 by norm W. Norm W is W transpose W square, norm W square is W transpose W. This is a simple coordinate geometry problem i am sure u can easily do this.

(Refer Slide Time: 19:08)

Margin of a hyperplane

$m = \frac{2}{\|w\|}$

Class 1, Class 2, $w^T x + b = 1$, $w^T x + b = -1$, $w^T x + b = 0$

NPTEL logo and a small video inset of a man speaking.

So, let us look at what this means see this is the W transpose x plus b is equal to 0 hyperplane. So, W transpose x equal to plus 1 and minus 1 are on either side of this hyperplane, because there they will be parallel. They will be just the hyperplane moved parallel to itself on either side. So, this one will be W transpose x plus b is equal to plus 1

other will be $W^T x + b = -1$. What is shown as m here that is the margin. Simple coordinate geometry will show you that this margin is $2 / \|W\|$, right? So, actually W is as you know a hyperplane, because hyperplane equation is $W^T x + b = 0$, W is a normal vector to the hyperplane, right?

So, that that is why how we shown the W there. Because the distance between the two hyperplanes is $2 / \|W\|$ distance between the hyperplane on the closest pattern is $1 / \|W\|$. Because we we are only guaranteeing that between $W^T x + b = 1$ and $W^T x + b = -1$. I will not be able to they will not be any training pattern of saying $W^T x + b$ can be just can be scaled up to this cannot be scaled beyond this. If I scale beyond this, it will go out. So, ultimately the closest patterns will be either on $W^T x + b = 1$ or $W^T x + b = -1$ as you, as as in this figure.

So, if I take this orientation of the hyperplane, then this is the separation I can get. So, this will be symmetric on either side. So, the distance between hyperplane and the closest pattern will be one by instead of $2 / \|W\|$. It will be $1 / \|W\|$. See in general, because both sides the patterns need not have to be symmetric, one side patterns may not actually be on this hyperplane by the time I move this hyperplane this much. I may have it, this patterns on this side the patterns may be far away, but I cannot hit any more because these two hyperplanes have to be $W^T x + b = 1$ and $W^T x + b = -1$.

So, ultimately if margin is $2 / \|W\|$, I know at a distance one by norm W , there is a pattern and no pattern is closer than $1 / \|W\|$, right? So, distance between the hyperplane and the closest pattern one by norm W intuitively more the margin better the chance of correct classification of new patterns, right? (()) if if the margin is more, I am leaving much more space for me for some more patterns to come here. So, given that these are i and d examples very unlikely that a pattern from the this class come this side, so more the margin intuitively better is my hyperplane.

So, essentially you want a hyperplane that that has maximum margin. So, we will define optimal hyperplane to be a separating hyperplane with maximum margin. So, we want to find a separating hyperplane W, b for which $1 / \|W\|$ is maximum or norm W is minimum or norm W^2 is minimum or $W^T W$ is minimum, right? I

want to maximize margin, so I have to maximize $1/\|W\|$, which is same as minimizing $\|W\|$ or which is same as minimizing $\|W\|^2$. So, I essentially want a separating hyperplane for which the margin is maximum or $\|W\|$ is minimum, okay?

(Refer Slide Time: 22:38)

The optimization problem

- Among all separating hyperplanes, the one with largest margin is the optimal hyperplane.
- So, the optimal hyperplane is a solution to the following optimization problem.
- Find $W \in \mathbb{R}^m, b \in \mathbb{R}$ to

$$\begin{aligned} &\text{minimize} && \frac{1}{2}W^TW \\ &\text{subject to} && y_i(W^T X_i + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$
- This is a constrained optimization problem with quadratic cost function and linear inequality constraints.

NPTEL PR NPTEL course - p.41/28

So, let us formulate this among all separating hyperplanes the one with the largest margin is the optimal hyperplane. So, we can now write this as an optimization problem what is the optimization as we seen. We want something that has separating hyperplane that has least margin. That is maximum margin or least $\|W\|$. So, finding an optimal hyperplane is same as finding a W on b W belong to \mathbb{R}^m and b belong to \mathbb{R} to minimize $\|W\|^2$ subject to the constraint that this W and b form a separating hyperplane for them to form a separating hyperplane. They have to satisfy $y_i(W^T X_i + b) \geq 1$ i is the index of the examples, so i goes from 1 to n .

So, for each example x_i y_i into W each example x_i y_i into $W^T x_i + b$ should be greater than or equal to 1. This tells me that the hyperplane given by W would not be correctly separates correctly. I mean no pattern in between $W^T x_i + b$ is equal to plus 1 and minus 1 among all. Say W, b we satisfy this, I want that one which has the least $\|W\|^2$ or maximum margin, right?

Such things are called constrained optimization problems, right? We want to minimize

this, but do not want to minimize this over all W and b , but only over those W and b which also satisfy this. So, among all W b that satisfy this constraints which is the one that minimizes this. Such optimization problems are called constrained optimization problems. These are unlike the optimization, we come across so far where when we are minimizing empirical risk is just some function to minimize is like finding maximum or minimum of a of a high dimensional function with just a function f of x .

You are asking which x will minimize this among all x in in the space is now no longer that right among all like that satisfies some constraints, which one will minimize this function such thing are called constrained optimization problems. Since, some of you may not have done constrained optimization earlier I will give you a very brief review of constrained optimization.

(Refer Slide Time: 24:52)

Constrained Optimization

- Consider the following optimization problem

$$\begin{aligned} &\text{minimize} && f(x) \\ &\text{subject to} && \mathbf{a}_j^T \mathbf{x} + b_j \leq 0, \quad j = 1, \dots, r \end{aligned}$$

where $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is a continuously differentiable function, and $\mathbf{a}_j \in \mathbb{R}^m, b_j \in \mathbb{R}, j = 1, \dots, r$.

- A point, $\mathbf{x} \in \mathbb{R}^m$, is called a **feasible point** (for this problem) if $\mathbf{a}_j^T \mathbf{x} + b_j \leq 0, j = 1, \dots, r$.

NPTEL PR NPTEL course - p.43/28

So, any general constraint optimization problem is of the following kind minimize some function f of x subject to some constraints. We are only interested in what are called linear constraints. We will take them to be linear constraint. So minimize f of x subject to \mathbf{a}_j transpose x plus b_j less than or equal to 0 j is equal to 1 to r . So, I have r constraints \mathbf{a}_1 transpose x plus b_1 should be less than 0 \mathbf{a}_2 transpose x plus b_2 should be less than or equal to 0 and so on so forth. f is a function that maps \mathbb{R}^m to \mathbb{R} and we assume it to be continuously differentiable.

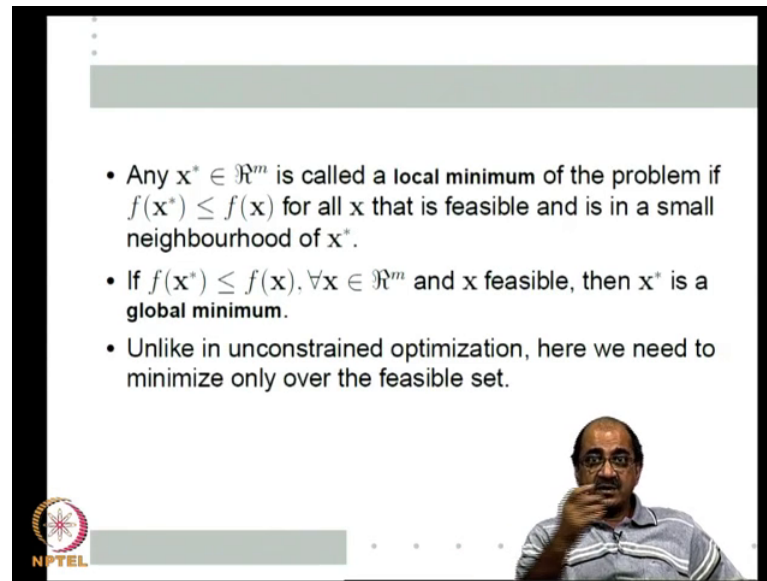
These \mathbf{a}_j vectors are also in \mathbb{R}^m , because I want to calculate \mathbf{a}_j transpose x . So, this is some

linear function of x a j affine function of x . Normally you call this as the objective function or cost function or criterion function. So, you want to minimize a cost or a criterion function objective function subject to these are called the constraints, right? So, what it means is search only among those x in \mathbb{R}^m , which satisfy all these constraints. Among all those x find the x that minimizes f that is the problem such a problem is called a constrained optimization problem, okay?

See in general in an unconditional optimization problem, we have no constraints from just finding an x over \mathbb{R}^m that locally minimizes f or globally minimizes f a simplest way of looking. It is sometimes there might be a solution for a constrained optimization problem, but there will be no solution to unconstrained optimization problem. Suppose, I want to minimize $3x$ now, there is no no minimum. We can take x as low as you want and I keep going down, so ultimately I hit minus infinity, right? This on the other hand, suppose I want to minimize $3x$ subject to x greater than 0 , then obviously the minimum with 0 minimum x is right.

So, essentially that is what is meant by when I am, when I want to minimize $3x$. I have to search over all possible x , so I have no no minimum, but when I say subject to minimize $3x$ subject to x greater than 0 , I am I need to only search over those x , which satisfy my constraints a point x in \mathbb{R}^m , which satisfies all the constraints is called a feasible point. So, given this constrained optimization problem or optimization problem any x that satisfies the constraints is called a feasible point. So, any x that satisfy a j transpose x plus b_j less than or equal to 0 j is equal to 1 to r . All the constraints have to be satisfied then it is called a feasible point. So, essentially we are going to minimize f over those x that are feasible, okay?

(Refer Slide Time: 27:32)



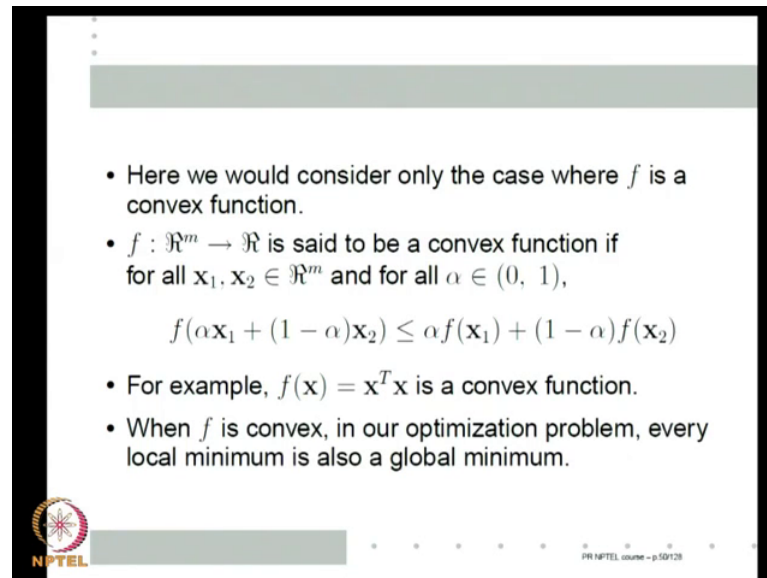
- Any $x^* \in \mathcal{R}^m$ is called a **local minimum** of the problem if $f(x^*) \leq f(x)$ for all x that is feasible and is in a small neighbourhood of x^* .
- If $f(x^*) \leq f(x), \forall x \in \mathcal{R}^m$ and x feasible, then x^* is a **global minimum**.
- Unlike in unconstrained optimization, here we need to minimize only over the feasible set.

Now, we can define what is a solution any x^* in \mathcal{R}^m is called a local minimum of this problem or local solution to this problem. If $f(x^*) \leq f(x)$ for all x that is feasible and is in a small neighborhood of x^* for an unconstrained problem, we define x^* to be a local minimum. If $f(x^*) \leq f(x)$ for all x in a neighborhood around x^* . Now, I want $f(x^*) \leq f(x)$ for all x in a neighborhood around x^* , but not every x in neighborhood need to satisfy this only all feasible x in a small neighborhood around x^* need to satisfy this.

Similarly, for a global solution I want this to be satisfied by all x in \mathcal{R}^m , but only those x which are feasible. So, normally with unconstrained minimization x^* is a global minimum if $f(x^*) \leq f(x)$ for all x in \mathcal{R}^m . Now, for the constrained problem which is a global minimum if $f(x^*) \leq f(x)$ for for any x in \mathcal{R}^m that is feasible x in \mathcal{R}^m that are not feasible. Do not have to satisfy this; that essentially what a constrained optimization problem will be unlike in unconstrained optimization here. We need to minimize only over the feasible set.

Now, how does one solve this? No ,unconstrained optimization problem, we can differentiate and equate to 0; that is what we have been doing, right? The gradient descent is a simple method to solve the unconstrained optimization problem, how do we solve the constrained optimization problem?

(Refer Slide Time: 29:01)



• Here we would consider only the case where f is a convex function.

• $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is said to be a convex function if for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^m$ and for all $\alpha \in (0, 1)$,

$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$$

• For example, $f(\mathbf{x}) = \mathbf{x}^T \mathbf{x}$ is a convex function.

• When f is convex, in our optimization problem, every local minimum is also a global minimum.

NPTEL PR NPTEL course - p.50128

Now, for us for our case we need to for the for the SVM algorithms that we are going to discuss. We need to look at only constrained optimization problems whose constraints are linear and whose cost function or criterion function is convex. So, we will restrict ourselves to only convex functions just in case some of you do not know what a convex function is? A function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is said to be convex if for every $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^m$ and all real numbers α in the open interval $(0, 1)$; that is α strictly between 0 and 1.

of course, it can even be close interval $[0, 1]$, but still does not matter $f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$. That is $f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2)$ is less than or equal to $\alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$ geometrically. What this means is, if I look at the graph of $f : \mathbb{R}^m \rightarrow \mathbb{R}$ function instead of $\mathbb{R}^m \rightarrow \mathbb{R}$, we look at the graph of f . If I take two points; $\mathbf{x}_1, \mathbf{x}_2$ and I join them on the graph. Then the actual graph should be below the line that joins is, but anyways this is the condition for convexity.

For example, $\mathbf{x}^T \mathbf{x}$ is a convex function you can easily verify this when f is convex. The the reason why convexity helps is in the optimization problem every local minimum is also a global minimum. You do not have to worry about local minimum problems when f is convex and the constraints are linear. In that case in the optimization problem every local minimum is also a global minimum. That is a that is a huge advantage when solving an optimization problem, okay?

(Refer Slide Time: 30:45)

• We now look at one method of solving the constrained optimization problem.

• Given our optimization problem, define

$$L(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^r \mu_j (\mathbf{a}_j^T \mathbf{x} + b_j)$$

NPTEL

Now, how do I solve this constrained optimization problem? Given the optimization problem, we define a function L is a function of your \mathbf{x} \mathbf{x} is an $r \times m$ vector and also a vector $\boldsymbol{\mu}$ $\boldsymbol{\mu}$ is has many components as there constraints. They are all constraints, so this vector $\boldsymbol{\mu}$ is μ_1, μ_2 .

(Refer Slide Time: 31:19)

Constrained Optimization

• Consider the following optimization problem

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{a}_j^T \mathbf{x} + b_j \leq 0, \quad j = 1, \dots, r \end{aligned}$$

where $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is a continuously differentiable function, and $\mathbf{a}_j \in \mathbb{R}^m, b_j \in \mathbb{R}, j = 1, \dots, r$.

• A point, $\mathbf{x} \in \mathbb{R}^m$, is called a **feasible point** (of this problem) if $\mathbf{a}_j^T \mathbf{x} + b_j \leq 0, j = 1, \dots, r$.

NPTEL

Just a small notational comment just for the optimization tutorials, so to say we will use the notation that all bold figures are \mathbf{x} . See so far we have not been distinguishing between vectors and scalars. But just for this optimization part just this three four slides

where I discussed optimization. I will use all bold face characters for vectors, so this \mathbf{x} 's are vectors \mathbf{a}_j 's are vectors, but b_j 's are scalars adjust. So, that is a little more easier when we revert back to our pattern recognition. Of course, there we know capital X is feature vector capital W is your weight vector and so on.

So, there we do not use any notation like earlier which are vectors say scalars is clear from context, because this an abstract optimization problem. I decided to use bold face for all vectors, you will see throughout as and when I am discussing this abstract optimization problem. So, this μ is another vector, it has r components $\mu_1, \mu_2, \dots, \mu_r$. So, you define a function L a function of \mathbf{x} and μ as f of \mathbf{x} my criterion function plus $\sum_{j=1}^r \mu_j (\mathbf{a}_j^T \mathbf{x} + b_j)$ is equal to 1 to all where r is the number of constraints μ_j times \mathbf{a}_j transpose \mathbf{x} plus b_j .

I am writing all my constraints as what are called inequality constraints my constraints are all of the form $\mathbf{a}_j^T \mathbf{x} + b_j \leq 0$. So, that L function I am writing is f of \mathbf{x} plus some μ_j times \mathbf{a}_j transpose \mathbf{x} plus b_j summed over j is equal to 1 to r such a L is called the Lagrangian for the problem

(Refer Slide Time: 33:00)

• We now look at one method of solving the constrained optimization problem.

• Given our optimization problem, define

$$L(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^r \mu_j (\mathbf{a}_j^T \mathbf{x} + b_j)$$

• The L is called the Lagrangian of the problem and the μ_j are called the Lagrange multipliers.

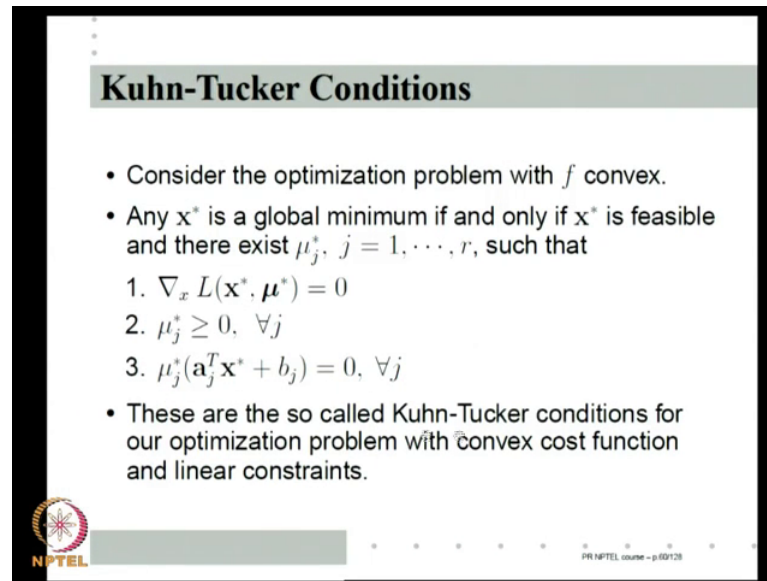
• Essentially, the constrained optimization problem can be solved through unconstrained optimization of L .

NPTEL

PR NPTEL course - p.54128

This μ_j 's are called the Lagrange multipliers, so the Lagrangian essentially adds the criterion function to some unspecified constant times the constraint. The nice thing about this Lagrangian method is that it essentially transform the constrained optimization problem into one of unconstrained optimization of the function L in a in a in a in a in a specific sense, okay? Let us look at this specific sense.

(Refer Slide Time: 33:32)



Kuhn-Tucker Conditions

- Consider the optimization problem with f convex.
- Any x^* is a global minimum if and only if x^* is feasible and there exist μ_j^* , $j = 1, \dots, r$, such that
 1. $\nabla_x L(x^*, \mu^*) = 0$
 2. $\mu_j^* \geq 0, \forall j$
 3. $\mu_j^*(a_j^T x^* + b_j) = 0, \forall j$
- These are the so called Kuhn-Tucker conditions for our optimization problem with convex cost function and linear constraints.

NPTEL PR NPTEL course - p.00128

So, we are considering the old optimization problem that is minimize f subject to linear constraints $a_j^T x + b_j \leq 0$ we are assuming f is convex now here is the theorem that is of interest to us any x^* is a global minimum if and only if x^* is feasible when the x^* satisfies all the constraints. There exist some constants μ_j^* or such constants, which satisfy the following one is that the gradient with respect to x the first component. The first vector variable of the Lagrangian Lagrangian is a function of x on the Lagrange multipliers μ .

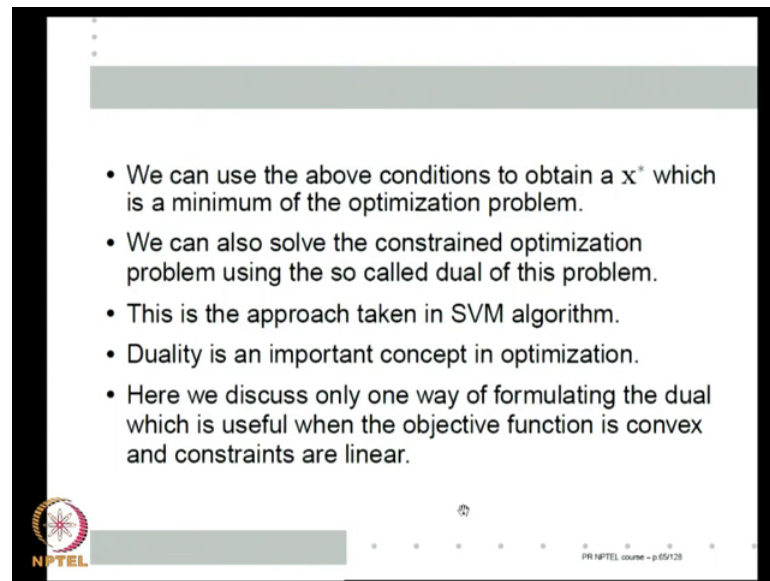
So, the gradient of l with respect to the optimization variables x evaluated x^* μ^* is 0. All the μ_j^* are non negative and $\mu_j^*(a_j^T x^* + b_j) = 0$. Now, we know μ_j^* is to be greater than or equal to 0. So, if μ_j^* is equal to 0, then anyway this is 0, but if μ_j^* is strictly greater than 0, then $a_j^T x^* + b_j$ has to be 0. I know because x^* is feasible $a_j^T x^* + b_j \leq 0$. It may be equal to 0 or it may be strictly less than 0, what the third constraint says is it cannot be strictly less than 0.

If μ_j^* is greater than 0, if μ_j^* is greater than 0, the constraint has to be satisfied by equality. This third condition is often known as a complementary slackness condition, so whenever the Lagrange multiplier μ_j^* are positive. Then the corresponding constraint has to be satisfied by equality right? So, these are called Kuhn Tucker conditions. So, given this optimization problem where the cost function is convex

and the constraints are linear. These are the conditions for any x^* to be a global minimum what are the conditions x^* should be feasible.

There must exist Lagrange multipliers μ_j^* such that gradient with respect to x of L of x^* μ^* is equal to 0. All the μ_j^* are positive and if any μ_j^* is strictly greater than 0, then $a_j^T x^* + b_j$ should be equal to 0, right? This is called the complementary slackness.

(Refer Slide Time: 36:15)



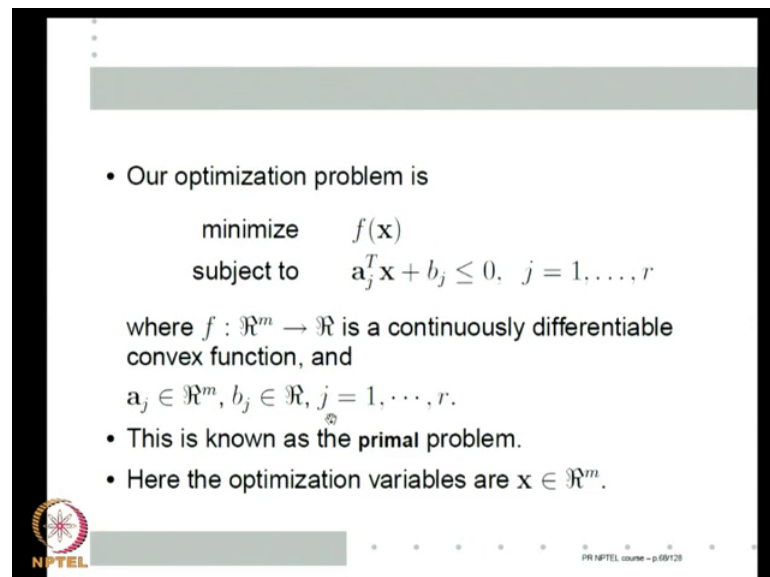
Essentially, because we know all the a_j 's and b_j 's and we can f is given to us. I can calculate L , I can actually you know find the gradient of L equate it to 0. Then make some guesses about which of the constraints are active which are inactive right active meaning when is $a_j^T x^* + b_j$ equally 0 L is strictly less than 0, right? So, I make guesses about which μ_j 's are 0, which if I am decide this μ_j 's are 0, right? Then for the remaining μ_j 's, I have the extra for remaining indices j . I have extra equations $a_j^T x^* + b_j$ is equal to 0, so I can solve $a_j^T x + b_j$ is equal to 0.

This this gradient equal to zero equation together to obtain all the μ_j 's and the x 's, okay? So, there is a way of solving the constrained optimization problem by by finding an x^* , which satisfies this Kuhn Tucker conditions. There are there many other methods there are just like gradient descent. There are other numerical techniques for solving this, but we will look at one particular technique which is used for the support vector machines, which is solving it using. The so called dual of the problem duality is a

is a very very important and a deep concept in optimization.

Once again in this course, we will not be able to teach duality, but those of you may not have gone through an optimization course, we will just look at duality in the very very limited context, right? We will just look at one way of formulating another optimization problem, which we will call the dual of this optimization problem, but even that method will work only for optimization problems, where the cost function is convex and the constraints are linear. So, only in that, only for that case we will we will look at how to find duals, okay?

(Refer Slide Time: 38:21)



• Our optimization problem is

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{a}_j^T \mathbf{x} + b_j \leq 0, \quad j = 1, \dots, r \end{aligned}$$

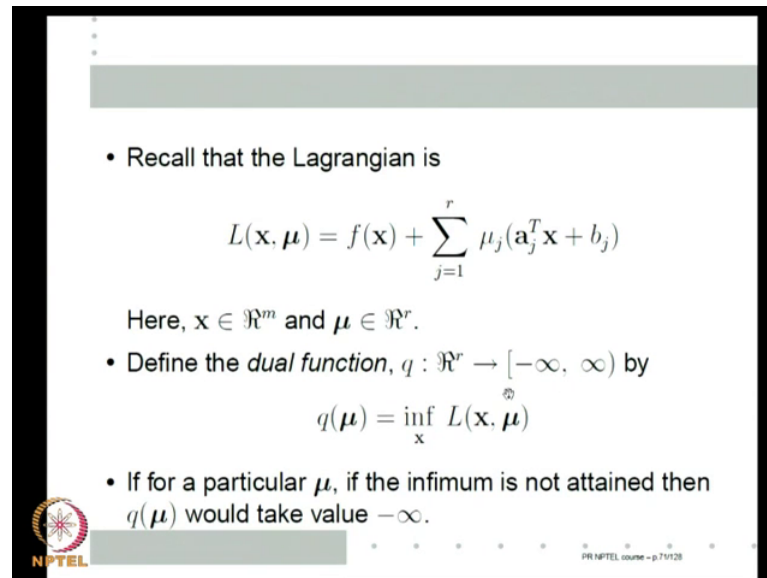
where $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is a continuously differentiable convex function, and
 $\mathbf{a}_j \in \mathbb{R}^m, b_j \in \mathbb{R}, j = 1, \dots, r$.

- This is known as the primal problem.
- Here the optimization variables are $\mathbf{x} \in \mathbb{R}^m$.

NPTEL logo and footer text: PR NPTEL course - p.00120

So, the given optimization problem is to recall minimize f of \mathbf{x} subject to $\mathbf{a}_j^T \mathbf{x} + b_j \leq 0$. We are assuming f is continuously differentiable and also convex now. \mathbf{a}_j 's are vectors in \mathbb{R}^m , so this is we are given r such constraint. This given optimization problem is called the primal problem. In this optimization problem the minimization is over \mathbf{x} . \mathbf{x} is in \mathbb{R}^m , so the optimization variables are \mathbf{x} belonging to \mathbb{R}^m . This is what we are minimizing over so the dimensionality of this optimization problem is m . We are you are searching over \mathbb{R}^m to find the minimum.

(Refer Slide Time: 39:03)



• Recall that the Lagrangian is

$$L(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^r \mu_j (\mathbf{a}_j^T \mathbf{x} + b_j)$$

Here, $\mathbf{x} \in \mathbb{R}^m$ and $\boldsymbol{\mu} \in \mathbb{R}^r$.

• Define the *dual function*, $q : \mathbb{R}^r \rightarrow [-\infty, \infty)$ by

$$q(\boldsymbol{\mu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu})$$

• If for a particular $\boldsymbol{\mu}$, if the infimum is not attained then $q(\boldsymbol{\mu})$ would take value $-\infty$.

NPTEL PRE NPTEL course - p.7/128

Now, this is the Lagrangian $L(\mathbf{x}, \boldsymbol{\mu})$ we have for \mathbf{x} in \mathbb{R}^m and $\boldsymbol{\mu}$ in \mathbb{R}^r this small r is the number of constraints. We have $L(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^r \mu_j (\mathbf{a}_j^T \mathbf{x} + b_j)$ given the Lagrangian. Let us define a function q which we call the dual function, which maps r dimensional real vectors, where r is the number of constraints to not real line. Slightly extended real line minus infinity to infinity minus infinity included infinity not included in real line, even my symphony is not included, it will just allow the q function to sometimes take minus infinity also, as values where q is defined.

Now, q 's domain is \mathbb{R}^r , so q is defined on vectors of the type $\boldsymbol{\mu}$. So, $q(\boldsymbol{\mu})$ is taken to be infimum over all \mathbf{x} in \mathbb{R}^m of $L(\mathbf{x}, \boldsymbol{\mu})$ where $L(\mathbf{x}, \boldsymbol{\mu})$ is this. So, given any $\boldsymbol{\mu}$ in \mathbb{R}^r define the value of $q(\boldsymbol{\mu})$ to be the infimum or minimum over all possible \mathbf{x} of $L(\mathbf{x}, \boldsymbol{\mu})$. The difference between infimum minimum is minima has to attained. So, infimum is the greatest lower bound, we are asking as \mathbf{x} goes over all possible vectors in \mathbb{R}^m . If I take this numbers $L(\mathbf{x}, \boldsymbol{\mu})$, what is the greatest lower bound of this?

So, sometimes it that that is a number that is actually as the value of $L(\mathbf{x}, \boldsymbol{\mu})$ for a particular \mathbf{x} sometimes, it might be a limit, because it can be a limit sometime the infimum may not be attained. Suppose, $L(\mathbf{x}, \boldsymbol{\mu})$ just happens to be minus b times \mathbf{x} or summation summation $\mathbf{x}_i b$ times summation \mathbf{x}_i or suppose is equal to \mathbf{x}_1 just like minimizing $3x$, there will be no minimum. This kind of infimum may not be attained,

right? So, if this is linear in any component of x , this is if there is a linear term in any component of x , then infimum will be minus infinity.

That is the reason we allowed q to take minus infinity as one of the possible values that if a particular μ the infimum is not attained. Then $q(\mu)$ would take value minus infinity just to allow that we define q to be a function from \mathbb{R}^r to \mathbb{R} , because we defined it as infimum. We have to allow the possibility that for some μ $q(\mu)$ may take value minus infinity. Now, so you give me this primal optimization problem, I know f I know a_j b_j . So, I can calculate my Lagrangian. Once I can calculate my Lagrangian for every single μ , I can calculate $q(\mu)$.

(Refer Slide Time: 41:47)

The Dual problem

- The dual problem is:

$$\begin{array}{ll} \text{maximize} & q(\mu) \\ \text{subject to} & \mu_j \geq 0, \quad j = 1, \dots, r \end{array}$$
- This is also a constrained optimization problem.
- Here the optimization is over \mathbb{R}^r and $\mu \in \mathbb{R}^r$ are the optimization variables.
- There is a nice connection between the primal and dual problems.

NPTEL logo and footer: PR NPTEL course - p.75128

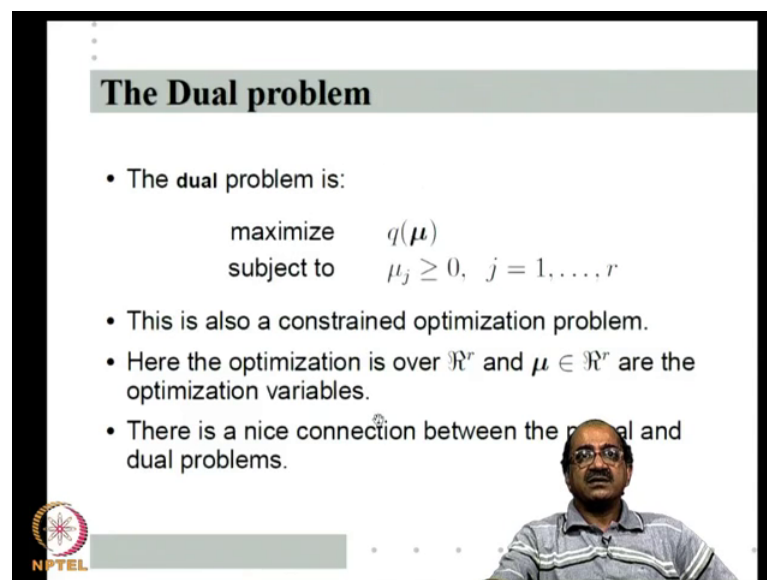
I know the function $q(\mu)$, so now I formulate another optimization problem maximize. This term $q(\mu)$ subject to all the μ 's greater than or equal to 0. This is also a constrained optimization problem, what is a optimization variables here is μ 's. So, the here optimization is over \mathbb{R}^r where little r is the number of constraints we had in the primal. So, here optimization is over \mathbb{R}^r and the optimization variables are μ . This problem is called the dual to the other problem, right?

So, given this optimization problem using these f and a_j 's and b_j 's I formulate the q . Then I formulate the q , then I formulate a new optimization problem, which is also a constrained optimization problem over \mathbb{R}^r with μ as the optimization variables. This problem is called the dual to that primal problem and there are some very nice

connections between the primal and the dual problems. What is the connections if the primal problem has a solution?

Then dual also has a solution and the two optimal values are same, what does that mean? This is the primal problem x^* could be a solution. Then if it is a solution which satisfy the Kuhn Tucker conditions and the optimal values f of x^* here, right? This is the dual problem if the solution exist. Of course, this will have its own Kuhn Tucker conditions with respect to the optimization variables μ . If it has a solution, it could be μ^* and the optimal value will be $q(\mu^*)$.

(Refer Slide Time: 43:35)



The Dual problem

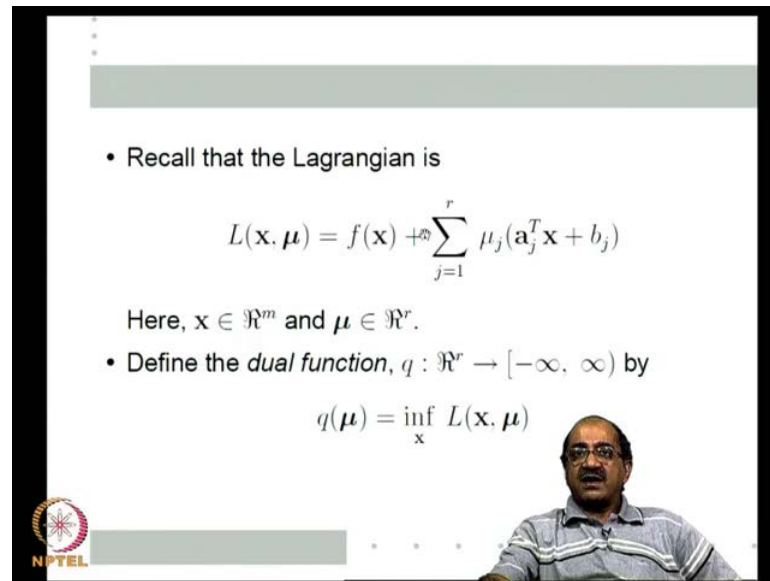
- The dual problem is:
$$\begin{array}{ll} \text{maximize} & q(\boldsymbol{\mu}) \\ \text{subject to} & \mu_j \geq 0, \quad j = 1, \dots, r \end{array}$$
- This is also a constrained optimization problem.
- Here the optimization is over \mathbb{R}^r and $\boldsymbol{\mu} \in \mathbb{R}^r$ are the optimization variables.
- There is a nice connection between the primal and dual problems.

NPTEL

What the theorem says is that if the primal has a solution x^* , then there will also be a μ^* which is a solution to the dual and f of x^* will be equal to q of μ^* . The optimal values are equal 2 for x^* to be an optimal solution for the primal and μ^* to be an optimal for the dual. These are the necessary and sufficient conditions x^* is optimal for primal and μ^* is optimal for dual, if and only if the the two conditions are satisfied.

One is x^* is feasible for primal and μ^* is feasible for the dual obviously it has to be the second condition is f of x^* should be equal to l of x^* μ^* should be equal to $\min_x l(x, \mu^*)$. Yes, while the proof is a little complicated, we can at least see by where the condition comes from now essentially.

(Refer Slide Time: 44:37)



• Recall that the Lagrangian is

$$L(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^r \mu_j (\mathbf{a}_j^T \mathbf{x} + b_j)$$

Here, $\mathbf{x} \in \mathbb{R}^m$ and $\boldsymbol{\mu} \in \mathbb{R}^r$.

• Define the *dual function*, $q : \mathbb{R}^r \rightarrow [-\infty, \infty)$ by

$$q(\boldsymbol{\mu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu})$$

f of \mathbf{x}^* is equal to l of \mathbf{x}^* $\boldsymbol{\mu}^*$, what does this condition mean? If f of \mathbf{x}^* should be equal to l of \mathbf{x}^* $\boldsymbol{\mu}^*$ at \mathbf{x}^* $\boldsymbol{\mu}^*$ this term should be 0, because of feasibility all μ 's are of the same sign $\mathbf{a}_j^T \mathbf{x} + b_j$ are of the same signs. So, nothing can cancel. So, the only way f of \mathbf{x}^* should be equal to l of \mathbf{x}^* $\boldsymbol{\mu}^*$ is if each term here is 0 μ_j^* into $\mathbf{a}_j^T \mathbf{x}^* + b_j$ is equal to 0, which is my complementary slackness condition, right?

Similarly, if I say f of \mathbf{x}^* is equal to minimum over \mathbf{x} l of \mathbf{x} $\boldsymbol{\mu}^*$ that means q $\boldsymbol{\mu}^*$ will be l of \mathbf{x}^* $\boldsymbol{\mu}^*$, which will be equal to f of \mathbf{x}^* , right? So, at least we can see why this conditions are interesting, so this is the primal dual relationship. The primal has a solution so does the dual and the optimal values are equal for \mathbf{x}^* to be optimal for the primal.

$\boldsymbol{\mu}^*$ to be optimal for the dual \mathbf{x}^* should be feasible for the primal $\boldsymbol{\mu}^*$ is feasible for the dual and f of \mathbf{x}^* should be equal to l of \mathbf{x}^* $\boldsymbol{\mu}^*$ should be equal to minimum over \mathbf{x} l of \mathbf{x} $\boldsymbol{\mu}^*$. We would be using the dual formulation for the SVM, okay?

(Refer Slide Time: 45:58)

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

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

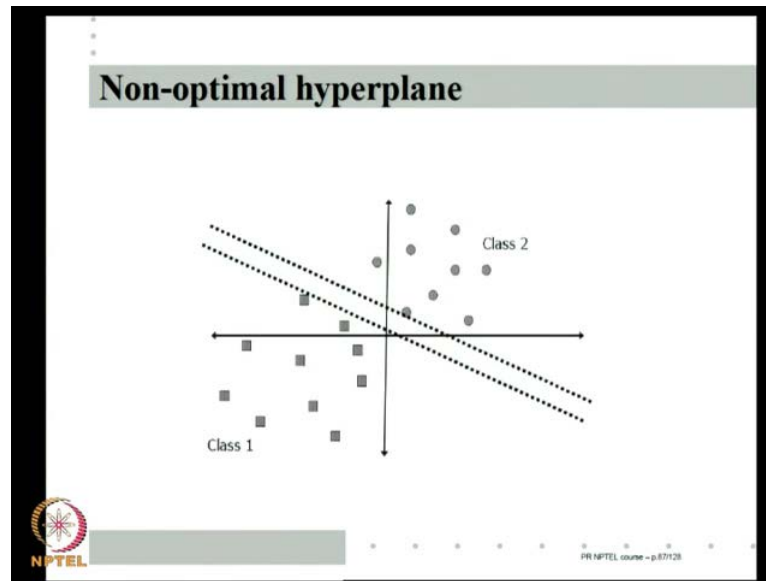
- Quadratic cost function and linear (inequality) constraints.
- Kuhn-Tucker conditions are necessary and sufficient. Every local minimum is global minimum.

NPTEL PR NPTEL course - p.05/128

Now, let us get back to the SVM optimization problem say optimal hyperplane is a solution of what this contained optimal? This is the primal problem minimize half and W transpose W I have written. Now, this in the my standard form I am writing my constraints as less than or equal to 0. So, I wrote it as $1 - y_i(W^T X_i + b) \leq 0$. There n constraints because there are n example. The W has the same dimension as the vector x , which is m . So, I want to minimize half W transpose W , which is a convex function.

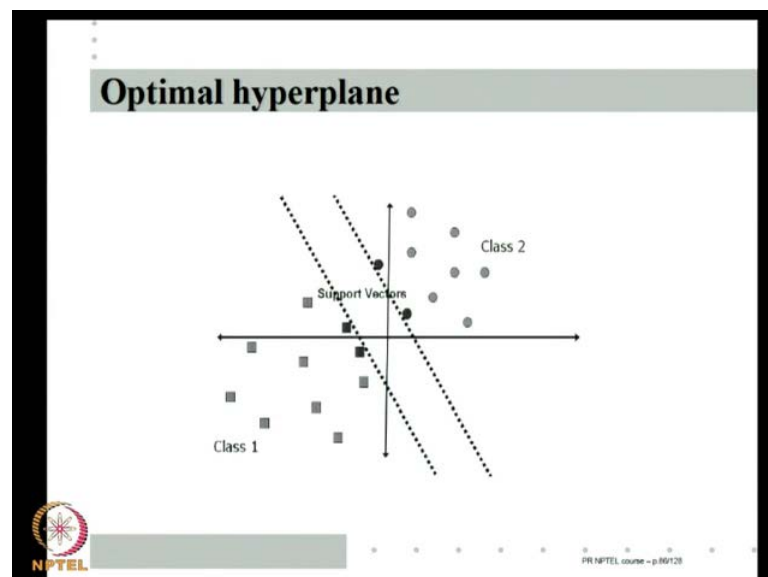
So, my what I called f in the abstract optimization problem is this which is convex and this is our linear constraints $1 - y_i(W^T X_i + b) \leq 0$. There are n constraints, this is a quadratic cost function convex quadratic function is a convex function, linear constraints. Hence, Kuhn Tucker conditions are necessary and sufficient every local minimum is a global minimum. So, essentially to appreciate, what this means is what we are asking for is minimization of W transpose W , which is same as maximization of the margin.

(Refer Slide Time: 47:14)



So, if you have a general thing margin is essentially you put your hyperplane somewhere in the center, so on either side of hyperplane you try to make a margin till you hit one of the patterns. So, in this direction, I get so much margin. On the other hand if I change my direction, there is only so much margin possible, right? So, this is a way of looking at what optimality means in terms of maximizing margin.

(Refer Slide Time: 47:45)



Hyperplane in this direction has low margin as hyperplane in this direction has high margin. So, I am trying to find that separating hyperplane, which has the higher margin

by the, by there are something written as support vectors here. Do not worry about this, we will come back to this figure again and then I will explain that support vectors, okay?

(Refer Slide Time: 48:13)

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

$$\text{minimize } \frac{1}{2}W^T W$$

$$\text{subject to } 1 - y_i(W^T X_i + b) \leq 0, \quad i = 1, \dots, n$$
- Quadratic cost function and linear (inequality) constraints.
- Kuhn-Tucker conditions are necessary and sufficient. Every local minimum is global minimum.

NPTEL PR NPTEL course - p.05128

So, this is my optimization problem. This is f x and these are my constraints.

(Refer Slide Time: 48:19)

- 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

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

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

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

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

NPTEL PR NPTEL course - p.05128

So, my Lagrangian will be my optimization variables are W and b mu's will give me a Lagrangian multipliers. This is the f x, this is the cost function these are the constraints mu i times 1 minus y i into W transpose x i plus b. This is my Lagrangian, so if I now apply my Kuhn Tucker conditions, what does the Kuhn Tucker conditions give first is

that the gradient of l with respect to the optimization variable should be 0. That is gradient with respect to W should be 0 and derivative with respect to partial (\cdot) with respect to b should be 0.

So, first take gradient with respect to W of l this term will give me a W , just so that this term will not give me $2W$, but W . I put a half here right from this term. There is only 1 W term here that is $\sum_i \mu_i y_i W^T X_i$, so gradient of that with respect to W will be $\sum_i \mu_i y_i X_i$, right? That will come with a minus term, so if I equate it to 0, I get W the optimal W is summation over i is equal to 1 to n $\mu_i^* y_i X_i$.

So, by equating gradient of l with respect (\cdot) to 0 I get $w^* = \sum_{i=1}^n \mu_i^* y_i X_i$. If I equate derivative with respect to b to 0 is only one term here that gives me b . That will be $\sum_{i=1}^n \mu_i y_i b$ derivative of that with respect to b will be summation $\mu_i y_i$. So, at optimal Lagrangian multiplies summation i equal to 1 to n $\mu_i^* y_i$ should be equal to 0.

Then Kuhn Tucker conditions also demand feasibility. So, for feasibility I need to have $1 - y_i (W^T X_i + b^*) \geq 0$ for all i and I need complementary slackness. So, $\mu_i^* (1 - y_i (W^T X_i + b^*)) = 0$. All μ_i should be greater than or equal to 0. These are my Kuhn Tucker conditions, okay?

(Refer Slide Time: 50:40)

- 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.

PR NPTEL course - p100128

So, W^* is given by $\mu_i^* y_i x_i$ and this is my complementary slackness condition. Let us say this define the set of indices i for which μ_i^* the corresponding Lagrange multipliers are strictly positive by the set S by complementary slackness. What is complementary slackness give me? $\mu_i^* (1 - y_i) x_i^T W^* + b^* = 0$, so for any i if μ_i^* is strictly greater than 0.

Then $(1 - y_i) x_i^T W^* + b^*$ should be equal to 0. Now, for every i in the set S μ_i^* is strictly greater than 0. That is the definition of the set S , so i belongs to S implies μ_i^* is greater than 0. Hence, $(1 - y_i) x_i^T W^* + b^* = 0$, which is same as $y_i x_i^T W^* + b^* = 1$.

What does that mean? $x_i^T W^* + b^*$ is equal to y_i by multiplying y_i on both sides by y_i . Note that y_i has a plus 1 or minus 1, so y_i^2 is always 1. It essentially means $x_i^T W^* + b^*$ is plus 1, if x_i is in the plus 1 class minus 1. If x_i is in the minus 1 class, so actually x_i is on the hyperplane. $x_i^T W^* + b^*$ is equal to plus 1 or minus 1, which means if μ_i^* is greater than 0, then the corresponding x_i is closest to the separating hyperplane.

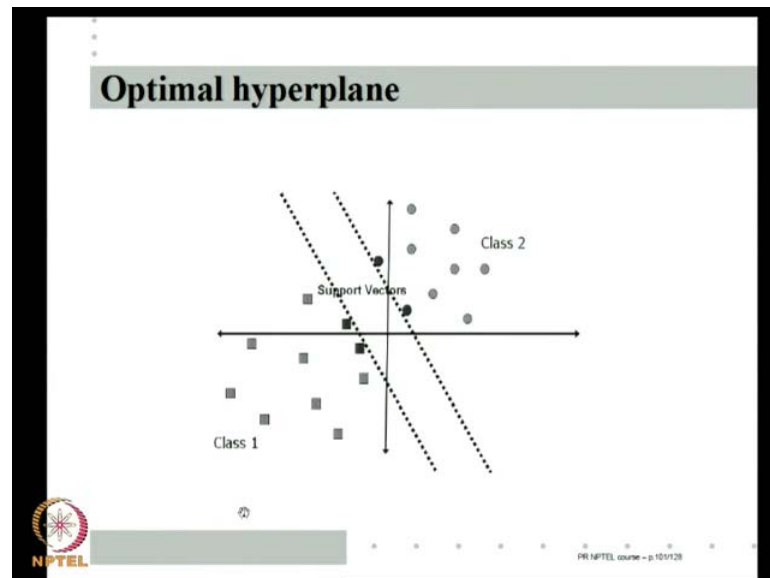
So, we call x_i where i belongs to set S as support vectors see there are as many constraints as there are training examples. So, corresponding each x_i there is a constraint. Hence, corresponding each x_i there is a Lagrange multiplier μ_i^* because there are as many Lagrange multipliers as their constraints. So, we can actually say each Lagrange multiplier is for one example. So, if the Lagrange multiplier is strictly greater than 0, the corresponding example a corresponding x_i is the one closest to the separating hyperplane.

These are called the support vectors and W^* . We know is summation over i $\mu_i^* y_i x_i$ if μ_i^* is equal to 0. Anyway that term does not contribute, so you can actually take W^* to be summation over i belonging to S $\mu_i^* y_i x_i$, which means the optimal W is a linear combination of the support vectors. Of course, when we did perceptron, we know we knew that essentially a separating hyperplane will be a linear combination of the x_i of the examples. But we know that the optimal hyperplane now is a specific linear combination of only some of the examples not all the examples.

These examples are called the support vectors and those happen to be the ones closest to the hyperplane. As a matter of fact support vectors constitute a very useful output of this

method. Apart from the W star identifying which of the example that are closest to the separating class separating boundary is also a useful by product of the SVM method. Now, because these x i's are called the support vectors and the optimal hyperplane is can be written as a linear combination of the support vectors. The method is called a support vector machine, okay?

(Refer Slide Time: 53:49)



That is just take a look at it now. So, this is if this is my optimal hyperplane, right? So, these the ones which are dark here right these two and these two these two these two of class two and these two of class one patterns; they happened to be the support vectors they are the ones closest to the separating hyperplane, okay?

(Refer Slide Time: 54:12)

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 μ_i^* , $i = 1, \dots, n$.

NPTEL

PR NPTEL course - p.105/128

So, let us sum it up the optimal hyperplane W^* b^* is given by W^* is summation over i $\mu_i^* y_i X_i$, which is same as i belongs to S $\mu_i^* y_i X_i$ where S is the set of all indices i for which the corresponding Lagrange multipliers such strictly positive. If I know all the Lagrange multipliers, I can also find b^* because the complementary slackness says for every j y_j into $X_j^T W^* + b^*$ into μ_j^* is 0. So, if there is any j 's has that μ_j^* is greater than 0, then y_j into $X_j^T W^* + b^*$ should be equal to 1.

If I multiply both sides by μ_j^* and take this term on that side I get b^* is equal to y_j minus $X_j^T W^*$ for any j 's, such that μ_j^* is greater than 0. So, every j such that μ_j^* is greater than 0, I can do this. I will get the same value that is what optimization theory tells me. So, if I know all the Lagrange multipliers μ 's, then I can calculate both W^* and b^* . Thus W^* and b^* are determined by μ_i^* $i = 1$ to n .

If I know all the Lagrange multipliers, then I know the support vectors W^* is a specific linear combination of support vectors, where the weights in the linear combination depend on Lagrange multipliers. b^* can also be given in terms of the W^* . So, we can obtain the, we can now solve the dual of the optimization problem to get μ stars, okay?

(Refer Slide Time: 55:36)

Dual optimization problem for SVM

- The dual function is

$$q(\boldsymbol{\mu}) = \inf_{W,b} \left\{ \frac{1}{2}W^TW + \sum_{i=1}^n \mu_i[1 - y_i(W^TX_i + b)] \right\}$$

- If $\sum \mu_i y_i \neq 0$ then $q(\boldsymbol{\mu}) = -\infty$.
- Hence we need to maximize q only over those $\boldsymbol{\mu}$ s.t. $\sum \mu_i y_i = 0$.
- Infimum w.r.t. W is attained at $W = \sum \mu_i y_i X_i$.
- We obtain the dual by substituting $W = \sum \mu_i y_i X_i$ and imposing $\sum \mu_i y_i = 0$.

NPTEL PRE NPTEL course - p.111128

So, for this problem, what is my dual function $q(\boldsymbol{\mu})$ is infimum over W, b of my Lagrangian is half $W^T W$ plus this all the constraint μ_i into this. Now, if I look into this, I have a $\mu_i y_i b$ right summation i is equal to 1 to n $\mu_i y_i$ whole into b is 1 term minus. Now, I am taking infimum both with respect to W and b , so as you already know it is like minimizing $3x$. It is some b into some summation i is equal to 1 to n $\mu_i y_i$, so if $\mu_i y_i$ summation $\mu_i y_i$ is not equal to 0.

Then $q(\boldsymbol{\mu})$ is equal to minus infinity because infimum will be only minus infinity, which means ultimately because I want to maximize q . There is there is there is no I need to only look at those μ for which $\mu_i y_i$ is equal to 0. Now, how do I get infinity with respect to W ? I want to minimize this with respect to W , so take the gradient with respect to W and equate to 0. That we know already know the solution infimum with respect to W is attained by taking W is equal to $\mu_i y_i X_i$, right?

(Refer Slide Time: 56:58)

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

- 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 PR NPTEL course - p.119128

So, essentially I get my dual function by substituting W is equal to $\mu_i y_i x_i$ and imposing $\mu_i y_i$ is equal to 0. So, if we did this we will get the dual of the SVM optimization problem. We will we will look at the details once again next class, but right now I will show you the final expression. This turns out to be the dual of the SVM optimization problem, right? Once again in terms of μ this is a linear terms, this is a quadratic term. So, it is once again a quadratic cost function and the constraints there are so many inequality constraints and one equality constraints.

But all constraints are still linear, so it is the dual is once again a quadratic is a constrained optimization problem with quadratic cost function linear constraints. Another very important thing for us we will, I will once again stress it in the next class is that the training examples come only as inner products here. $x_i^T x_j$ that is the only way the training examples enter into this optimization problem, which which happens to be very useful to us later on. We will remember it because the dual problem the maximization with μ μ is the vector of Lagrange multipliers.

There are as many Lagrange multipliers as there are constraints in the primal problem. There is one constraint per example in the primal problem, so the optimization is over \mathbb{R}^n where n is a number of examples irrespective of the dimension of x . No matter what is the dimension of x is this optimization problem is only over x . So, this is what we will remember. So, we have introduced the optimization problem for the SVM. We looked at

the primal problem and we will we will look at the dual problem in the next class. Then see how to generalize it to a linearly non separable case.

Thank you.