

Introduction to Machine Learning
Prof. Sudeshna Sarkar
Department of Computer Science and Engineering
Indian Institute of Technology, Kharagpur

Module - 5
Lecture - 21
Introduction Support Vector Machine

Good morning, so we start the second part of the lecture of this week. Today, we will give an introduction to support vector machine. We have studied several classifiers, and they have different properties. Support vector machine is one of the most effective classifiers among those, which has sort of linear. It has a very you know good mathematical intuition behind the support vector machine and we are able to handle certain cases where there is non-linearity by using non-linear basis functions or in particular we will see these are called kernel functions.

Now why support vector machine is so popular, we will see that support vector machines have a clever way to prevent over fitting. And we can work with relatively larger number of features without requiring too much computation.

(Refer Slide Time: 01:31)



So, in order to motivate support vector machine, let us look back at logistic regression, which I talked about in the last class. In logistic regression, we have try to find out the probability of y equal to 1 given the input x . And we have model this $h(x)$ as the logistic function σ or g of $\beta^T x$. Now we predict one when this $h(x)$ is greater than 0.5, so we predict one when $h(x)$ is greater than 0.5 and 0 when $h(x)$ is less than 0.5. And the probability value is higher probability of confidence in the output is higher when $h(x)$ is much larger than 0.5, because we are using this sigmoid function which is s shaped you know the function has a shape like this, so at this point $h(x)$ is 0.5.

But as we go towards the side $h(x)$ has value closer to one our probability or confidence in the output would be higher. So, the larger the value of $h(x)$ or this smaller the value of $h(x)$ higher is our confidence in the output. So, we have more confident on the predictions which have further from the decision surface. So, near the decision surface, so logistic regression actually gives you linear classifier.

Let us say just look at, suppose these are your features x_1, x_2 . And you have different point so suppose these are you are positive points, these are you are positive points, and let us say these are you are negative points. Now suppose this is you are decision surface now if this is you are decision surface, then you have less confidence if these is a point here, you have less confidence in you are output of a point, which is closer to the decision surface. You have a higher confidence for those points, which you are further from the decision surface. If you have a large number of features, you could use Bayesian learning.

In Bayesian learning, we will look at all possible classifiers, wait them by their a posteriori probability and combined them to get the answer. And as we have discussed Bayesian learning is computationally intractable. We want to a good alternative to Bayesian learning, which is computationally attractive. So, this brings us to support vector machines.

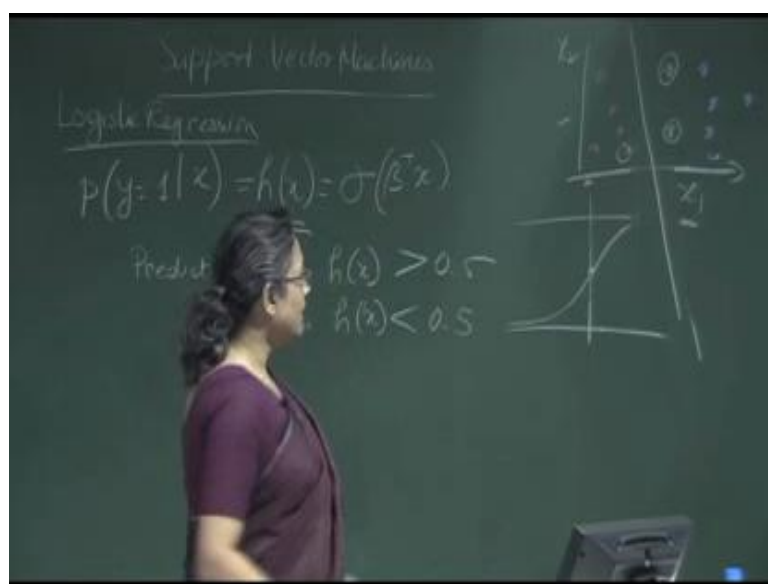
In support vector machine, we want a classifier, so that it maximizes the separation of between the points and the decision surface. So, for that we can define, what we mean by the margin so or you know if you look at this point, we can find the distance of this point

from the decision surface. We can find the distance of this point from the decision surface, this point, this point etcetera for each point; we can find the distance from the decision surface.

Now, if you look at the negative points, this is the closest negative point to the decision surface. In fact, this is the closest point over all to the decision surface and this is the distance of the closest point to the decision surface can be defined as the margin. And these two points these two positive points are closes from the decision surface.

Now we could have shifted this line here so that the distance from the closes positive points and the closes negative point is the same. In that case the margin width will be slightly large, so the width of this band will be twice this. But there are other decision surfaces for which there can be other margins, and among all possible decision surfaces we want to choose that one for which the margin width is highest. So my margin is denoted by the minimum distance of a training instance from the decision surface, and we want to choose that decision surface for which the margin width is highest. All this is under the assumption that the positive and the negative points are linearly separable by a linear decision surface.

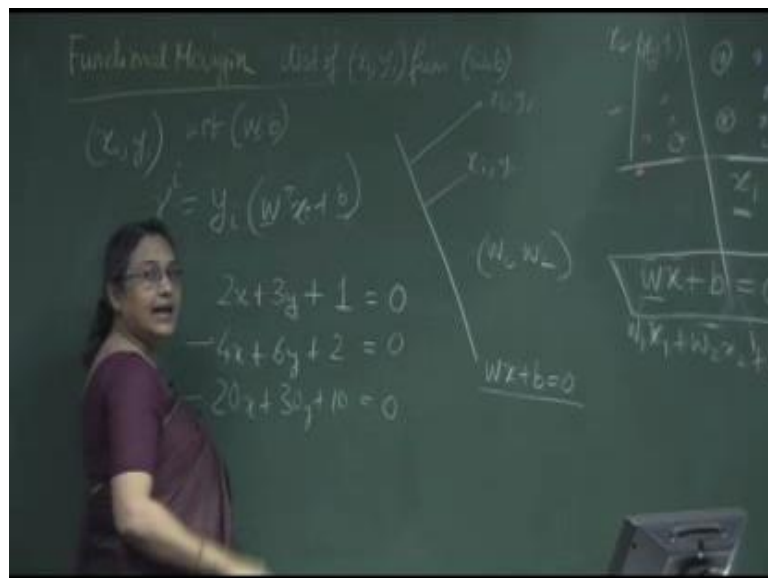
(Refer Slide Time: 07:14)



Now, you know if you consider a particular decision surface, and then you considered the points, which are close to the decision surface. Now these three points are close to the decision surface and then it creates distance from the decision surface. These points are called the support vectors. So, in the case where you are maximizing the margin width, the distance of the closest negative point to this line and the closest positive point to this line will be the same.

And there are a minimum of two support vectors, but there can be more than two support vectors, but typically the number of support vectors is extremely small. And if you notice these support vectors are the ones which actually determine the equation of the line, and typically the support vector number will be very small. And this particular decision surface decides which point given a test point, you can do the test which side of the decision surface it is based on that you can classify the points as positive or negative. Now, with this introduction let us define what we mean by margin formally.

(Refer Slide Time: 08:49)



So, first we will talk about what we call functional margin. Now functional margin of a point you take a point x_i, y_i is an arbitrary point; suppose this is a point x_i, y_i . So, this point with respect to a decision surface, now suppose the equation of this decision surface is given by $w \cdot x + b = 0$. So W is a vector, x is a vector.

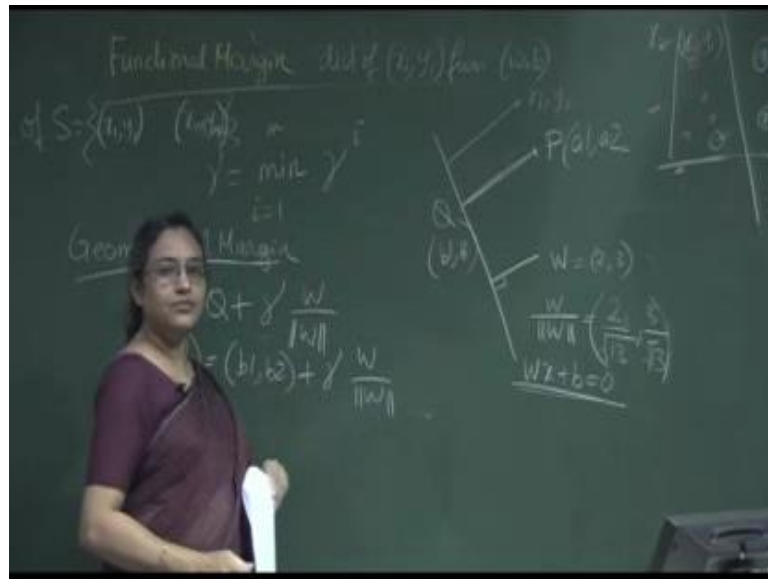
Suppose x_1 and x_2 are the two features, so this is $W_1 x_1 + W_2 x_2 + b = 0$, this is the equation of the line, and it can be written as $W \cdot x + b = 0$. So, with respect to this line $W \cdot x + b = 0$ which is parameterized by W, b the functional margin is measured by the distance of the point x_i, y_i from the decision boundary. So, if the distance of x_i, y_i from the decision boundary, it is given by $W \cdot b$.

So, now let us see how to find this distance. Suppose, this is a straight line, whose equation is given by $W \cdot x + b = 0$. Now the vector which is normal to this line has the normal to this line has, so this vector W_1, W_2 is normal to this straight-line. So, we can define the distance of x_i, y_i from W, b as let us say γ_i and $\gamma_i = y_i (W_1 x_i + W_2 y_i + b)$. So, γ_i will be given by $y_i (W_1 x_i + W_2 y_i + b)$ which is the normal distance from this point x_i, y_i to the decision surface, so this is given by this.

Now, this functional margin you know among another if I take another point, let us say x_j, y_j , it will have another function margin. So, the functional margin of this is more than the functional margin of this. So, if the point is here, further away from the decision surface, we have higher confidence in the classification of this point. So, larger functional margin means more confidence in predicting the class of that point.

The problem of defining the margin as this distance is as follows. This W, b can be arbitrarily scaled. Suppose, the equation of this line is $2x + 3y + 1 = 0$, so w is $2, 3$, b is 1 . Now I can write the same equation as $4x + 6y + 2 = 0$; and for that matter as $20x + 30y + 10 = 0$. So, if we scale this equation the equation of the line remain same, but the functional margin becomes larger. So, the functional margin does not depend on the actual line, it also depends on the coefficients that we are using. Therefore, we need to use some normalization, so that we can look at normalized distance, we will come to that presently. But before that, we have seen what is the functional margin of a point at one point from the line.

(Refer Slide Time: 13:37)



But let us define the functional margin of set of points so we have a set of points $x_1, y_1, x_2, y_2, \dots, x_m, y_m$. So, we have a set of points $x_1, y_1, x_2, y_2, \dots, x_m, y_m$ or rather let us write x, y these are my training points. And with respect to them, the functional margin is defined to be $\gamma = \min_{i=1}^m \gamma_i$. So, among all the functional margins, I want to choose the smallest functional margin that is a smallest functional margin of the different points is the find to be the functional margin of the set of points.

Now, let us come to as I said definition of a functional margin suffers from these drawbacks. So, we want to use form which does not have this drawback and we come to the concept of geometrical margin, which is invariant to the scaling of the equation. Now again let us consider the decision surface given by $Wx + b = 0$. And let us look at a normal to this decision surface. So, the normal to this decision surface is given by W - the vector w is normal to this decision surface.

And the unit vector normal to the decision surface, so this is the vector w and the unit vector normal to the decision surface is given by W divided by $\|W\|$. So, if the equation of the line is $2x + 3y + 1 = 0$, so $W = (2, 3)$. And W by $\|W\|$

is $\frac{2}{\sqrt{2^2 + 3^2}}$ that is $\frac{2}{\sqrt{13}}$, so this is the unit vector, which is normal to the equation normal to this decision surface.

Now, let us look at this point p. So, suppose p has a value a_1, a_2 . And if you draw normal to the decision surface suppose it needs a decision surface at q. And let us say that this has this is a_1, a_2 ; and q has the coordinates b_1, b_2 . Now if you want to find the distance of p from q, the distance is in the direction of the normal vector. So, we can write p equal to q plus gamma times W by W hat, W by W hat is the vector which is normal to this. So, the coordinates of Q plus if gamma is the distance of p from this decision surface the coordinates of P is given by coordinate of Q plus gamma times W by W hat that is we can write other way round a_1, a_2 equal to b_1, b_2 plus gamma times w by w hat. And given w we can write this equation and from that we can find gamma. Let me just rub this; please keep this basic diagram in mind.

(Refer Slide Time: 18:44)



So, from this, we can find out that w transports a_1, a_2 minus gamma w by w hat plus b equal to c , because this point b_1, b_2 . b_1, b_2 , if you just recall, let me just draw the diagram again. So, this is the decision surface this is p - a_1, a_2 ; and this is q given b_1, b_2 . So, q given b_1, b_2 lies on the line $W x$ plus b equal to 0, so this is b_1, b_2 , so it lies on this line, so W into this plus b equal to 0. So, from this, we get W equal to we can solve

for γ and we can find γ equal to $\frac{W \cdot (a_1 a_2 + b)}{\|W\|}$, which is equal to $\frac{W \cdot (a_1 a_2 + b)}{\|W\|}$. So, from this, we find γ equal to $\frac{y}{\|W\|}$ times $W \cdot (a_1 a_2 + b)$ is plus 1 or minus 1.

So, for geometric margin, what we will do is that we will scale W , you know w is the weights of the line $W \cdot b$, we can scale $w \cdot b$ arbitrarily by dividing all of them by some same number or multiplying them by the same number, we will scale W so that $\|w\|$ equal to 1. We will scale w so that $\|w\|$ equal to 1, and then we will find the geometric margin. So, geometric margin, so will scale W , so that $\|w\|$ equal to 1 and then the geometric margin will be given by γ equal to $\frac{y}{\|W\|}$ times $W \cdot (a_1 a_2 + b)$. So, this is the geometric margin, which we get after normalization.

(Refer Slide Time: 22:22)

Geometric Margin

- For a decision surface (w, b)
- the vector orthogonal to it is given by w .
- The unit length orthogonal vector is $\frac{w}{\|w\|}$
- $P = Q + \gamma \frac{w}{\|w\|}$

And as in the previous case if you have a set of points $x_1, y_1, x_2, y_2, \dots, x_m, y_m$, we can find out the geometric margin as the one which is smallest.

(Refer Slide Time: 22:25)

Geometric Margin

$$P = Q + \gamma \frac{w}{\|w\|}$$

$$(b1, b2) = (a1, a2) - \gamma \frac{w}{\|w\|}$$

$$\rightarrow w^T \left((a1, a2) - \gamma \frac{w}{\|w\|} \right) + b = 0$$

$$\rightarrow \gamma = \frac{w^T (a1, a2) + b}{\|w\|}$$

$$= \frac{w^T}{\|w\|} (a1, a2) + \frac{b}{\|w\|}$$

$$= \frac{w^T}{\|w\|} (a1, a2) + \frac{b}{\|w\|}$$

$$\gamma = \gamma \left(\frac{w^T}{\|w\|} (a1, a2) + \frac{b}{\|w\|} \right)$$

Geometric margin : $\|w\| = 1$
 Geometric margin of (w,b) wrt $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$
 -- smallest of the geometric margins of individual points.

(Refer Slide Time: 22:31)

Functional Margin: set of (x, y) from (w,b)

$P(a1, a2)$

$Q(b1, b2)$

w

$\gamma = \min \gamma$

Geometric Margin

$$w^T \left((a1, a2) - \gamma \frac{w}{\|w\|} \right) + b = 0$$

$$\Rightarrow \gamma = \frac{w^T (a1, a2) + b}{\|w\|}$$

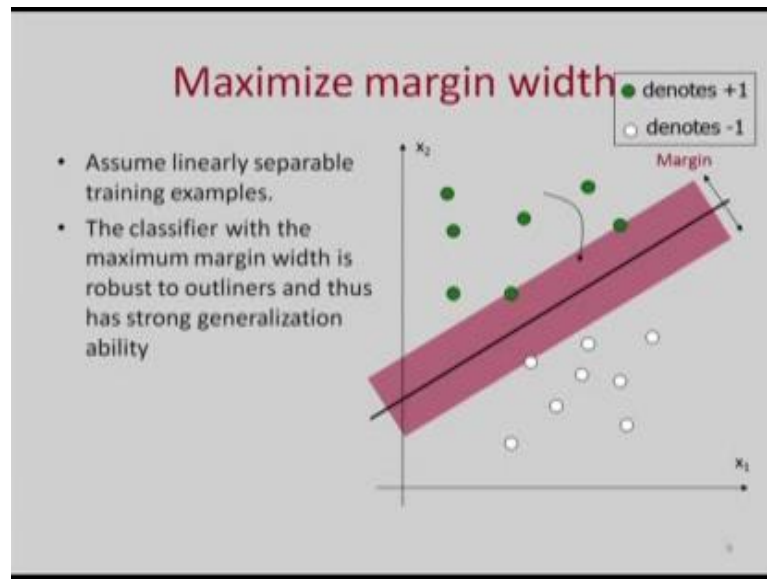
$$= \frac{w^T}{\|w\|} (a1, a2) + \frac{b}{\|w\|}$$

$$\gamma = \gamma \left(\frac{w^T}{\|w\|} (a1, a2) + \frac{b}{\|w\|} \right)$$

$\|w\|=1 \Rightarrow \gamma = \gamma (w^T (a1, a2) + b)$

So, geometric margin will be as before the minimum of γ that will be the geometric margin of set of points.

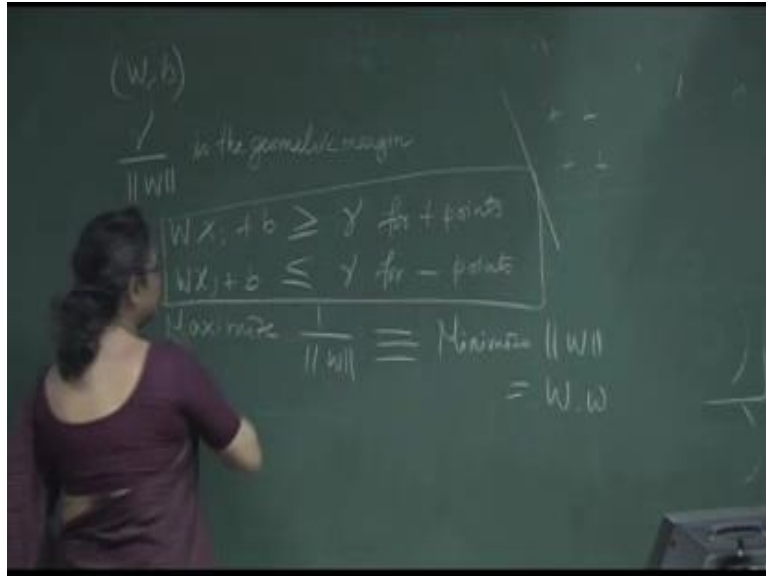
(Refer Slide Time: 22:51)



Now, we will look at how what we really want to do. If we look at this diagram, we assume that the training examples are linearly separable, and let us say this line is our decision surface, and this red band is the margin. And the green points here these two green points lie on the margin this white point also lie on the margin, these are the support vectors. So, classify with the maximum margin width is what we want it is robust to outliers and it has strong generalization ability.

Now once we have defined this we have seen that γ is our geometric margin, and we want to maximize this margin. You know if we without normalization, we get if you have a γ by norm of W is the geometric margin. And we need to maximize γ by W subject to constrain. So, we pose the our optimization problem as follows given a set of training examples labeled as positive and negative.

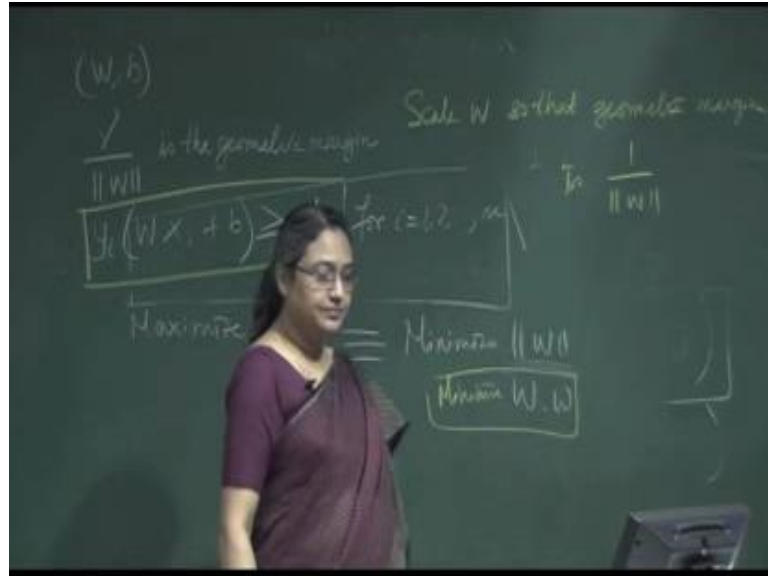
(Refer Slide Time: 24:18)



If $W \cdot b$ characterizes the decision surface then γ by W is the geometric margin. And we want to learn the values of W, b , so that this geometric margin is largest subject to constrain. What are the constrains, we have the positive points on this side, for each positive point, $W \cdot x$ plus b will be greater than equal to γ so all positive points will be at a distance of greater than equal to γ from the margin. So, $W \cdot x_i$ plus b will be greater than equal to γ for plus points; and $W \cdot x_j$ plus b will be less than equal to $-\gamma$ for minus points, so and γ by W is what we want to maximize.

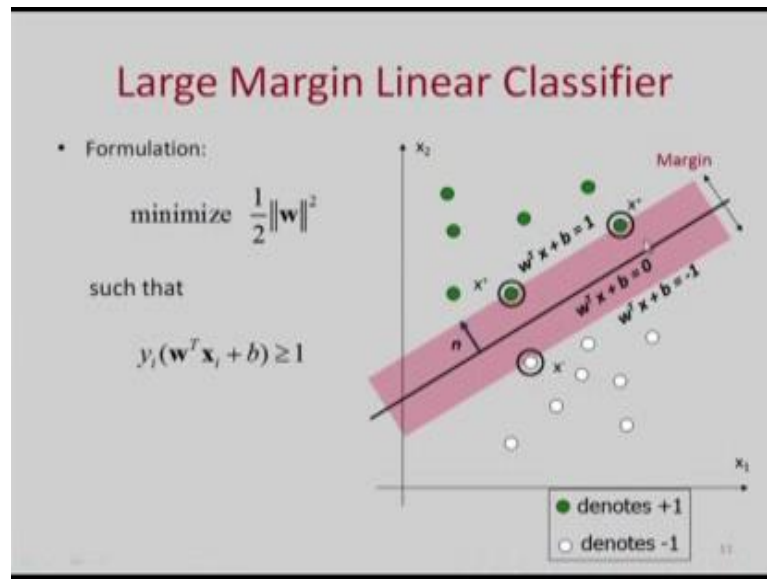
Now, if we so we can say we can maximize so we want to scales, so that so that this γ is 1. We can scale so that this γ is 1, and we can say maximize 1 by W , which is the same as equivalent to minimize W , and W is $w \cdot w$ or $w^T \cdot w$, so this is the $w \cdot w$ dot product. So, we want to minimize $W \cdot W$, subject to this constrains. And now these two constrains can be combined and written in the same form as y_i times $W \cdot x_i$ plus b greater than equal to γ for all points.

(Refer Slide Time: 26:45)



And we will normalize the values of W , so that actually I was wrong we do not normalize w to be 1, we keep w we normalize so that γ equal to 1. So, the geometric margin is γ by W , we will change W , so that so we will scale W for normalization this is what we will do, please pay careful attention. We will scale W , so that geometric margin is 1 by W , so we scale it so that geometric margin is 1 by W . And we want to maximize this margin subject to this constraints $y_i, Wx_i + b$ greater than equal to γ which will become greater than equal to 1, for all training instances right and equivalently we can say minimize W or minimize $W \cdot W$, so minimize $W \cdot W$.

(Refer Slide Time: 28:35)



So, if you look at the slide now so this is the formulation of our optimization problem we want to minimize $\mathbf{W} \cdot \mathbf{W}$. Or we can say minimize \mathbf{W}^2 that is \mathbf{W}_1^2 plus \mathbf{W}_2^2 minimize \mathbf{W}^2 , so that these constrains $y_i \text{ times } \mathbf{W} \cdot \mathbf{x}_i + b$ is greater than equal to 1. What is translates to is that suppose this is you are decision surface which has the equation $\mathbf{W} \cdot \mathbf{x} + b$ equal to 0. And the line parallel to it and this side is $\mathbf{W} \cdot \mathbf{x} + b$ equal to minus 1, the line parallel to this, and this side is $\mathbf{W} \cdot \mathbf{x} + b$ equal to 1. And the margin has been scaled, so that the geometric margin has width 1, and the positive points on the margin have the satisfy the equation $\mathbf{W} \cdot \mathbf{x} + b$ equal to 1, the negative points on the margin satisfy the equation $\mathbf{W} \cdot \mathbf{x} + b$ equal to minus 1.

(Refer Slide Time: 29:40)

Solving the Optimization Problem

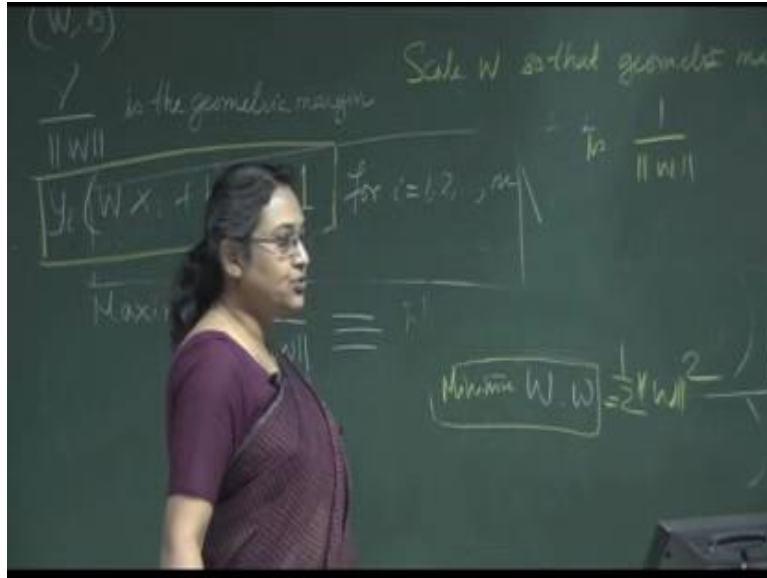
$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{s.t.} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \end{aligned}$$

- Optimization problem with convex quadratic objectives and linear constraints
- Can be solved using QP.
- Lagrange duality to get the optimization problem's dual form,
 - Allow us to use kernels to get optimal margin classifiers to work efficiently in very high dimensional spaces.
 - Allow us to derive an efficient algorithm for solving the above optimization problem that will typically do much better than generic QP software.

17

So, how to solve this optimization problem this is a quadratic programming problem. We have a quadratic objective function, and we have convex quadratic objectives. We have convex quadratic objective function, and we have a set of m constrains m inequality constrains which are linear. So, this optimization problem can be solved using a commercial quadratic programming solver. So, by solving this, we get our classifier and that should be the end of the story. However, there are few other things that we can do by this very powerful formulation and this is what we will discuss now.

(Refer Slide Time: 30:48)



Now, we can straight away solve this problem using, so we are we are minimizing W dot W or W square. And we are we can also say minimize half W square this half is a so that we have a nice form, when we get the solution we could have solve this problem straight away, but we are going to covert this problem into dual formulation. And we will see that the dual formulation have certain properties.

Now, if you want to solve this optimization problem the method to solve this is by using Lagrange multiplier, you can use Lagrangian, and you can get a formulation of this problem, which you can solve and you can get the Lagrange duality to get the dual of this optimization problem. And what we gain by using this is that we can you know this formulation has some nice properties which we will are we able to see later.

With this, we come to the end of this introductory part of the support vector machine. In the next class, we will talk about the dual formulation.

Thank you.