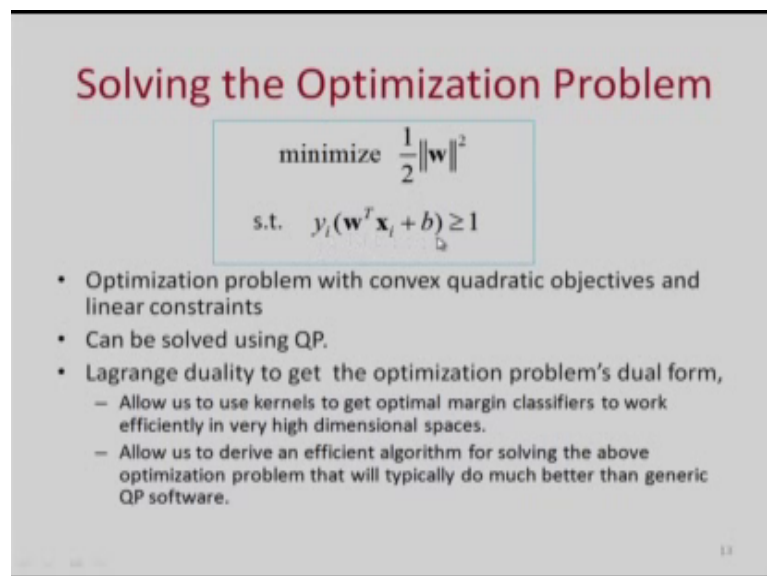


Introduction to Machine Learning
Prof. Sudeshna Sarkar
Department of Computer Science and Engineering
Indian Institute of Technology, Kharagpur

Module - 5
Lecture - 22
SVM: The Dual Formulation

Good morning. We will now talk about Part C of this lecture, where we will look at The Dual Formulation of support vector machine.

(Refer Slide Time: 00:28)



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.

13

In the last class, we looked at the formulation of the optimization problem corresponding to support vector machine where we have to minimize half w square, this is the convex quadratic optimization function subject to this linear constraints $y_i \mathbf{w}^T \mathbf{x}_i + b$ greater than equal to 1 for all examples.

(Refer Slide Time: 00:55)

Lagrangian Duality in brief

The Primal Problem

$$\begin{array}{ll} \min_w & f(w) \\ \text{s.t.} & g_i(w) \leq 0, \quad i = 1, \dots, k \\ & h_i(w) = 0, \quad i = 1, \dots, l \end{array}$$

The generalized Lagrangian:

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

the α 's ($\alpha_i \geq 0$) and β 's are called the Lagrange multipliers

Lemma:

$$\max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{otherwise} \end{cases}$$

A re-written Primal:

$$\min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta)$$

So, before we look at how to get the dual of this particular formulation, let us very briefly talk about Lagrangian Duality. Suppose, we take a general primal general problem and its primal formulation is given by, this is an optimization problem where you want to minimize $f(w)$, w are the parameters. You want to find values of w , so that is to minimize $f(w)$ and you have a set of linear constraints.

There are two type of linear constraints, equality constraints - we have l equality constraints $h_i(w)$ equal to 0 and k inequality constraints $g_i(w)$ less than equal to 0, all these constraints are linear. Corresponding to this problem the generalized lagrangian is given by a function of w, α, β , $f(w)$ plus summation over $i=1$ to k for all the number of non-linear constraints $i=1$ to k $\alpha_i g_i(w)$ plus summation $i=1$ to l $\beta_i h_i(w)$, where $h_i(w)$ are the equality constraints, $g_i(w)$ are the inequality constraints. The alphas and betas are called Lagrange multipliers and the value of alpha is greater than equal to 0.

So, this is the lagrangian of this optimization function. Now, what we want to do is that, we want to; so this particular lagrangian if you take the values of w, α and β such that if the primal constraints are not satisfied then the value of this lagrangian will be infinity. If the constraints are not satisfied the value of the lagrangian is infinity and it

will be equal to f^* if the constraints are satisfied right. So, we want to find out the values of α, β for which L ; so if you look at the maximum of the value of L w α, β then if w satisfies the primal constraints, it will be equal to f^* and it will be infinity otherwise.

And we can rewrite the primal, as finding the value of this expression $\max_{\alpha, \beta} L(w, \alpha, \beta)$, we want to find out those. So, we first do maximize keeping w fixed for a particular w , we can maximize over α, β and find the expression $\max_{\alpha, \beta} L(w, \alpha, \beta)$ and $\min_w \max_{\alpha, \beta} L(w, \alpha, \beta)$ is either f^* or infinity and if we find the minimum of this, it will be giving us the solution of this. Because, otherwise it has a value of infinity, so the primal can be rewritten as, you take the lagrangian and you find out α, β for a fixed w you can take maximization over α, β and then you can do minimum over w . So, minimum over w maximum over α, β , $L(w, \alpha, \beta)$ is a rewriting of the primal formulation.

(Refer Slide Time: 04:42)

Lagrangian Duality, cont.

The Primal Problem: $p^* = \min_w \max_{\alpha, \beta, a_i \geq 0} L(w, \alpha, \beta)$

The Dual Problem: $d^* = \max_{\alpha, \beta, a_i \geq 0} \min_w L(w, \alpha, \beta)$

Theorem (weak duality):
 $d^* = \max_{\alpha, \beta, a_i \geq 0} \min_w L(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta, a_i \geq 0} L(w, \alpha, \beta) = p^*$

Theorem (strong duality):
 If there exist a saddle point of $L(w, \alpha, \beta)$, we have $d^* = p^*$

So, this is the primal problem and it has a solution let us say the solution is p^* . p^* is the solution of the primal problem - minimum over w maximum over α, β of the lagrangian and the dual problem is we are just putting the minimum here and maximum here. So, we take first, we take max of our α, β minimum over w , $L(w, \alpha, \beta)$

this is the dual formulation and it has the solution d^* . So, this is the primal problem solution, this is the dual problem solution and we have two theorems. This theorem says that first of all if you change the order of max min and min max, it is general expression that $\max \min$ of this expression, any expression is less than equal to $\min \max$ of this expression, right.

And d^* is the max of min of this expression and therefore, d^* is less than equal to p^* . So, d^* is always less than equal to p^* . Now, if there exist a saddle point of this lagrangian where they are equal that is called the saddle point and that is the optimum value of both. So, the optimum value of the primal formulation and the optimum value of the dual formulation will be identical when there is a saddle point.

(Refer Slide Time: 06:32)

The KKT conditions

If there exists some saddle point of L , then it satisfies the following "Karush-Kuhn-Tucker" (KKT) conditions:

$$\frac{\partial}{\partial w_i} L(w, \alpha, \beta) = 0, \quad i = 1, \dots, k$$

$$\frac{\partial}{\partial \beta_i} L(w, \alpha, \beta) = 0, \quad i = 1, \dots, l$$

$$\alpha_i g_i(w) = 0, \quad i = 1, \dots, m$$

$$g_i(w) \leq 0, \quad i = 1, \dots, m$$

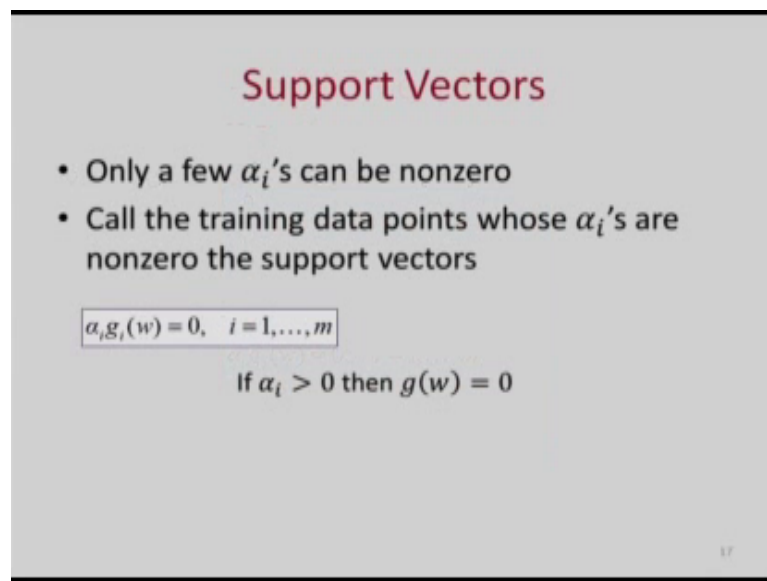
$$\alpha_i \geq 0, \quad i = 1, \dots, m$$

Theorem: If w^* , a^* and b^* satisfy the KKT condition, then it is also a solution to the primal and the dual problems.

And if the saddle point exists, then the saddle point satisfies the following condition called KKT condition or Karush-Kuhn-Tucker condition. Now, the condition says that the partial derivative of this lagrangian, with respect to w_i and with respect to β_i will be equal to 0, according to the KKT condition. And from these two, you will find out that what you get is that $\alpha_i g_i(w)$ will be equal to 0, for i equal to 1 to m , $g_i(w)$ is less than equal to 0 and α_i greater than equal to 0. So, these are the conditions that you get when the saddle point exists.

And the theorem says if w^* , α^* , and β^* satisfy the KKT condition then it is also a solution to the primal and dual problems. With this brief description, brief outline of Lagrangian duality let us go back to SVM and see how it can be applied there. The details of this theory are beyond the scope of this class and you can read some material on convex optimization if you want to learn more about this.

(Refer Slide Time: 07:59)



Support Vectors

- Only a few α_i 's can be nonzero
- Call the training data points whose α_i 's are nonzero the support vectors

$\alpha_i g_i(w) = 0, \quad i = 1, \dots, m$

If $\alpha_i > 0$ then $g_i(w) = 0$

Now, if we look at our SVM formulation what we have is we have $f(w)$ as half w square and we have the $g(w)$ as $y_i w^T x_i + b \geq 1$. We do not have the h , we do not have the equality constraints, we have only the objective function and the g_i constraints. So, we are only dealing with α_i , not the β_i . So, we are dealing with $f(w)$ plus $\alpha_i g_i(w)$.

So, this KKT conditions also says that this $\alpha_i g_i(w)$ equal to 0 and $g_i(w)$ is less than 0, it is because $\alpha_i g_i(w)$ is 0 only when α_i is 0 then $g_i(w)$ can be non-zero, and otherwise, $g_i(w)$ is . And which says that only the few of the α_i 's can be non-zero and the training data points whose α_i 's are non-zero are called the support vectors. So, some of the α_i are non-zero and the training data corresponding to data support α_i is greater than equal to 0, if α_i greater than 0, $g_i(w)$ will be equal to 0.

(Refer Slide Time: 09:45)

Solving the Optimization Problem

Quadratic programming with linear constraints

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

Lagrangian Function

$$\begin{aligned} &\text{minimize } L_p(\mathbf{w}, b, \alpha_i) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1) \\ &\text{s.t. } \alpha_i \geq 0 \end{aligned}$$

18

Now, let us see the implication. So, this is the original optimization problem and when we take, this is the SVM optimization problem we take the lagrangian which gives us minimization of $\frac{1}{2} \|\mathbf{w}\|^2$ b alpha, we have written L_p - p denotes primal. So, minimize $L_p(\mathbf{w}, b, \alpha_i)$, half \mathbf{w} square minus sigma alpha i $y_i \mathbf{w}^T \mathbf{x}_i + b - 1$, subject to the constraints alpha i greater than equal to 0. This is by getting the lagrangian of the optimization problem.

(Refer Slide Time: 10:25)

Solving the Optimization Problem

$$\begin{aligned} \text{minimize } L_p(\mathbf{w}, b, \alpha_i) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \\ \text{s.t. } \alpha_i &\geq 0 \end{aligned}$$

Minimize wrt \mathbf{w} and b for fixed α

$$\frac{\partial L_p}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L_p}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$L_p(\mathbf{w}, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) - b \sum_{i=1}^m \alpha_i y_i$$

$$L_p(\mathbf{w}, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

Now, if we take the partial derivative of this L_p with respect to \mathbf{w} and b what we get here is, \mathbf{w} equal to $\sum \alpha_i y_i \mathbf{x}_i$, and from the second one by taking partial derivative of L_p with respect to b and setting it to 0 we get $\sum \alpha_i y_i$ equal to 0.

So, what it means is that if I substitute this value of \mathbf{w} in this expression here, if I substitute this value of $\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i$ in this expression here what I get is L_p with b and α equal to $\sum \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) - b \sum \alpha_i y_i$. I am sorry. So, we put \mathbf{w} is here, so this \mathbf{w} becomes half of $\alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ minus $b \sum \alpha_i y_i$. So, this is my L_p when I substitute this value of \mathbf{w} . But we know that $\sum \alpha_i y_i$ equal to 0 from this constraint. So, this expression on the right side can be ignored and finally, we get L_p with b and α as this expression.

So, L_p with b and α is $\sum \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$. Now, this very important formulation and we will look at the properties of this formulation to get certain properties of the support vector machine algorithm.

(Refer Slide Time: 12:20)

The Dual problem

Now we have the following dual opt problem:

$$\begin{array}{ll} \max_{\alpha} J(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \text{s.t.} & \alpha_i \geq 0, \quad i=1, \dots, k \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{array}$$

This is a **quadratic programming** problem.
— A global maximum of α_i can always be found.

So, the Dual problem, before we go to that let us look at the dual problem the dual problem is maximizing of $J(\alpha)$ where $J(\alpha)$ is the expression we saw earlier and these are the constraints $\alpha_i \geq 0$ and $\sum \alpha_i y_i = 0$. This is the dual problem, which is a quadratic programming problem and from this quadratic programming problem we can solve and find the global maximum value of α_i . We can find out the values of α_i , by solving this quadratic programming problem.

(Refer Slide Time: 12:59)

Support vector machines

- Once we have the Lagrange multipliers $\{\alpha_i\}$ we can reconstruct the parameter vector w as a weighted combination of the training examples:

$$w = \sum_{i=1}^n \alpha_i y_i x_i \quad w = \sum_{i \in SV} \alpha_i y_i x_i$$

- For testing with a new data z
 - Compute

$$w^T z + b = \sum_{i \in SV} \alpha_i y_i (x_i^T z) + b$$

and classify z as class 1 if the sum is positive, and class 2 otherwise

Note: w need not be formed explicitly

22

And this quadratic programming problem is much easier to solve than the primal formulation. This is much simpler, because the constraints are simpler and we will see it has certain nice properties.

So, once we solve and get the lagrange multipliers α , we can reconstruct the parameter vectors. We can find w as $\sum \alpha_i y_i x_i$. In fact, we noted that α_i is non-zero only for few of the examples. Those examples are the one, once which are the support vectors. So, w is obtained from $\sum \alpha_i y_i x_i$ where i ranges among the support vectors and usually the support vectors are few in number and w can be computed from the coordinates of those support vectors.

Also when we get a new data point z in order to find out the output corresponding to this we can compute $w^T z + b$, which is $\sum \alpha_i y_i x_i^T z + b$ and we classify z as class 1 if the sum is positive and class 2 if otherwise. Now, you note that w need not be found explicitly we can just use this expression and this expression has a very nice property when you put z what you are doing is - this α_i , this is y_i , this is $x_i^T z$. So, you are taking the dot product of the support vector with your test point.

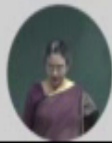
(Refer Slide Time: 14:59)

Solving the Optimization Problem

- The discriminant function is:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{i \in \text{SV}} \alpha_i \mathbf{x}_i^T \mathbf{x} + b$$

- It relies on a *dot product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i
- Solving the optimization problem involved computing the *dot products* $\mathbf{x}_i^T \mathbf{x}_j$ between all pairs of training points
- The optimal \mathbf{w} is a linear combination of a small number of data points.



12

So, the discriminant function is given by this dot product of \mathbf{x}_i^T and \mathbf{x} right. So, the computation reduces to mainly finding these dot products. So, you have the dot product between the test point \mathbf{x} and the support vector \mathbf{x}_i . Why is this, such an exciting thought? Now, \mathbf{x}_i is a vector and this can be a high dimensional vector, if you take the dot product of these two linear vectors what you get is a scalar. So, the dot product is a scalar value.

And we will look at what are the implications later, when we solve the optimization problem also, if we look at this formulation where we solve the optimization problem. In here you see what we have is the dot product of the training points, so $\alpha_i \alpha_j \mathbf{y}_i \mathbf{y}_j$ is either plus 1 minus 1. So, these are very simple to compute multiply and $\mathbf{x}_i^T \mathbf{x}_j$ the dot product of \mathbf{x}_i and \mathbf{x}_j . So, when we solve the optimization problems it involved computing the dot products between all the pairs of training points and the optimal \mathbf{w} is linear combination of a small number of data points. So, these are some of the features about this SVM formulation.

We stop here today, in the next class we will look at certain properties of SVM and how these properties can be used for those formulations of SVM. With this I end today's lecture.

Thank you.