

## CS229 Lecture notes

Andrew Ng

### Part V

# Support Vector Machines

This set of notes presents the Support Vector Machine (SVM) learning algorithm. SVMs are among the best (and many believe is indeed the best) “off-the-shelf” supervised learning algorithm. To tell the SVM story, we’ll need to first talk about margins and the idea of separating data with a large “gap.” Next, we’ll talk about the optimal margin classifier, which will lead us into a digression on Lagrange duality. We’ll also see kernels, which give a way to apply SVMs efficiently in very high dimensional (such as infinite-dimensional) feature spaces, and finally, we’ll close off the story with the SMO algorithm, which gives an efficient implementation of SVMs.

## 1 Margins: Intuition

We’ll start our story on SVMs by talking about margins. This section will give the intuitions about margins and about the “confidence” of our predictions; these ideas will be made formal in Section 3.

Consider logistic regression, where the probability  $p(y = 1|x; \theta)$  is modeled by  $h_{\theta}(x) = g(\theta^T x)$ . We would then predict “1” on an input  $x$  if and only if  $h_{\theta}(x) \geq 0.5$ , or equivalently, if and only if  $\theta^T x \geq 0$ . Consider a positive training example ( $y = 1$ ). The larger  $\theta^T x$  is, the larger also is  $h_{\theta}(x) = p(y = 1|x; \theta)$ , and thus also the higher our degree of “confidence” that the label is 1. Thus, informally we can think of our prediction as being a very confident one that  $y = 1$  if  $\theta^T x \gg 0$ . Similarly, we think of logistic regression as making a very confident prediction of  $y = 0$ , if  $\theta^T x \ll 0$ . Given a training set, again informally it seems that we’d have found a good fit to the training data if we can find  $\theta$  so that  $\theta^T x^{(i)} \gg 0$  whenever  $y^{(i)} = 1$ , and

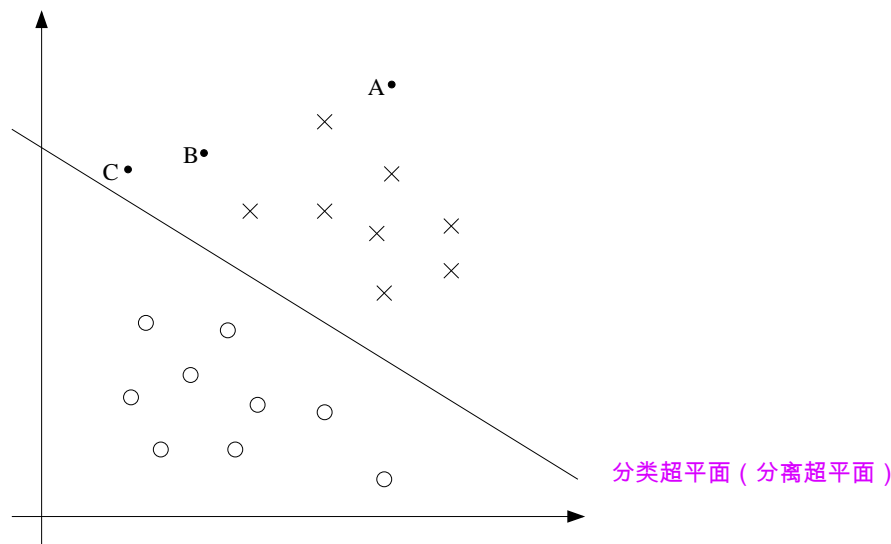
非正式的定义：直观上的解释，当样本为正类的时候，我们希望 $\theta^T x \gg 0$ ，当样本为负类的时候，我们希望 $\theta^T x \ll 0$

如果我们只从出发，希望模型达到的目标就是让训练数据中 $y=1$ 的特征，而是 $y=0$ 的特征。Logistic回归就是要学习得到，使得正例的特征远大于0，负例的特征远小于0，而且要在全部训练实例上达到这个目标。

2

$\theta^T x^{(i)} \ll 0$  whenever  $y^{(i)} = 0$ , since this would reflect a very confident (and correct) set of classifications for all the training examples. This seems to be a nice goal to aim for, and we'll soon formalize this idea using the notion of functional margins.

For a different type of intuition, consider the following figure, in which  $x$ 's represent positive training examples,  $o$ 's denote negative training examples, a decision boundary (this is the line given by the equation  $\theta^T x = 0$ , and is also called the **separating hyperplane**) is also shown, and three points have also been labeled A, B and C.



Notice that the point A is very far from the decision boundary. If we are asked to make a prediction for the value of  $y$  at A, it seems we should be quite confident that  $y = 1$  there. Conversely, the point C is very close to the decision boundary, and while it's on the side of the decision boundary on which we would predict  $y = 1$ , it seems likely that just a small change to the decision boundary could easily have caused our prediction to be  $y = 0$ . Hence, we're much more confident about our prediction at A than at C. The point B lies in-between these two cases, and more broadly, we see that if a point is far from the separating hyperplane, then we may be significantly more confident in our predictions. Again, **informally we think it'd be nice if, given a training set, we manage to find a decision boundary that allows us to make all correct and confident (meaning far from the decision boundary) predictions on the training examples.** We'll formalize this later using the notion of geometric margins.

符号定义

## 2 Notation

label :  $y$  的取值用  $\{-1, 1\}$ , 代替  $\{0, 1\}$ 

To make our discussion of SVMs easier, we'll first need to introduce a new notation for talking about classification. We will be considering a linear classifier for a binary classification problem with labels  $y$  and features  $x$ . From now, we'll use  $y \in \{-1, 1\}$  (instead of  $\{0, 1\}$ ) to denote the class labels. Also, rather than parameterizing our linear classifier with the vector  $\theta$ , we will use parameters  $w, b$ , and write our classifier as

$$h_{w,b}(x) = g(w^T x + b).$$

Here,  $g(z) = 1$  if  $z \geq 0$ , and  $g(z) = -1$  otherwise. This " $w, b$ " notation allows us to explicitly treat the intercept term  $b$  separately from the other parameters. (We also drop the convention we had previously of letting  $x_0 = 1$  be an extra coordinate in the input feature vector.) Thus,  $b$  takes the role of what was previously  $\theta_0$ , and  $w$  takes the role of  $[\theta_1 \dots \theta_n]^T$ .

Note also that, from our definition of  $g$  above, our classifier will directly predict either 1 or  $-1$  (cf. the perceptron algorithm), without first going through the intermediate step of estimating the probability of  $y$  being 1 (which was what logistic regression did).

## 3 Functional and geometric margins

函数间隔和几何间隔

Lets formalize the notions of the functional and geometric margins. Given a training example  $(x^{(i)}, y^{(i)})$ , we define the **functional margin** of  $(w, b)$  with respect to the training example

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x + b).$$

函数间隔的定义

首先函数间隔  $> 0$  的话, 则代表分类正确了

其次, 我们希望在分类正确的前提下最大化函数间隔

Note that if  $y^{(i)} = 1$ , then for the functional margin to be large (i.e., for our prediction to be confident and correct), then we need  $w^T x + b$  to be a large positive number. Conversely, if  $y^{(i)} = -1$ , then for the functional margin to be large, then we need  $w^T x + b$  to be a large negative number. Moreover, if  $y^{(i)}(w^T x + b) > 0$ , then our prediction on this example is correct. (Check this yourself.) Hence, **a large functional margin represents a confident and a correct prediction.**

函数间隔应该大于0并且应该尽量大

For a linear classifier with the choice of  $g$  given above (taking values in  $\{-1, 1\}$ ), there's one property of the functional margin that makes it not a very good measure of confidence, however. Given our choice of  $g$ , we note that if we replace  $w$  with  $2w$  and  $b$  with  $2b$ , then since  $g(w^T x + b) = g(2w^T x + 2b)$ ,

上面的理论存在不足之处, 我们可以通过倍增  $w, b$  达到目的, 倍增  $w, b$ , 函数并没有改变 (分类超平面并没有改变)

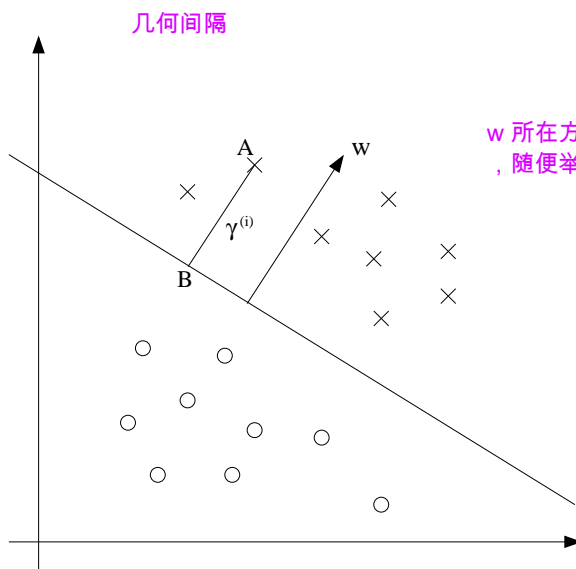
this would not change  $h_{w,b}(x)$  at all. I.e.,  $g$ , and hence also  $h_{w,b}(x)$ , depends only on the sign, but not on the magnitude, of  $w^T x + b$ . However, replacing  $(w, b)$  with  $(2w, 2b)$  also results in multiplying our functional margin by a factor of 2. Thus, it seems that by exploiting our freedom to scale  $w$  and  $b$ , we can make the functional margin arbitrarily large without really changing anything meaningful. Intuitively, it might therefore make sense to impose some sort of normalization condition such as that  $\|w\|_2 = 1$ ; i.e., we might replace  $(w, b)$  with  $(w/\|w\|_2, b/\|w\|_2)$ , and instead consider the functional margin of  $(w/\|w\|_2, b/\|w\|_2)$ . We'll come back to this later.

Given a training set  $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$ , we also define the function margin of  $(w, b)$  with respect to  $S$  as the smallest of the functional margins of the individual training examples. Denoted by  $\hat{\gamma}$ , this can therefore be written:

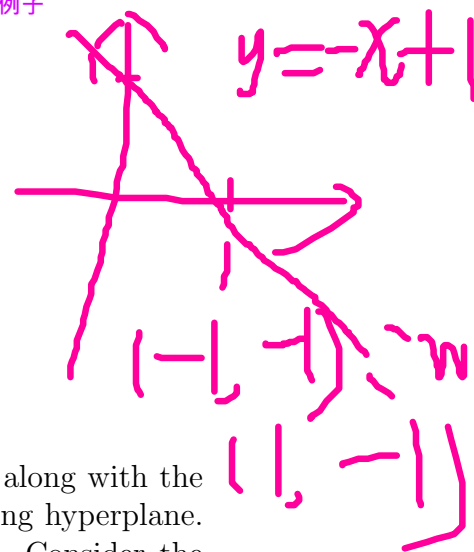
$$\hat{\gamma} = \min_{i=1, \dots, m} \hat{\gamma}^{(i)}.$$

给出训练集上最小函数间隔的定义

Next, let's talk about **geometric margins**. Consider the picture below:



$w$  所在方向与超平面所在方向是正交的，随便举个例子



The decision boundary corresponding to  $(w, b)$  is shown, along with the vector  $w$ . Note that  $w$  is orthogonal (at  $90^\circ$ ) to the separating hyperplane. (You should convince yourself that this must be the case.) Consider the point at A, which represents the input  $x^{(i)}$  of some training example with label  $y^{(i)} = 1$ . Its distance to the decision boundary,  $\gamma^{(i)}$ , is given by the line segment AB.

How can we find the value of  $\gamma^{(i)}$ ? Well,  $w/\|w\|$  is a unit-length vector pointing in the same direction as  $w$ . Since A represents  $x^{(i)}$ , we therefore

find that the point  $B$  is given by  $x^{(i)} - \gamma^{(i)} \cdot w / \|w\|$ . But this point lies on the decision boundary, and all points  $x$  on the decision boundary satisfy the equation  $w^T x + b = 0$ . Hence,

$$w^T \left( x^{(i)} - \gamma^{(i)} \frac{w}{\|w\|} \right) + b = 0.$$

Solving for  $\gamma^{(i)}$  yields

$$\gamma^{(i)} = \frac{w^T x^{(i)} + b}{\|w\|} = \left( \frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|}.$$

分类为正的时候，这个值是正的，分类为负的时候，这个值为负的

This was worked out for the case of a positive training example at  $A$  in the figure, where being on the “positive” side of the decision boundary is good. More generally, we define the geometric margin of  $(w, b)$  with respect to a training example  $(x^{(i)}, y^{(i)})$  to be

$$\gamma^{(i)} = y^{(i)} \left( \left( \frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right).$$

我们乘上 $y(i)$ ，让几何间隔恒为正

Note that if  $\|w\| = 1$ , then the functional margin equals the geometric margin—this thus gives us a way of relating these two different notions of margin. Also, the geometric margin is invariant to rescaling of the parameters; i.e., if we replace  $w$  with  $2w$  and  $b$  with  $2b$ , then the geometric margin does not change. This will in fact come in handy later. Specifically, because of this invariance to the scaling of the parameters, when trying to fit  $w$  and  $b$  to training data, we can impose an arbitrary scaling constraint on  $w$  without changing anything important; for instance, we can demand that  $\|w\| = 1$ , or  $|w_1| = 5$ , or  $|w_1 + b| + |w_2| = 2$ , and any of these can be satisfied simply by rescaling  $w$  and  $b$ .

Finally, given a training set  $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$ , we also define the geometric margin of  $(w, b)$  with respect to  $S$  to be the smallest of the geometric margins on the individual training examples:

$$\gamma = \min_{i=1, \dots, m} \gamma^{(i)}.$$

给出训练集上最小几何间隔的定义

## 4 The optimal margin classifier

最优间隔分类器

Given a training set, it seems from our previous discussion that a natural desideratum is to try to find a decision boundary that maximizes the (geometric) margin, since this would reflect a very confident set of predictions

on the training set and a good “fit” to the training data. Specifically, this will result in a classifier that separates the positive and the negative training examples with a “gap” (geometric margin).

For now, we will assume that we are given a training set that is linearly separable; i.e., that it is possible to separate the positive and negative examples using some separating hyperplane. How we we find the one that achieves the maximum geometric margin? We can pose the following optimization problem:

$$\begin{aligned} \max_{\gamma, w, b} \quad & \gamma \quad \text{几何间隔} \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq \gamma, \quad i = 1, \dots, m \\ & \|w\| = 1. \end{aligned}$$

到目前为止，都是训练集线性可分，能用一个超平面将数据集划分开。  
我们限制 $\|w\|=1$ ，那么 $y^{(i)}(w^T x^{(i)} + b) / \|w\|$ 也就是：  
 $y^{(i)}(w^T x^{(i)} + b)$ ，也就是说当 $\|w\|=1$ 时，函数间隔就是几何间隔

I.e., we want to maximize  $\gamma$ , subject to each training example having functional margin at least  $\gamma$ . The  $\|w\| = 1$  constraint moreover ensures that the functional margin equals to the geometric margin, so we are also guaranteed that all the geometric margins are at least  $\gamma$ . Thus, solving this problem will result in  $(w, b)$  with the largest possible geometric margin with respect to the training set.

If we could solve the optimization problem above, we’d be done. But the “ $\|w\| = 1$ ” constraint is a nasty (non-convex) one, and this problem certainly isn’t in any format that we can plug into standard optimization software to solve. So, let’s try transforming the problem into a nicer one. Consider:

$$\begin{aligned} \max_{\gamma, w, b} \quad & \frac{\hat{\gamma}}{\|w\|} \quad \text{函数间隔}/\|w\|, \text{仍然是几何间隔, 仍然是最大化几何间隔} \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, \quad i = 1, \dots, m \end{aligned}$$

Here, we’re going to maximize  $\hat{\gamma}/\|w\|$ , subject to the functional margins all being at least  $\hat{\gamma}$ . Since the geometric and functional margins are related by  $\gamma = \hat{\gamma}/\|w\|$ , this will give us the answer we want. Moreover, we’ve gotten rid of the constraint  $\|w\| = 1$  that we didn’t like. The downside is that we now have a nasty (again, non-convex) objective  $\frac{\hat{\gamma}}{\|w\|}$  function; and, we still don’t have any off-the-shelf software that can solve this form of an optimization problem.

Let’s keep going. Recall our earlier discussion that we can add an arbitrary scaling constraint on  $w$  and  $b$  without changing anything. This is the key idea we’ll use now. We will introduce the scaling constraint that the functional margin of  $w, b$  with respect to the training set must be 1:

$$\hat{\gamma} = 1.$$

$\|w\| = 1$ 这个约束导致这个最优化问题是一个非凸优化问题

我们利用几何间隔和函数间隔的关系去掉了  $\|w\| = 1$  这个约束条件

但是这仍然是一个非凸优化问题

我们前面知道我们可以任意等比例缩放  $w$  和  $b$ ，分离超平面不会有任何改变，我们就要用到这种思想。

我的理解：函数间隔： $y^{(i)}(w^T x^{(i)} + b)$ ，我们又可以任意的缩放  $w$  和  $b$ ，那我们一定可以将最小函数间隔缩放（改变  $w$  和  $b$  即可）为 1。那我们就假设最小函数间隔为 1。原来的问题是最大化函数间隔  $\|w\|$ ，现在函数间隔为 1 了（总可以调整  $w, b$  实现），那就是最大化  $1/\|w\|$ ，也就是最小化  $1/2 * \|w\|^2$

Since multiplying  $w$  and  $b$  by some constant results in the functional margin being multiplied by that same constant, this is indeed a scaling constraint, and can be satisfied by rescaling  $w, b$ . Plugging this into our problem above, and noting that maximizing  $\hat{\gamma}/\|w\| = 1/\|w\|$  is the same thing as minimizing  $\|w\|^2$ , we now have the following optimization problem:

$$\begin{aligned} \min_{\gamma, w, b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, m \end{aligned}$$

We've now transformed the problem into a form that can be efficiently solved. The above is an optimization problem with a convex quadratic objective and only linear constraints. Its solution gives us the **optimal margin classifier**. This optimization problem can be solved using commercial quadratic programming (QP) code.<sup>1</sup>

While we could call the problem solved here, what we will instead do is make a digression to talk about Lagrange duality. This will lead us to our optimization problem's dual form, which will play a key role in allowing us to use kernels to get optimal margin classifiers to work efficiently in very high dimensional spaces. The dual form will also allow us to derive an efficient algorithm for solving the above optimization problem that will typically do much better than generic QP software.

## 5 Lagrange duality

拉格朗日乘数法  
拉格朗日对偶

原始问题  
对偶问题  
KKT

Lets temporarily put aside SVMs and maximum margin classifiers, and talk about solving constrained optimization problems.

Consider a problem of the following form:

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned} \quad \text{注意约束条件是 } = 0$$

Some of you may recall how the method of Lagrange multipliers can be used to solve it. (Don't worry if you haven't seen it before.) In this method, we define the **Lagrangian** to be

可以用拉格朗日乘数法  
解决这个问题

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

<sup>1</sup>You may be familiar with linear programming, which solves optimization problems that have linear objectives and linear constraints. QP software is also widely available, which allows convex quadratic objectives and linear constraints.

拉格朗日乘数

Here, the  $\beta_i$ 's are called the **Lagrange multipliers**. We would then find and set  $\mathcal{L}$ 's partial derivatives to zero:

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0; \quad \frac{\partial \mathcal{L}}{\partial \beta_i} = 0,$$

解出来以后, 如果有多组解, 那么带入各组  $w$  和  $\beta$ , 再经过比较选出使  $L(w, \beta)$  最小的  $w, \beta$  即可

and solve for  $w$  and  $\beta$ .

In this section, we will generalize this to constrained optimization problems in which we may have inequality as well as equality constraints. Due to time constraints, we won't really be able to do the theory of Lagrange duality justice in this class,<sup>2</sup> but we will give the main ideas and results, which we will then apply to our optimal margin classifier's optimization problem.

Consider the following, which we'll call the **primal** optimization problem:

原始问题

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

注意这里是小于等于 (不等式)

To solve it, we start by defining the **generalized Lagrangian**

广义的拉格朗日乘法

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

Here, the  $\alpha_i$ 's and  $\beta_i$ 's are the Lagrange multipliers. Consider the quantity

拉格朗日乘数

$$\theta_{\mathcal{P}}(w) = \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta).$$

Here, the " $\mathcal{P}$ " subscript stands for "primal." Let some  $w$  be given. If  $w$  violates any of the primal constraints (i.e., if either  $g_i(w) > 0$  or  $h_i(w) \neq 0$  for some  $i$ ), then you should be able to verify that

如果给定的  $w$  不满足约束条件, 那么通过调整  $\alpha$  或者  $\beta$  总能使  $\theta_{\mathcal{P}}(w)$  取到正无穷

$$\theta_{\mathcal{P}}(w) = \max_{\alpha, \beta: \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w) \quad (1)$$

$$= \infty. \quad (2)$$

Conversely, if the constraints are indeed satisfied for a particular value of  $w$ , then  $\theta_{\mathcal{P}}(w) = f(w)$ . Hence,

如果, 给定  $w$  满足约束条件, 那么  $\theta_{\mathcal{P}}(w)$  的最大值一定是  $f(w)$

$$\theta_{\mathcal{P}}(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{otherwise.} \end{cases}$$

综上, 得到这个等式

<sup>2</sup>Readers interested in learning more about this topic are encouraged to read, e.g., R. T. Rockafeller (1970), *Convex Analysis*, Princeton University Press.



在一堆约束条件下，优化  $f(w)$ ，转化为在  $\alpha \geq 0$  的条件下，优化  $\theta_P(w)$ ，并且相当于先给定  $w$ ，然后调整  $\alpha, \beta$ ，使  $L(w, \alpha, \beta)$  取最大值，整体上使最大值的最小值

9

推出上面那个等式的目的是什么呢？说明只要  $w$  满足约束条件，那么  $\theta_P(w)$  和  $f(w)$  是一样的

Thus,  $\theta_P$  takes the same value as the objective in our problem for all values of  $w$  that satisfies the primal constraints, and is positive infinity if the constraints are violated. Hence, if we consider the minimization problem

$$\min_w f(w) = \min_w \theta_P(w) = \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta),$$

最小化  $f(w)$  也就是最小化  $\theta_P(w)$

we see that it is the same problem (i.e., and has the same solutions as) our original, primal problem. For later use, we also define the optimal value of the objective to be  $p^* = \min_w \theta_P(w)$ ; we call this the **value** of the primal problem.

Now, let's look at a slightly different problem. We define

$$\theta_D(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta).$$

Here, the “ $\mathcal{D}$ ” subscript stands for “dual.” Note also that whereas in the definition of  $\theta_P$  we were optimizing (maximizing) with respect to  $\alpha, \beta$ , here we are minimizing with respect to  $w$ .

We can now pose the **dual** optimization problem:

先固定  $\alpha, \beta, \dots$ ，整体是最小值的最大

$$\max_{\alpha, \beta: \alpha_i \geq 0} \theta_D(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta).$$

This is exactly the same as our primal problem shown above, except that the order of the “max” and the “min” are now exchanged. We also define the optimal value of the dual problem's objective to be  $d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \theta_D(w)$ .

How are the primal and the dual problems related? It can easily be shown that

$$d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta) = p^*.$$

(You should convince yourself of this; this follows from the “max min” of a function always being less than or equal to the “min max.”) However, under certain conditions, we will have

$$d^* = p^*,$$

so that we can solve the dual problem in lieu of the primal problem. Let's see what these conditions are.

Suppose  $f$  and the  $g_i$ 's are convex,<sup>3</sup> and the  $h_i$ 's are affine.<sup>4</sup> Suppose further that the constraints  $g_i$  are (strictly) feasible; this means that there exists some  $w$  so that  $g_i(w) < 0$  for all  $i$ .

<sup>3</sup>When  $f$  has a Hessian, then it is convex if and only if the Hessian is positive semi-definite. For instance,  $f(w) = w^T w$  is convex; similarly, all linear (and affine) functions are also convex. (A function  $f$  can also be convex without being differentiable, but we won't need those more general definitions of convexity here.)

<sup>4</sup>I.e., there exists  $a_i, b_i$ , so that  $h_i(w) = a_i^T w + b_i$ . “Affine” means the same thing as linear, except that we also allow the extra intercept term  $b_i$ .

Under our above assumptions, there must exist  $w^*, \alpha^*, \beta^*$  so that  $w^*$  is the solution to the primal problem,  $\alpha^*, \beta^*$  are the solution to the dual problem, and moreover  $p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$ . Moreover,  $w^*, \alpha^*$  and  $\beta^*$  satisfy the **Karush-Kuhn-Tucker (KKT) conditions**, which are as follows:

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, n \quad (3)$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l \quad (4)$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k \quad (5)$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k \quad (6)$$

$$\alpha_i^* \geq 0, \quad i = 1, \dots, k \quad (7)$$

Moreover, if some  $w^*, \alpha^*, \beta^*$  satisfy the KKT conditions, then it is also a solution to the primal and dual problems.

We draw attention to Equation (5), which is called the KKT **dual complementarity** condition. Specifically, it implies that if  $\alpha_i^* > 0$ , then  $g_i(w^*) = 0$ . (I.e., the “ $g_i(w) \leq 0$ ” constraint is **active**, meaning it holds with equality rather than with inequality.) Later on, this will be key for showing that the SVM has only a small number of “support vectors”; the KKT dual complementarity condition will also give us our convergence test when we talk about the SMO algorithm.

## 6 Optimal margin classifiers

正式开始讨论最优间隔分类器

Previously, we posed the following (primal) optimization problem for finding the optimal margin classifier:

$$\begin{aligned} \min_{\gamma, w, b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, m \end{aligned}$$

为什么引入对偶问题？

一是对偶问题往往更容易求解

二是可以自然的引入核函数，进而推广到非线性分类问题。

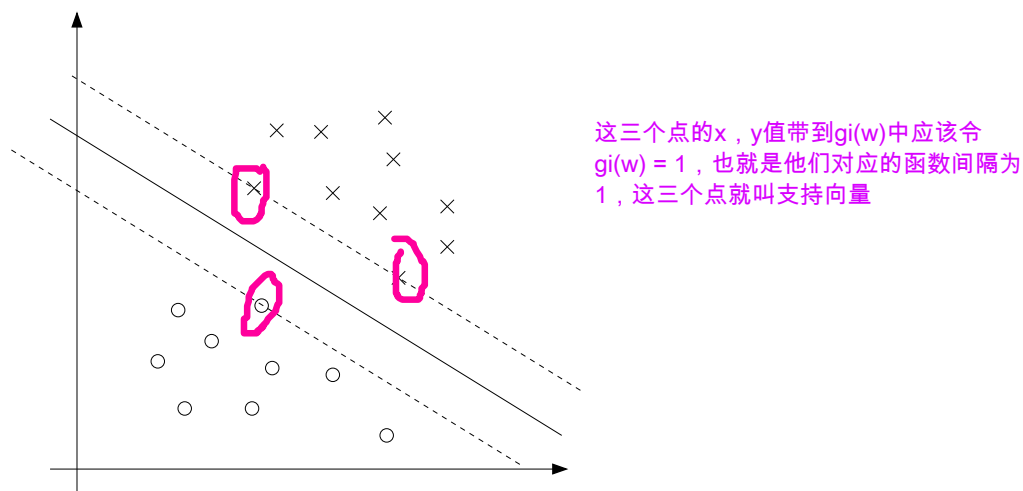
We can write the constraints as

$$g_i(w) = -y^{(i)}(w^T x^{(i)} + b) + 1 \leq 0.$$

We have one such constraint for each training example. Note that from the KKT dual complementarity condition, we will have  $\alpha_i > 0$  only for the training examples that have functional margin exactly equal to one (i.e., the ones

如果  $\alpha_i > 0$ ，那么  $g_i(w) = 0$ ， $g_i(w) = 0$  意味着函数间隔 = 1

corresponding to constraints that hold with equality,  $g_i(w) = 0$ ). Consider the figure below, in which a maximum margin separating hyperplane is shown by the solid line.



The points with the smallest margins are exactly the ones closest to the decision boundary; here, these are the three points (one negative and two positive examples) that lie on the dashed lines parallel to the decision boundary. Thus, only three of the  $\alpha_i$ 's—namely, the ones corresponding to these three training examples—will be non-zero at the optimal solution to our optimization problem. These three points are called the **support vectors** in this problem. The fact that the number of support vectors can be much smaller than the size the training set will be useful later.

Lets move on. Looking ahead, as we develop the dual form of the problem, one key idea to watch out for is that we'll try to write our algorithm in terms of only the inner product  $\langle x^{(i)}, x^{(j)} \rangle$  (think of this as  $(x^{(i)})^T x^{(j)}$ ) between points in the input feature space. The fact that we can express our algorithm in terms of these inner products will be key when we apply the kernel trick.

When we construct the Lagrangian for our optimization problem we have:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1]. \quad (8)$$

Note that there're only " $\alpha_i$ " but no " $\beta_i$ " Lagrange multipliers, since the problem has only inequality constraints.

Lets find the dual form of the problem. To do so, we need to first minimize  $\mathcal{L}(w, b, \alpha)$  with respect to  $w$  and  $b$  (for fixed  $\alpha$ ), to get  $\theta_{\mathcal{D}}$ , which we'll do by  
 固定 $\alpha$ , 调整  $w$ 和 $b$ , 最小化  $L(w, b, \alpha)$

setting the derivatives of  $\mathcal{L}$  with respect to  $w$  and  $b$  to zero. We have:

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0$$

This implies that

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}. \quad (9)$$

As for the derivative with respect to  $b$ , we obtain

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i y^{(i)} = 0. \quad (10)$$

If we take the definition of  $w$  in Equation (9) and plug that back into the Lagrangian (Equation 8), and simplify, we get

$$\mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} - b \underbrace{\sum_{i=1}^m \alpha_i y^{(i)}}_{=0}.$$

But from Equation (10), the last term must be zero, so we obtain

$$\mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)}.$$

Recall that we got to the equation above by minimizing  $\mathcal{L}$  with respect to  $w$  and  $b$ . Putting this together with the constraints  $\alpha_i \geq 0$  (that we always had) and the constraint (10), we obtain the following dual optimization problem:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle. \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0, \end{aligned}$$

You should also be able to verify that the conditions required for  $p^* = d^*$  and the KKT conditions (Equations 3–7) to hold are indeed satisfied in our optimization problem. Hence, we can solve the dual in lieu of solving the primal problem. Specifically, in the dual problem above, we have a maximization problem in which the parameters are the  $\alpha_i$ 's. We'll talk later

about the specific algorithm that we're going to use to solve the dual problem, but if we are indeed able to solve it (i.e., find the  $\alpha$ 's that maximize  $W(\alpha)$  subject to the constraints), then we can use Equation (9) to go back and find the optimal  $w$ 's as a function of the  $\alpha$ 's. Having found  $w^*$ , by considering the primal problem, it is also straightforward to find the optimal value for the intercept term  $b$  as

$$b^* = -\frac{\max_{i:y^{(i)}=-1} w^{*T} x^{(i)} + \min_{i:y^{(i)}=1} w^{*T} x^{(i)}}{2}. \quad (11)$$

(Check for yourself that this is correct.)

Before moving on, let's also take a more careful look at Equation (9), which gives the optimal value of  $w$  in terms of (the optimal value of)  $\alpha$ . Suppose we've fit our model's parameters to a training set, and now wish to make a prediction at a new point input  $x$ . We would then calculate  $w^T x + b$ , and predict  $y = 1$  if and only if this quantity is bigger than zero. But using (9), this quantity can also be written:

我们写出了分类函数的定义

可以看出来，对于新点的预测，我们只要

计算新点与训练集数据点的内积即可。

然后就是支持向量也在这里显现出来，

所有非支持向量的  $\alpha$  都为0，所以对于内积的计算

只要针对少量的支持向量而不是数据集中的所有点

$$w^T x + b = \left( \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x + b \quad (12)$$

$$= \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b. \quad (13)$$

Hence, if we've found the  $\alpha_i$ 's, in order to make a prediction, we have to calculate a quantity that depends only on the inner product between  $x$  and the points in the training set. Moreover, we saw earlier that the  $\alpha_i$ 's will all be zero except for the support vectors. Thus, many of the terms in the sum above will be zero, and we really need to find only the inner products between  $x$  and the support vectors (of which there is often only a small number) in order to calculate (13) and make our prediction.

By examining the dual form of the optimization problem, we gained significant insight into the structure of the problem, and were also able to write the entire algorithm in terms of only inner products between input feature vectors. In the next section, we will exploit this property to apply the kernels to our classification problem. The resulting algorithm, **support vector machines**, will be able to efficiently learn in very high dimensional spaces.

## 7 Kernels 核函数

Back in our discussion of linear regression, we had a problem in which the input  $x$  was the living area of a house, and we considered performing regres-

假设数据在低维空间是不可分的，那么可以将数据映射到高维空间，这个时候可以找一个映射  $f$  将  $x$  映射到高维空间，然后计算映射之后的高维空间中的向量的内积。

但是核函数可以做到直接在原来的低维空间中进行计算，而不需要显式地写出映射后的结果。

核函数能简化映射空间中的内积运算——刚好“碰巧”的是，在我们的 SVM 里需要计算的地方数据向量总是以内积的形式出现的。

sion using the features  $x$ ,  $x^2$  and  $x^3$  (say) to obtain a cubic function. To distinguish between these two sets of variables, we'll call the "original" input value the input **attributes** of a problem (in this case,  $x$ , the living area). When that is mapped to some new set of quantities that are then passed to the learning algorithm, we'll call those new quantities the input **features**. (Unfortunately, different authors use different terms to describe these two things, but we'll try to use this terminology consistently in these notes.) We will also let  $\phi$  denote the **feature mapping**, which maps from the attributes to the features. For instance, in our example, we had

$$\text{特征映射函数} \quad \phi(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}. \quad \text{本来就一个特征}x\text{，通过映射得到一个三维空间的向量}$$

Rather than applying SVMs using the original input attributes  $x$ , we may instead want to learn using some features  $\phi(x)$ . To do so, we simply need to go over our previous algorithm, and replace  $x$  everywhere in it with  $\phi(x)$ .

Since the algorithm can be written entirely in terms of the inner products  $\langle x, z \rangle$ , this means that we would replace all those inner products with  $\langle \phi(x), \phi(z) \rangle$ . Specifically, given a feature mapping  $\phi$ , we define the corresponding **Kernel** to be

$$K(x, z) = \phi(x)^T \phi(z). \quad \text{先映射，再计算}$$

Then, everywhere we previously had  $\langle x, z \rangle$  in our algorithm, we could simply replace it with  $K(x, z)$ , and our algorithm would now be learning using the features  $\phi$ .

Now, given  $\phi$ , we could easily compute  $K(x, z)$  by finding  $\phi(x)$  and  $\phi(z)$  and taking their inner product. But what's more interesting is that often,  $K(x, z)$  may be very inexpensive to calculate, even though  $\phi(x)$  itself may be very expensive to calculate (perhaps because it is an extremely high dimensional vector). In such settings, by using in our algorithm an efficient way to calculate  $K(x, z)$ , we can get SVMs to learn in the high dimensional feature space space given by  $\phi$ , but without ever having to explicitly find or represent vectors  $\phi(x)$ .

即使映射函数的计算代价比较高（如果先算映射的话，那就要先显示表示出来，然后再计算一个高维向量的内积），但核函数的计算代价不是那么高（不需要显示表示，直接基于低维向量进行计算即可）

Lets see an example. Suppose  $x, z \in \mathbb{R}^n$ , and consider

$$K(x, z) = (x^T z)^2.$$

We can also write this as

$$\begin{aligned}
 K(x, z) &= \left( \sum_{i=1}^n x_i z_i \right) \left( \sum_{j=1}^n x_j z_j \right) \\
 &= \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j \\
 &= \sum_{i,j=1}^n (x_i x_j)(z_i z_j)
 \end{aligned}$$

Thus, we see that  $K(x, z) = \phi(x)^T \phi(z)$ , where the feature mapping  $\phi$  is given (shown here for the case of  $n = 3$ ) by

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}.$$

Note that whereas calculating the high-dimensional  $\phi(x)$  requires  $O(n^2)$  time, finding  $K(x, z)$  takes only  $O(n)$  time—linear in the dimension of the input attributes.

For a related kernel, also consider

$$\begin{aligned}
 K(x, z) &= (x^T z + c)^2 \\
 &= \sum_{i,j=1}^n (x_i x_j)(z_i z_j) + \sum_{i=1}^n (\sqrt{2c} x_i)(\sqrt{2c} z_i) + c^2.
 \end{aligned}$$

(Check this yourself.) This corresponds to the feature mapping (again shown

for  $n = 3$ )

$$\phi(x) = \begin{bmatrix} x_1x_1 \\ x_1x_2 \\ x_1x_3 \\ x_2x_1 \\ x_2x_2 \\ x_2x_3 \\ x_3x_1 \\ x_3x_2 \\ x_3x_3 \\ \sqrt{2c}x_1 \\ \sqrt{2c}x_2 \\ \sqrt{2c}x_3 \\ c \end{bmatrix},$$

and the parameter  $c$  controls the relative weighting between the  $x_i$  (first order) and the  $x_ix_j$  (second order) terms.

More broadly, the kernel  $K(x, z) = (x^T z + c)^d$  corresponds to a feature mapping to an  $\binom{n+d}{d}$  feature space, corresponding of all monomials of the form  $x_{i_1}x_{i_2}\dots x_{i_k}$  that are up to order  $d$ . However, despite working in this  $O(n^d)$ -dimensional space, computing  $K(x, z)$  still takes only  $O(n)$  time, and hence we never need to explicitly represent feature vectors in this very high dimensional feature space.

Now, lets talk about a slightly different view of kernels. Intuitively, (and there are things wrong with this intuition, but nevermind), if  $\phi(x)$  and  $\phi(z)$  are close together, then we might expect  $K(x, z) = \phi(x)^T \phi(z)$  to be large. Conversely, if  $\phi(x)$  and  $\phi(z)$  are far apart—say nearly orthogonal to each other—then  $K(x, z) = \phi(x)^T \phi(z)$  will be small. So, we can think of  $K(x, z)$  as some measurement of how similar are  $\phi(x)$  and  $\phi(z)$ , or of how similar are  $x$  and  $z$ .

Given this intuition, suppose that for some learning problem that you're working on, you've come up with some function  $K(x, z)$  that you think might be a reasonable measure of how similar  $x$  and  $z$  are. For instance, perhaps you chose

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right).$$

This is a resonable measure of  $x$  and  $z$ 's similarity, and is close to 1 when  $x$  and  $z$  are close, and near 0 when  $x$  and  $z$  are far apart. Can we use this definition of  $K$  as the kernel in an SVM? In this particular example, the answer is yes. (This kernel is called the **Gaussian kernel**, and corresponds



下面介绍的是怎么证明一个核是可用的核函数？

to an infinite dimensional feature mapping  $\phi$ .) But more broadly, given some function  $K$ , how can we tell if it's a valid kernel; i.e., can we tell if there is some feature mapping  $\phi$  so that  $K(x, z) = \phi(x)^T \phi(z)$  for all  $x, z$ ?

Suppose for now that  $K$  is indeed a valid kernel corresponding to some feature mapping  $\phi$ . Now, consider some finite set of  $m$  points (not necessarily the training set)  $\{x^{(1)}, \dots, x^{(m)}\}$ , and let a square,  $m$ -by- $m$  matrix  $K$  be defined so that its  $(i, j)$ -entry is given by  $K_{ij} = K(x^{(i)}, x^{(j)})$ . This matrix is called the **Kernel matrix**. Note that we've overloaded the notation and used  $K$  to denote both the kernel function  $K(x, z)$  and the kernel matrix  $K$ , due to their obvious close relationship.

Now, if  $K$  is a valid Kernel, then  $K_{ij} = K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^T \phi(x^{(j)}) = \phi(x^{(j)})^T \phi(x^{(i)}) = K(x^{(j)}, x^{(i)}) = K_{ji}$ , and hence  $K$  must be symmetric. Moreover, letting  $\phi_k(x)$  denote the  $k$ -th coordinate of the vector  $\phi(x)$ , we find that for any vector  $z$ , we have

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j \\ &= \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(j)}) z_j \\ &= \sum_i \sum_j z_i \sum_k \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j \\ &= \sum_k \sum_i \sum_j z_i \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j \\ &= \sum_k \left( \sum_i z_i \phi_k(x^{(i)}) \right)^2 \\ &\geq 0. \end{aligned}$$

The second-to-last step above used the same trick as you saw in Problem set 1 Q1. Since  $z$  was arbitrary, this shows that  $K$  is positive semi-definite ( $K \geq 0$ ).

Hence, we've shown that if  $K$  is a valid kernel (i.e., if it corresponds to some feature mapping  $\phi$ ), then the corresponding Kernel matrix  $K \in \mathbb{R}^{m \times m}$  is symmetric positive semidefinite. More generally, this turns out to be not only a necessary, but also a sufficient, condition for  $K$  to be a valid kernel (also called a Mercer kernel). The following result is due to Mercer.<sup>5</sup>

---

<sup>5</sup>Many texts present Mercer's theorem in a slightly more complicated form involving  $L^2$  functions, but when the input attributes take values in  $\mathbb{R}^n$ , the version given here is equivalent.

**Theorem (Mercer).** Let  $K : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$  be given. Then for  $K$  to be a valid (Mercer) kernel, it is necessary and sufficient that for any  $\{x^{(1)}, \dots, x^{(m)}\}$ , ( $m < \infty$ ), the corresponding kernel matrix is symmetric positive semi-definite.

Given a function  $K$ , apart from trying to find a feature mapping  $\phi$  that corresponds to it, this theorem therefore gives another way of testing if it is a valid kernel. You'll also have a chance to play with these ideas more in problem set 2.

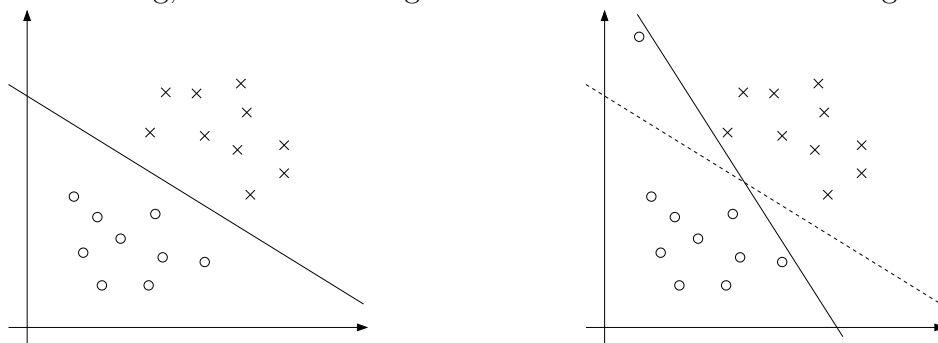
In class, we also briefly talked about a couple of other examples of kernels. For instance, consider the digit recognition problem, in which given an image (16x16 pixels) of a handwritten digit (0-9), we have to figure out which digit it was. Using either a simple polynomial kernel  $K(x, z) = (x^T z)^d$  or the Gaussian kernel, SVMs were able to obtain extremely good performance on this problem. This was particularly surprising since the input attributes  $x$  were just a 256-dimensional vector of the image pixel intensity values, and the system had no prior knowledge about vision, or even about which pixels are adjacent to which other ones. Another example that we briefly talked about in lecture was that if the objects  $x$  that we are trying to classify are strings (say,  $x$  is a list of amino acids, which strung together form a protein), then it seems hard to construct a reasonable, “small” set of features for most learning algorithms, especially if different strings have different lengths. However, consider letting  $\phi(x)$  be a feature vector that counts the number of occurrences of each length- $k$  substring in  $x$ . If we're considering strings of english alphabets, then there're  $26^k$  such strings. Hence,  $\phi(x)$  is a  $26^k$  dimensional vector; even for moderate values of  $k$ , this is probably too big for us to efficiently work with. (e.g.,  $26^4 \approx 460000$ .) However, using (dynamic programming-ish) string matching algorithms, it is possible to efficiently compute  $K(x, z) = \phi(x)^T \phi(z)$ , so that we can now implicitly work in this  $26^k$ -dimensional feature space, but without ever explicitly computing feature vectors in this space.

The application of kernels to support vector machines should already be clear and so we won't dwell too much longer on it here. Keep in mind however that the idea of kernels has significantly broader applicability than SVMs. Specifically, if you have any learning algorithm that you can write in terms of only inner products  $\langle x, z \rangle$  between input attribute vectors, then by replacing this with  $K(x, z)$  where  $K$  is a kernel, you can “magically” allow your algorithm to work efficiently in the high dimensional feature space corresponding to  $K$ . For instance, this kernel trick can be applied with the perceptron to derive a kernel perceptron algorithm. Many of the

algorithms that we'll see later in this class will also be amenable to this method, which has come to be known as the “kernel trick.”

## 8 Regularization and the non-separable case

The derivation of the SVM as presented so far assumed that the data is linearly separable. While mapping data to a high dimensional feature space via  $\phi$  does generally increase the likelihood that the data is separable, we can't guarantee that it always will be so. Also, in some cases it is not clear that finding a separating hyperplane is exactly what we'd want to do, since that might be susceptible to outliers. For instance, the left figure below shows an optimal margin classifier, and when a single outlier is added in the upper-left region (right figure), it causes the decision boundary to make a dramatic swing, and the resulting classifier has a much smaller margin.



引入松弛变量的目的有两个，一个是解决线性不可分的问题，一个是让超平面对异常点不那么敏感

To make the algorithm work for non-linearly separable datasets as well as be less sensitive to outliers, we reformulate our optimization (using  $\ell_1$  regularization) as follows:

$$\begin{aligned} \min_{\gamma, w, b} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & \gamma y^{(i)} (w^T x^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, m \\ & \xi_i \geq 0, \quad i = 1, \dots, m. \end{aligned}$$

其中  $\xi_i$  称为松弛变量 (slack variable)，对应数据点允许偏离的 functional margin 的量。当然，如果我们允许  $\xi_i$  任意大的话，那任意的超平面都是符合条件的了。所以，我们在原来的目标函数后面加上一项，使得这些松弛变量的总和也要最小

Thus, examples are now permitted to have (functional) margin less than 1, and if an example whose functional margin is  $1 - \xi_i$ , we would pay a cost of the objective function being increased by  $C\xi_i$ . The parameter  $C$  controls the relative weighting between the twin goals of making the  $\|w\|^2$  large (which we saw earlier makes the margin small) and of ensuring that most examples have functional margin at least 1.

As before, we can form the Lagrangian: 拉格朗日形式

$$\mathcal{L}(w, b, \xi, \alpha, r) = \frac{1}{2}w^T w + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y^{(i)}(x^T w + b) - 1 + \xi_i] - \sum_{i=1}^m r_i \xi_i.$$

Here, the  $\alpha_i$ 's and  $r_i$ 's are our Lagrange multipliers (constrained to be  $\geq 0$ ). We won't go through the derivation of the dual again in detail, but after setting the derivatives with respect to  $w$  and  $b$  to zero as before, substituting them back in, and simplifying, we obtain the following dual form of the problem:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle && \text{对偶形式} \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0, \end{aligned}$$

As before, we also have that  $w$  can be expressed in terms of the  $\alpha_i$ 's as given in Equation (9), so that after solving the dual problem, we can continue to use Equation (13) to make our predictions. Note that, somewhat surprisingly, in adding  $\ell_1$  regularization, the only change to the dual problem is that what was originally a constraint that  $0 \leq \alpha_i$  has now become  $0 \leq \alpha_i \leq C$ . The calculation for  $b^*$  also has to be modified (Equation 11 is no longer valid); see the comments in the next section/Platt's paper.

Also, the KKT dual-complementarity conditions (which in the next section will be useful for testing for the convergence of the SMO algorithm) are:

$$\alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1 \quad (14)$$

$$\alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1 \quad (15)$$

$$0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1. \quad (16)$$

Now, all that remains is to give an algorithm for actually solving the dual problem, which we will do in the next section.

再次总结引入对偶问题的目的：

一个是向量点积，引入核

一个是可以用 SMO 这样高校的算法

## 9 The SMO algorithm

The SMO (sequential minimal optimization) algorithm, due to John Platt, gives an efficient way of solving the dual problem arising from the derivation

of the SVM. Partly to motivate the SMO algorithm, and partly because it's interesting in its own right, let's first take another digression to talk about the coordinate ascent algorithm.

## 9.1 Coordinate ascent 坐标上升算法

Consider trying to solve the unconstrained optimization problem

$$\max_{\alpha} W(\alpha_1, \alpha_2, \dots, \alpha_m).$$

Here, we think of  $W$  as just some function of the parameters  $\alpha_i$ 's, and for now ignore any relationship between this problem and SVMs. We've already seen two optimization algorithms, gradient ascent and Newton's method. The new algorithm we're going to consider here is called **coordinate ascent**:

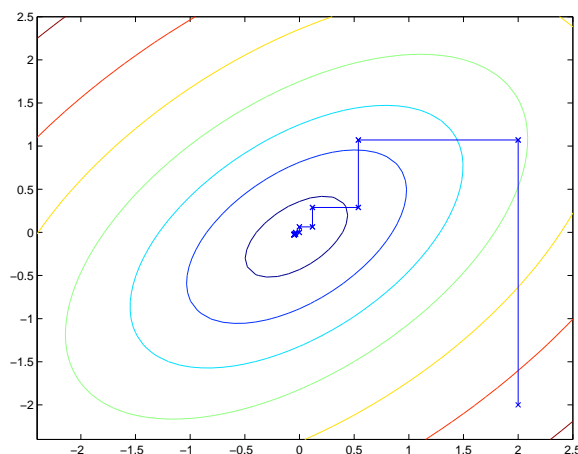
```

Loop until convergence: {
    For  $i = 1, \dots, m$ , {
         $\alpha_i := \arg \max_{\hat{\alpha}_i} W(\alpha_1, \dots, \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, \dots, \alpha_m)$ .
    }
}

```

Thus, in the innermost loop of this algorithm, we will hold all the variables except for some  $\alpha_i$  fixed, and reoptimize  $W$  with respect to just the parameter  $\alpha_i$ . In the version of this method presented here, the inner-loop reoptimizes the variables in order  $\alpha_1, \alpha_2, \dots, \alpha_m, \alpha_1, \alpha_2, \dots$ . (A more sophisticated version might choose other orderings; for instance, we may choose the next variable to update according to which one we expect to allow us to make the largest increase in  $W(\alpha)$ .)

When the function  $W$  happens to be of such a form that the “arg max” in the inner loop can be performed efficiently, then coordinate ascent can be a fairly efficient algorithm. Here's a picture of coordinate ascent in action:



总是会沿着和坐标轴平行的方向取到最小值

The ellipses in the figure are the contours of a quadratic function that we want to optimize. Coordinate ascent was initialized at  $(2, -2)$ , and also plotted in the figure is the path that it took on its way to the global maximum. Notice that on each step, coordinate ascent takes a step that's parallel to one of the axes, since only one variable is being optimized at a time.

## 9.2 SMO 序列最小化算法，最小化是指我们希望一次改变最小数目的 $\alpha_i$ ，这里一次选择改变两个 $\alpha_i$

We close off the discussion of SVMs by sketching the derivation of the SMO algorithm. Some details will be left to the homework, and for others you may refer to the paper excerpt handed out in class.

Here's the (dual) optimization problem that we want to solve:

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle. \quad (17)$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \quad (18)$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0. \quad (19)$$

Lets say we have set of  $\alpha_i$ 's that satisfy the constraints (18-19). Now, suppose we want to hold  $\alpha_2, \dots, \alpha_m$  fixed, and take a coordinate ascent step and reoptimize the objective with respect to  $\alpha_1$ . Can we make any progress? The answer is no, because the constraint (19) ensures that

$$\alpha_1 y^{(1)} = - \sum_{i=2}^m \alpha_i y^{(i)}.$$

Or, by multiplying both sides by  $y^{(1)}$ , we equivalently have

$$\alpha_1 = -y^{(1)} \sum_{i=2}^m \alpha_i y^{(i)}.$$

如果我们像坐标上升算法那样一次只改变一个 $\alpha$ 的值，那么这个等式将不会成立，也就是说我们要优化的对偶问题将不再满足约束条件，所以我们一次至少同时改变两个 $\alpha_i$

(This step used the fact that  $y^{(1)} \in \{-1, 1\}$ , and hence  $(y^{(1)})^2 = 1$ .) Hence,  $\alpha_1$  is exactly determined by the other  $\alpha_i$ 's, and if we were to hold  $\alpha_2, \dots, \alpha_m$  fixed, then we can't make any change to  $\alpha_1$  without violating the constraint (19) in the optimization problem.

Thus, if we want to update some subset of the  $\alpha_i$ 's, we must update at least two of them simultaneously in order to keep satisfying the constraints. This motivates the SMO algorithm, which simply does the following:

Repeat till convergence {

1. Select some pair  $\alpha_i$  and  $\alpha_j$  to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
2. Reoptimize  $W(\alpha)$  with respect to  $\alpha_i$  and  $\alpha_j$ , while holding all the other  $\alpha_k$ 's ( $k \neq i, j$ ) fixed.

}

To test for convergence of this algorithm, we can check whether the KKT conditions (Equations 14-16) are satisfied to within some *tol*. Here, *tol* is the convergence tolerance parameter, and is typically set to around 0.01 to 0.001. (See the paper and pseudocode for details.)

The key reason that SMO is an efficient algorithm is that the update to  $\alpha_i, \alpha_j$  can be computed very efficiently. Let's now briefly sketch the main ideas for deriving the efficient update.

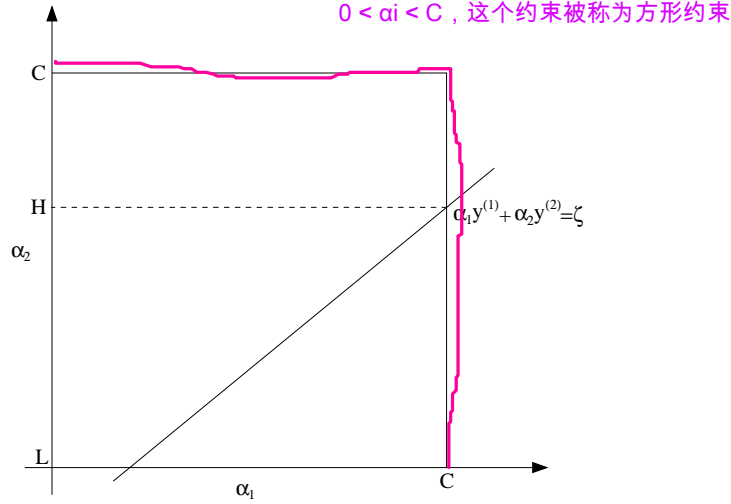
Let's say we currently have some setting of the  $\alpha_i$ 's that satisfy the constraints (18-19), and suppose we've decided to hold  $\alpha_3, \dots, \alpha_m$  fixed, and want to reoptimize  $W(\alpha_1, \alpha_2, \dots, \alpha_m)$  with respect to  $\alpha_1$  and  $\alpha_2$  (subject to the constraints). From (19), we require that

$$\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = - \sum_{i=3}^m \alpha_i y^{(i)}.$$

Since the right hand side is fixed (as we've fixed  $\alpha_3, \dots, \alpha_m$ ), we can just let it be denoted by some constant  $\zeta$ :

$$\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta. \tag{20}$$

We can thus picture the constraints on  $\alpha_1$  and  $\alpha_2$  as follows:



From the constraints (18), we know that  $\alpha_1$  and  $\alpha_2$  must lie within the box  $[0, C] \times [0, C]$  shown. Also plotted is the line  $\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta$ , on which we know  $\alpha_1$  and  $\alpha_2$  must lie. Note also that, from these constraints, we know  $L \leq \alpha_2 \leq H$ ; otherwise,  $(\alpha_1, \alpha_2)$  can't simultaneously satisfy both the box and the straight line constraint. In this example,  $L = 0$ . But depending on what the line  $\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta$  looks like, this won't always necessarily be the case; but more generally, there will be some lower-bound  $L$  and some upper-bound  $H$  on the permissible values for  $\alpha_2$  that will ensure that  $\alpha_1, \alpha_2$  lie within the box  $[0, C] \times [0, C]$ .

Using Equation (20), we can also write  $\alpha_1$  as a function of  $\alpha_2$ :

$$\alpha_1 = (\zeta - \alpha_2 y^{(2)}) y^{(1)}.$$

(Check this derivation yourself; we again used the fact that  $y^{(1)} \in \{-1, 1\}$  so that  $(y^{(1)})^2 = 1$ .) Hence, the objective  $W(\alpha)$  can be written

$$W(\alpha_1, \alpha_2, \dots, \alpha_m) = W((\zeta - \alpha_2 y^{(2)}) y^{(1)}, \alpha_2, \dots, \alpha_m).$$

Treating  $\alpha_3, \dots, \alpha_m$  as constants, you should be able to verify that this is just some quadratic function in  $\alpha_2$ . I.e., this can also be expressed in the form  $a\alpha_2^2 + b\alpha_2 + c$  for some appropriate  $a$ ,  $b$ , and  $c$ . If we ignore the “box” constraints (18) (or, equivalently, that  $L \leq \alpha_2 \leq H$ ), then we can easily maximize this quadratic function by setting its derivative to zero and solving. We'll let  $\alpha_2^{new, unclipped}$  denote the resulting value of  $\alpha_2$ . You should also be able to convince yourself that if we had instead wanted to maximize  $W$  with respect to  $\alpha_2$  but subject to the box constraint, then we can find the resulting value optimal simply by taking  $\alpha_2^{new, unclipped}$  and “clipping” it to lie in the



$[L, H]$  interval, to get

$$\alpha_2^{new} = \begin{cases} H & \text{if } \alpha_2^{new,unclipped} > H \\ \alpha_2^{new,unclipped} & \text{if } L \leq \alpha_2^{new,unclipped} \leq H \\ L & \text{if } \alpha_2^{new,unclipped} < L \end{cases}$$

Finally, having found the  $\alpha_2^{new}$ , we can use Equation (20) to go back and find the optimal value of  $\alpha_1^{new}$ .

There're a couple more details that are quite easy but that we'll leave you to read about yourself in Platt's paper: One is the choice of the heuristics used to select the next  $\alpha_i$ ,  $\alpha_j$  to update; the other is how to update  $b$  as the SMO algorithm is run.