

第七章 支持向量机

支持向量机（support vector machine，简称SVM）于1964年由Vapnik和Chervonenkis建立，在上世纪90年代获得快速发展并衍生出一系列改进和扩展算法，在人像识别、文本分类、手写字识别及生物信息学等领域获得广泛应用。

1 间隔与支持向量

以二分类问题为例，两个类别分别以“+1”和“-1”表示，对于训练样本集 $D=\{(x_0, y_0), (x_1, y_1), \dots, (x_{n-1}, y_{n-1})\}$, $y_i \in \{-1, +1\}$ 。要实现正确分类，就是要基于训练集 D 在样本空间中找到一个划分超平面，将不同类别的样本分开。对于样本有两个特征的情况，如图7-1所示，此时划分超平面为直线。能将训练样本分开的直线有很多，我们应该取哪一条呢？

显然应该去找位于两类训练样本“正中间”的直线，即图7-1中的粗线。因为该直线对训练样本局部扰动的容忍性最好。

例如，由于训练集的局限性或噪声的因素，训练集外的样本可能比图7-1中的样本更接近两个类的分隔界，这将使许多划分直线出现错误，而粗线受影响最小。换言之，这条直线对未见样本的泛化能力最强。

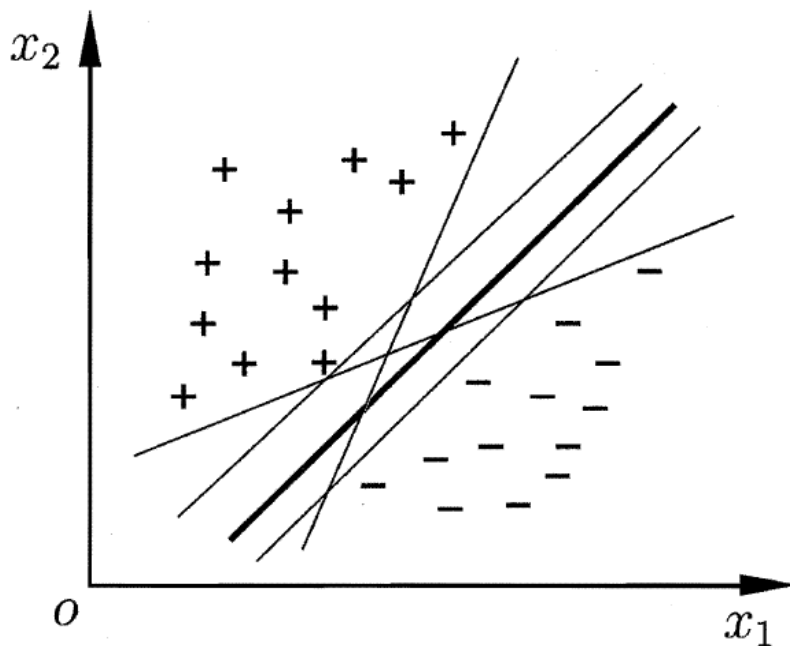


图7-1 将两类训练样本正确分类的多个超平面

在样本空间中，划分超平面可通过如下线性方程来描述：

$$\mathbf{w}'^T \mathbf{x} + b' = w'_0 x_0 + w'_1 x_1 + \cdots + w'_{m-1} x_{m-1} + b' = 0$$

其中 \mathbf{w}' 为法向量， b' 为位移项。显然划分超平面可被法向量和位移确定，我们记该超平面为 (\mathbf{w}', b') 。样本空间中任意点 \mathbf{x} 到超平面 (\mathbf{w}', b') 的距离可写为：

$$r = \frac{|\mathbf{w}'^T \mathbf{x} + b'|}{\|\mathbf{w}'\|}$$

假设超平面 (\mathbf{w}', b') 能将训练样本正确分类，过正类样本点中距离分隔超平面最近的点，作与超平面相平行的超平面，其方程为：

$$\mathbf{w}'^T \mathbf{x} + b' = d$$

相应地，由于超平面位于样本点中间，过负类样本点中距离分隔超平面最近的点，与超平面相平行的超平面方程为：

$$\mathbf{w}'^T \mathbf{x} + b' = -d$$

这三个超平面方程，都同时除以一个数d，不会改变它们的位置，因此，三个超平面方程转变为：

$$\text{分隔超平面} \quad \mathbf{w}^T \mathbf{x} + b = 0$$

$$\text{正类样本一侧} \quad \mathbf{w}^T \mathbf{x} + b = 1$$

$$\text{负类样本一侧} \quad \mathbf{w}^T \mathbf{x} + b = -1$$

对于任意样本点 $(\mathbf{x}_i, y_i) \in D$ ，有：

$$\begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 1, y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq -1, y_i = -1 \end{cases}$$

距离分隔超平面最近的训练样本点使上式中等号成立，它们被称为“支持向量”，两个异类支持向量到超平面的距离之和为：

$$\gamma = \frac{2}{\|\mathbf{w}\|}$$

γ 称为间隔（margin）。

要找到间隔最大的划分超平面，也就是要找到能满足约束的参数 \mathbf{w} 和 b ，使得 γ 最大，即：

$$\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|}$$

$$\text{s. t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, i = 0, 1, \dots, n-1$$

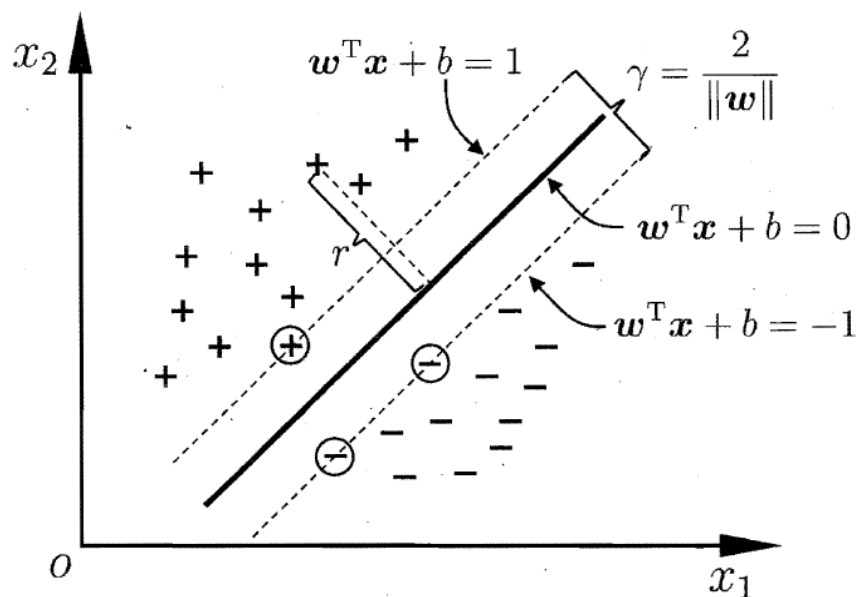


图7.2 支持向量与间隔

$$\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|}$$

显然，最大化 $\|\mathbf{w}\|^{-1}$ 等价于最小化 $\|\mathbf{w}\|$ 或 $\|\mathbf{w}\|^2$

因此，优化问题可重写为：

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s. t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, i = 0, 1, \dots, n-1$$

这就是支持向量机的基本形式。

可以看出，支持向量机最终转化为一个有约束的优化问题，虽然可以利用一些现成的优化程序包求解，但通常采用拉格朗日乘子法将其转化为对偶问题求解。

为此，我们先来介绍有约束优化问题的KKT条件。

2 KKT条件

我们知道，对于无约束优化问题 $\min f(\mathbf{x})$ ，令函数对所有自变量的偏导为0，即函数梯度 $\nabla_{\mathbf{x}} f(\mathbf{x})=0$ ，求解即可。

对于具有等式约束条件的优化问题：

$$\min f(\mathbf{x})$$

$$\text{s.t. } h_m(\mathbf{x}) = 0, m = 0, 1, \dots, M-1$$

如果约束方程足够简单，可通过消元将其转变为无约束问题。但如果约束方程较为复杂，可以利用拉格朗日乘子法解决。

定义拉格朗日函数：

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{m=0}^{M-1} \lambda_m h_m(\mathbf{x})$$

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{m=0}^{M-1} \lambda_m h_m(\mathbf{x})$$

然后令 L 对 \mathbf{x} 及 $\boldsymbol{\lambda}$ 的梯度为0:

$$\begin{cases} \nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) = 0 \\ \nabla_{\boldsymbol{\lambda}} L(\mathbf{x}, \boldsymbol{\lambda}) = 0 \end{cases}$$

解上式求得 \mathbf{x} 及 $\boldsymbol{\lambda}$ 的值。

为了解释这样做的原因，我们以二维问题 $f(\mathbf{x}, \mathbf{y})$ 的优化为例，假设只有一个等式约束条件 $h(\mathbf{x}, \mathbf{y})=0$ 。

图7-3画出了 f 的等高线（虚线）及约束条件线（实线）。

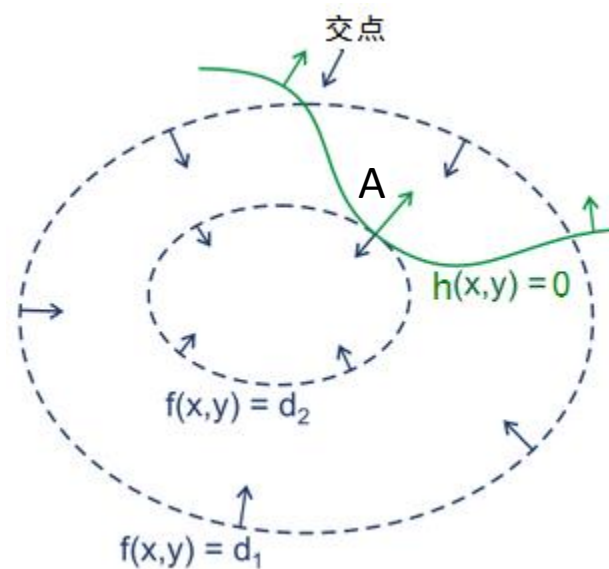


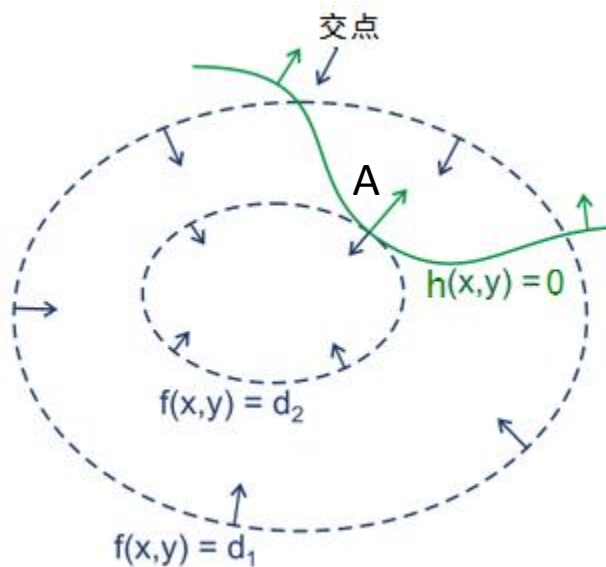
图7-3 具有等式约束条件的优化问题示意图

显然，可行解一定在约束线上并使得 f 最小的点处，这个点一定在约束线与某条等高线相切的位置，即图中A点处，此点处约束线梯度 ∇h 与 f 的梯度 ∇f 方向相同或相反，即有：

$$\nabla_x f(\mathbf{x}) + \lambda \nabla_x h(\mathbf{x}) = 0$$

满足上式且满足约束条件 $h(\mathbf{x})=0$ 的点即为可能的解，联立起来，正好就是拉格朗日乘子法的结果：

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{m=0}^{M-1} \lambda_m h_m(\mathbf{x})$$



$$\begin{cases} \nabla_x L(\mathbf{x}, \boldsymbol{\lambda}) = 0 \\ \nabla_{\boldsymbol{\lambda}} L(\mathbf{x}, \boldsymbol{\lambda}) = 0 \end{cases}$$

对于即包含等式约束又包含不等式约束的优化问题：

$$\min f(\mathbf{x})$$

$$\text{s.t. } h_m(\mathbf{x}) = 0, m = 0, 1, \dots, M - 1$$

$$g_k(\mathbf{x}) \leq 0, k = 0, 1, \dots, K - 1$$

KKT（Karush-Kuhn-Tucker）条件给出了最优解的必要条件，KKT条件是拉格朗日乘子法的推广。先定义拉格朗日函数：

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \sum_{m=0}^{M-1} \lambda_m h_m(\mathbf{x}) + \sum_{k=0}^{K-1} \alpha_k g_k(\mathbf{x})$$

KKT条件如下：

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = 0$$

$$h_m(\mathbf{x}) = 0, m = 0, 1, \dots, M - 1$$

$$\alpha_k g_k(\mathbf{x}) = 0, k = 0, 1, \dots, K - 1$$

$$\alpha_k \geq 0$$

可以看到，KKT条件中对等式约束的处理与上面的拉格朗日乘子法完全相同，下面我们以二维问题 $f(x_1, x_2)$ 的优化为例，对不等式约束条件的处理进行解释，为简单起见，考虑只有一个不等式约束的情况。如图7-4所示，依据 $f(x_1, x_2)$ 无约束极小值点 x^* 的位置分两种情况

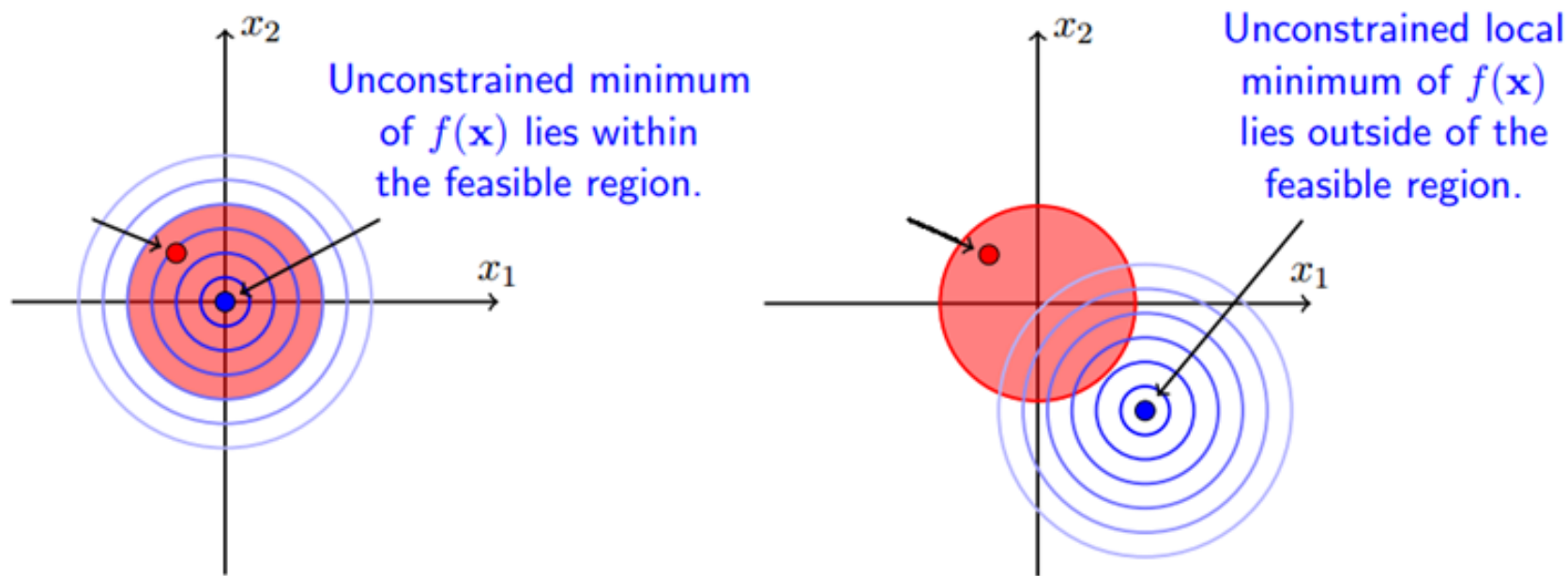


图7-4 不等式约束条件优化问题图示

(1) x^* 位于 $g(x)<0$ 的区域内（图7-4a），此时必然存在 x^* 的一个邻域全部位于 $g(x)\leq 0$ 的约束区域之内，这时按无约束问题求解即可，约束条件 $g(x)\leq 0$ 不起作用。

(2) x^* 位于 $g(x)<0$ 的区域之外（图7-4b），此时可行的极小值点必然在 $g(x)\leq 0$ 区域内且使得 $f(x_1, x_2)$ 最小的位置，即处于 $g(x)=0$ 且与 f 的某条等高线相切的位置，这与等式约束的情况相同。

综合以上两种情况，可得：

$$\alpha g(x) = 0$$

$\alpha=0$ 对应于第(1)种情况，此时约束条件不起作用，令 $\alpha=0$ 消去约束条件即可；

$g(x)=0$ 对应于第(2)种情况，此时可行解落在约束边界上。

最后解释为什么要求 $\alpha \geq 0$ 。对于第(1)种情况， $\alpha=0$ 显然满足。对于第(2)种情况，极小值出现在 $g(x)=0$ 与 f 的某条等高线相切的位置，该点满足：

$$\nabla_x f(x) + \alpha \nabla_x g(x) = 0$$

将该式重写为：

$$-\nabla_x f(x) = \alpha \nabla_x g(x)$$

由于梯度表示一个函数增速最大的方向，相应地， $-\nabla_x f(x)$ 表示 f 减小最快的方向。如果 $\alpha < 0$ ，则 $\alpha \nabla_x g(x)$ 也指向 g 减小最快的方向，此方向指向符合约束条件的区域，从切点沿此方向前进，必然存在一点即在可行区域，又使得 f 减小，该点 f 值小于切点，这与切点为极小值点的假设矛盾。因此必须 $\alpha > 0$ 。

总之，对(1)、(2)两种情况，都有 $\alpha \geq 0$ 。

3 对偶问题

一个优化问题可以从两个角度考察，即“主问题” (primal problem)和“对偶问题” (dual problem)。对于上节的主问题：

$$\begin{aligned} \min f(\mathbf{x}) \\ \text{s. t. } h_m(\mathbf{x}) = 0, m = 0, 1, \dots, M-1 \\ g_k(\mathbf{x}) \leq 0, k = 0, 1, \dots, K-1 \end{aligned}$$

基于其拉格朗日函数：

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \sum_{m=0}^{M-1} \lambda_m h_m(\mathbf{x}) + \sum_{k=0}^{K-1} \alpha_k g_k(\mathbf{x})$$

定义其对偶函数为：

$$\Gamma(\boldsymbol{\lambda}, \boldsymbol{\alpha}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \sum_{m=0}^{M-1} \lambda_m h_m(\mathbf{x}) + \sum_{k=0}^{K-1} \alpha_k g_k(\mathbf{x})$$

$$\Gamma(\lambda, \alpha) = \inf_x L(x, \lambda, \alpha) = f(x) + \sum_{m=0}^{M-1} \lambda_m h_m(x) + \sum_{k=0}^{K-1} \alpha_k g_k(x)$$

其中 \inf 表示下确界，即最大下界。若 \tilde{x} 为主问题可行域中的点，则对任意 $\alpha \geq 0$ 和 λ 都有

$$\sum_{m=0}^{M-1} \lambda_m h_m(\tilde{x}) + \sum_{k=0}^{K-1} \alpha_k g_k(\tilde{x}) \leq 0$$

进而有：

$$\Gamma(\lambda, \alpha) = \inf_x L(x, \lambda, \alpha) \leq L(\tilde{x}, \lambda, \alpha) \leq f(\tilde{x})$$

若主问题的最小值为 p^* ，则对任意 $\alpha \geq 0$ 和 λ 都有：

$$\Gamma(\lambda, \alpha) \leq p^*$$

$$\Gamma(\lambda, \alpha) \leq p^*$$

因此对偶函数给出了主问题最小值的下界。显然这个下界取决于 α 和 λ 的值。于是一个很自然的问题是：基于对偶函数能获得的最大下界是什么？这就引出了优化问题：

$$\max_{\lambda, \alpha} \Gamma(\lambda, \alpha), \text{ s. t. } \alpha \geq 0$$

上式就是前述主问题的对偶问题，其中 λ 和 α 称为“对偶变量”。

针对于支持向量机的基本形式：

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

将约束条件重写为 $g(\mathbf{x}) \leq 0$ 的形式：

$$1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0$$

拉格朗日函数为：

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=0}^{n-1} \alpha_i [1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)]$$

令 $L(\mathbf{w}, b, \alpha)$ 对 \mathbf{w} 和 b 的偏导为0：

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial w_j} = w_j - \sum_{i=0}^{n-1} \alpha_i y_i x_{ij} = 0$$

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = - \sum_{i=0}^{n-1} \alpha_i y_i = 0$$

因而有：

$$\mathbf{w} = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i \quad \sum_{i=0}^{n-1} \alpha_i y_i = 0$$

$$\mathbf{w} = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i \quad \sum_{i=0}^{n-1} \alpha_i y_i = 0$$

将上面两式代回拉格朗日函数：

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=0}^{n-1} \alpha_i [1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)]$$

$$= \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=0}^{n-1} \alpha_i - \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - b \sum_{i=0}^{n-1} \alpha_i y_i$$

$$= \sum_{i=0}^{n-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

再考虑到 $\sum_{i=0}^{n-1} \alpha_i y_i = 0$ 的约束，从而得到对偶问题：

$$\max_{\alpha} \sum_{i=0}^{n-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{s. t. } \sum_{i=0}^{n-1} \alpha_i y_i = 0$$

$$\alpha_i \geq 0, i = 0, 1, \dots, n-1$$

解对偶问题得到 α 后，将其代入 $\mathbf{w} = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i$ 即可求出 \mathbf{w} 。

\mathbf{b} 可通过任一支持向量 (\mathbf{x}_s, y_s) 求出，因为在支持向量处满足 $y_s(\mathbf{w}^T \mathbf{x}_s + \mathbf{b}) = 1$ 。现实任务中常采用一种更稳健的做法：使用所有支持向量求解的平均值：

$$b = \frac{1}{|S|} \sum_{s \in S} \left(\frac{1}{y_s} - \mathbf{w}^T \mathbf{x}_s \right) = \frac{1}{|S|} \sum_{s \in S} \left(y_s - \sum_{i \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_s \right)$$

$$b = \frac{1}{|S|} \sum_{s \in S} \left(\frac{1}{y_s} - \mathbf{w}^T \mathbf{x}_s \right) = \frac{1}{|S|} \sum_{s \in S} \left(y_s - \sum_{i \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_s \right)$$

其中， $S=\{i|\alpha_i>0, i=0,1,\cdots,n-1\}$ 为所有支持向量的下标集， $|S|$ 表示 S 集合的长度。

解出 \mathbf{w} 和 b 即可得到模型：

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_j + b$$

解得的任何 α_i 都对应一个样本 (\mathbf{x}_i, y_i) ，根据KKT条件，如果 $\alpha_i=0$ ，则该样本不会对 $f(\mathbf{x})$ 有任何影响；若 $\alpha_i>0$ ，则必有 $y_i f(\mathbf{x}_i)=1$ ，所对应的样本是一个支持向量。这是支持向量机的一个重要性质：训练完成后，大部分的训练样本都不需保留，最终模型仅与支持向量有关。

4. 软间隔支持向量机

前面的讨论中假定训练样本是线性可分的，即存在一个超平面能将不同类的样本完全分开。但现实问题也许不是线性可分的，解决该问题的办法之一是允许支持向量机在一些样本上出错，为此引入“软间隔”（soft margin）的概念，即允许某些样本不满足 $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$ 的约束条件。相应地，把前面要求所有样本都满足该约束条件即所有样本都必须分类正确的情况称为“硬间隔”（hard margin）。

对软间隔的支持向量机，在最大化间隔的同时，希望不满足约束的样本应尽可能少，于是优化目标可写为：

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=0}^{n-1} \ell_{0/1}(y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1)$$

其中 $C>0$ 是一个常数, $\ell_{0/1}$ 是"0/1损失函数":

$$\ell_{0/1}(z) = \begin{cases} 1, & z < 0 \\ 0, & z \geq 0 \end{cases}$$

显然, 当 C 为无穷大时, 优化目标函数将迫使所有样本均满足约束, 于是软间隔等价于硬间隔, 当 C 取有限值时, 软间隔允许一些样本不满足约束。

引入“松弛变量” $\xi_i \geq 0$, 可将上式重写为:

$$\min_{\mathbf{w}, b, \xi_i} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=0}^{n-1} \xi_i$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0, i = 0, 1, \dots, n-1$$

这就是常用的“软间隔支持向量机”。显然, 每个样本都有一个对应的松弛变量, 用以表征该样本不满足约束的程度。

对软间隔支持向量机，定义拉格朗日函数：

$$L(\mathbf{w}, b, \alpha, \xi, \mu) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=0}^{n-1} \xi_i + \sum_{i=0}^{n-1} \alpha_i (1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b)) - \sum_{i=0}^{n-1} \mu_i \xi_i$$

令 $L(\mathbf{w}, b, \alpha, \xi, \mu)$ 对 \mathbf{w} , b , ξ_i 的偏导为零可得：

$$\mathbf{w} = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i \quad \sum_{i=0}^{n-1} \alpha_i y_i = 0 \quad C = \alpha_i + \mu_i$$

将上面三式代入拉格朗日函数得到对偶问题：

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=0}^{n-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s. t.} \quad & \sum_{i=0}^{n-1} \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, i = 0, 1, \dots, n-1 \end{aligned}$$

将软间隔的对偶问题与硬间隔比较，两者唯一的差别就在于对偶变量的约束不同：软间隔是 $0 \leq \alpha_i \leq C$ ，而硬间隔是 $0 \leq \alpha_i$ 。对于软间隔支持向量机，KKT条件要求：

$$\begin{cases} \alpha_i \geq 0, \mu_i \geq 0 \\ y_i f(\mathbf{x}_i) - 1 + \xi_i \geq 0 \\ \alpha_i (y_i f(\mathbf{x}_i) - 1 + \xi_i) = 0 \\ \xi_i \geq 0, \mu_i \xi_i = 0 \end{cases}$$

对于任意训练样本 (\mathbf{x}_i, y_i) ，总有 $\alpha_i=0$ 或 $y_i f(\mathbf{x}_i)=1-\xi_i$ 。

- 若 $\alpha_i=0$ ，则该样本不会对 $f(\mathbf{x})$ 有任何影响；
- 若 $\alpha_i>0$ ，则必有 $y_i f(\mathbf{x}_i)=1-\xi_i$ ，即该样本是支持向量：又由于 $C=\alpha_i+\mu_i$ ，若 $\alpha_i<C$ ，则 $\mu_i>0$ ，由于 $\mu_i \xi_i=0$ 则 $\xi_i=0$ ，即该样本恰好落在最大间隔边界上；若 $\alpha_i=C$ ，则有 $\mu_i=0$ ，此时若 $\xi_i \leq 1$ 则该样本落在最大间隔内部，若 $\xi_i>1$ 则该样本被错误分类。

$$\max_{\alpha} \sum_{i=0}^{n-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

5. SMO算法

对偶问题是一个二次规划问题，可使用通用的二次规划算法来求解，但对于实际任务会有很大的计算开销，因此，人们根据问题本身的特性，开发了很多高效算法，应用最广泛的是SMO（sequential minimal optimization）算法。

SMO算法的核心思想是将一系列 α 的优化问题转化为一个个小的优化问题求解，即每一步选择两个 α 进行优化，其它 α 值暂时固定不变。之所以选择两个 α ，是因为要符合 $\sum \alpha_i y_i = 0$ 的约束，无法只选择一个 α 进行优化。然后对两个 α 的二次优化问题可以得到其解析解，从而加速整个优化过程。

不失一般性，假设选出 α_0 和 α_1 来进行优化，为了方便，我们暂时以 $K_{ij}=\mathbf{x}_i^T\mathbf{x}_j$ 表示两个向量 \mathbf{x}_i 与 \mathbf{x}_j 的内积。将对偶问题写为极小值问题，并展开为与 α_0 、 α_1 相关的部分和无关的部分（在本次优化中可视为常数，用 c_1 表示）：

$$\min_{\alpha} \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j K_{ij} - \sum_{i=0}^{n-1} \alpha_i$$

$$\min_{\alpha} \frac{1}{2} K_{00} \alpha_0^2 + \frac{1}{2} K_{11} \alpha_1^2 + y_0 y_1 K_{01} \alpha_0 \alpha_1$$

$$+ y_0 \alpha_0 \sum_{i=2}^{n-1} \alpha_i y_i K_{0i} + y_1 \alpha_1 \sum_{i=2}^{n-1} \alpha_i y_i K_{1i} - (\alpha_0 + \alpha_1) + c_1$$

利用约束条件可消去一个 α ：

$$y_0 \alpha_0 + y_1 \alpha_1 = c_2 \quad \Longrightarrow \quad \alpha_0 = y_0 (c_2 - y_1 \alpha_1)$$

其中 c_2 为常数，在本次优化中固定，另外用到了 $y_i^2=1$ 。

将 α_0 的表达式代入优化问题消去 α_0 得到单变量优化问题：

$$\begin{aligned} \min \Psi(\alpha_1) = & \frac{1}{2} K_{00} (c_2 - y_1 \alpha_1)^2 + \frac{1}{2} K_{11} \alpha_1^2 + y_1 K_{01} \alpha_1 (c_2 - y_1 \alpha_1) \\ & + (c_2 - y_1 \alpha_1) \sum_{i=2}^{n-1} \alpha_i y_i K_{0i} + y_1 \alpha_1 \sum_{i=2}^{n-1} \alpha_i y_i K_{1i} - y_0 (c_2 - y_1 \alpha_1) - \alpha_1 + c_1 \end{aligned}$$

令 $d\Psi/(d\alpha_1)=0$ ，得到：

$$(K_{00} + K_{11} - 2K_{01})\alpha_1 = y_1 K_{00} c_2 - y_1 K_{01} c_2 + y_1 \sum_{i=2}^{n-1} \alpha_i y_i K_{0i} - y_1 \sum_{i=2}^{n-1} \alpha_i y_i K_{1i} - y_0 y_1 + 1$$

这样就得到了 α_1 的更新公式，其右侧可通过当前值（old）来计算：

$$c_2 = y_0 \alpha_0^{\text{old}} + y_1 \alpha_1^{\text{old}}$$

$$f(\mathbf{x}_0) = \mathbf{w}^T \mathbf{x}_0 + b = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_0 + b = \sum_{i=0}^{n-1} \alpha_i y_i K_{0i} + b \xrightarrow{\text{yields}}$$

$$\sum_{i=2}^{n-1} \alpha_i y_i K_{0i} = f(\mathbf{x}_0) - \alpha_0^{\text{old}} y_0 K_{00} - \alpha_1^{\text{old}} y_1 K_{01} - b$$

$$f(\mathbf{x}_1) = \mathbf{w}^T \mathbf{x}_1 + b = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_1 + b = \sum_{i=0}^{n-1} \alpha_i y_i K_{1i} + b \xrightarrow{\text{yields}}$$

$$\sum_{i=2}^{n-1} \alpha_i y_i K_{1i} = f(\mathbf{x}_1) - \alpha_0^{\text{old}} y_0 K_{01} - \alpha_1^{\text{old}} y_1 K_{11} - b$$

将 c_2 、 $\sum \alpha_i y_i K_{0i}$ 、 $\sum \alpha_i y_i K_{1i}$ 表达式代入迭代式右端得到：

$$(K_{00} + K_{11} - 2K_{01})\alpha_1^{\text{new}} = (K_{00} + K_{11} - 2K_{01})\alpha_1^{\text{old}} + y_1[(f(\mathbf{x}_0) - y_0) - (f(\mathbf{x}_1) - y_1)]$$

$$\text{令：} \quad E_i = f(\mathbf{x}_i) - y_i \quad \eta = K_{00} + K_{11} - 2K_{01}$$

$$\text{则：} \quad \alpha_1^{\text{new}} = \alpha_1^{\text{old}} + \frac{y_1(E_0 - E_1)}{\eta}$$

利用上式计算得到 α_1^{new} 还需要检验它是否满足约束条件，因为它应该介于0到C之间，但对于某一次迭代中选出的特定的 α_0 和 α_1 ，还可以更精确地限定其取值范围。首先，新产生的 α_0^{new} 和 α_1^{new} 应该处于图7-6中 $[0, C] \times [0, C]$ 所限定的“盒子”内。如果 $y_0 \neq y_1$ （图7.6a），则有： $\alpha_0 - \alpha_1 = k$ (k 为常数)

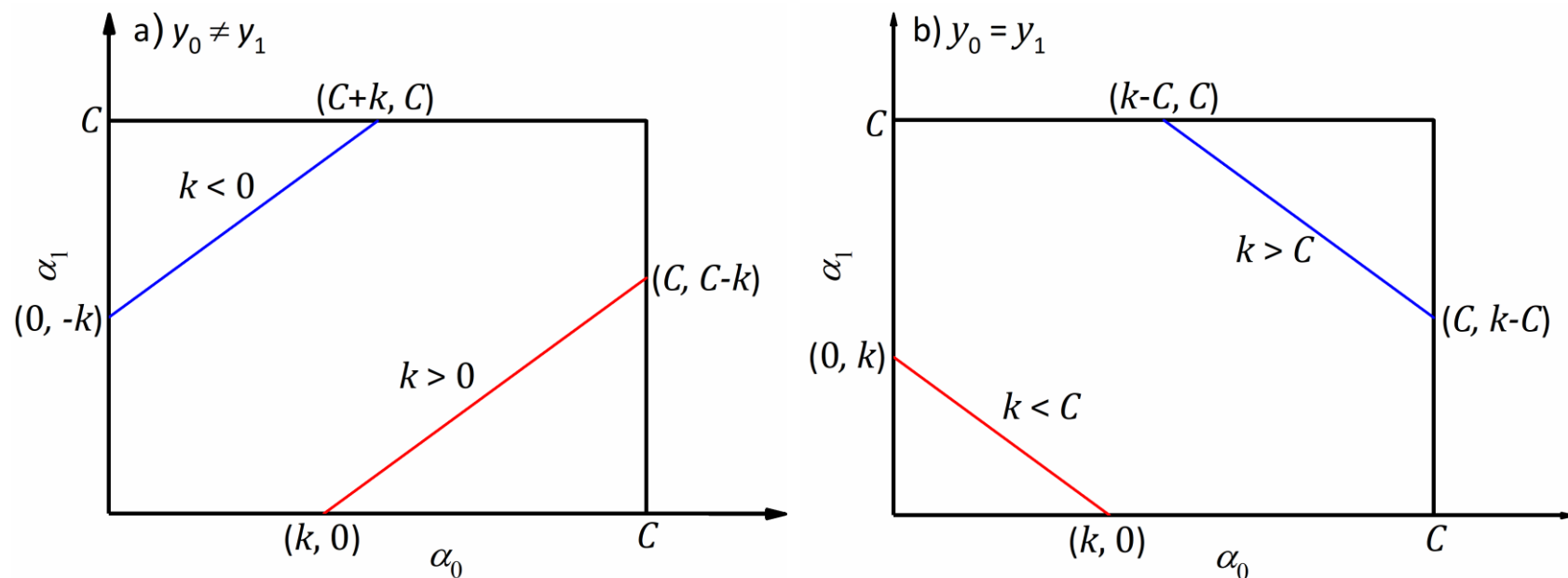
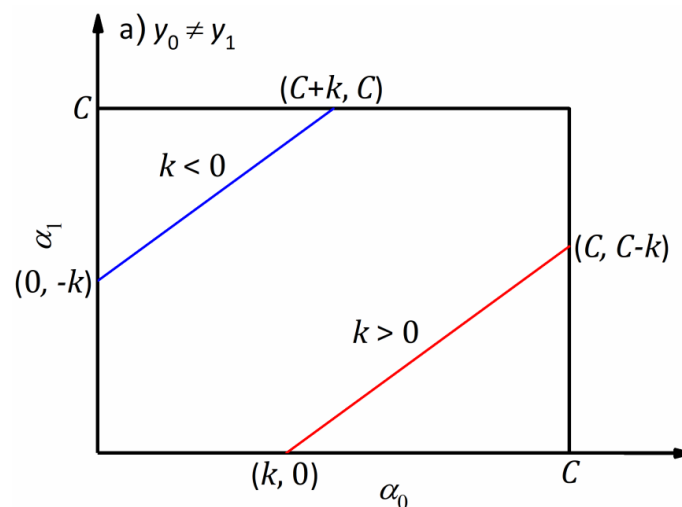


图7-6 拉格朗日乘子更新时的可行域

$$\alpha_0 - \alpha_1 = k \implies \alpha_1 = \alpha_0 - k$$

上式为一直线方程，当 $k > 0$ 时该直线位于对角线下侧，如图7-6a中红线所示，该直线与“盒子”的交点在 $(k, 0)$ 和 $(C, C-k)$ ，显然 α_1 只能在0到 $C-k$ 之间变化。



由于当前值 $(\alpha_0^{\text{old}}, \alpha_1^{\text{old}})$ 也在这条线上，即 $\alpha_0^{\text{old}} - \alpha_1^{\text{old}} = k$ ，得到 α_1 的上下限分别为：

$$L = 0, H = C - k = C + \alpha_1^{\text{old}} - \alpha_0^{\text{old}}$$

当 $k < 0$ 时该直线位于对角线上侧，如图7-6a中蓝线所示，该直线与“盒子”的交点在 $(0, -k)$ 和 $(C+k, C)$ ，此时 α_1 介于 $-k$ 到 C 之间变化， α_1 的上下限分别为：

$$L = -k = \alpha_1^{\text{old}} - \alpha_0^{\text{old}}, H = C$$

$$L = 0, H = C - k = C + \alpha_1^{\text{old}} - \alpha_0^{\text{old}}$$

$$L = -k = \alpha_1^{\text{old}} - \alpha_0^{\text{old}}, H = C$$

综合两种情况，得到 $y_0 \neq y_1$ 时 α_1 的上下限：

$$\begin{cases} L = \max(0, \alpha_1^{\text{old}} - \alpha_0^{\text{old}}) \\ H = \min(C, C + \alpha_1^{\text{old}} - \alpha_0^{\text{old}}) \end{cases}$$

类似地可以得到 $y_0 = y_1$ 时（参见图7-6b） α_1 的上下限：

$$\begin{cases} L = \max(0, \alpha_0^{\text{old}} + \alpha_1^{\text{old}} - C) \\ H = \min(C, \alpha_0^{\text{old}} + \alpha_1^{\text{old}}) \end{cases}$$

因此，计算得到的 α_1^{new} 需要根据其取值范围进行修剪，得到本轮优化的结果：

$$\begin{cases} \alpha_1^{\text{new,clip}} = H, \text{ while } \alpha_1^{\text{new}} \geq H \\ \alpha_1^{\text{new,clip}} = L, \text{ while } \alpha_1^{\text{new}} \leq L \\ \alpha_1^{\text{new,clip}} = \alpha_1^{\text{new}}, \text{ while } L < \alpha_1^{\text{new}} < H \end{cases}$$

然后根据 $\alpha_1^{\text{new,clip}}$ 再计算 α_0^{new} :

$$y_0 \alpha_0^{\text{new}} + y_1 \alpha_1^{\text{new,clip}} = y_0 \alpha_0^{\text{old}} + y_1 \alpha_1^{\text{old}} \implies$$

$$\alpha_0^{\text{new}} = \alpha_0^{\text{old}} + y_0 y_1 (\alpha_1^{\text{old}} - \alpha_1^{\text{new,clip}})$$

在每次完成两个 α 的优化之后，需要重新计算位移项 b 。当 $0 < \alpha_0^{\text{new}} < C$ 时，则样本 \mathbf{x}_0 为支持向量且落在分隔边界上，满足 $y_0 f(\mathbf{x}_0) - 1 = 0$ ，即：

$$y_0 = f(\mathbf{x}_0) = \sum_{i=0}^{n-1} \alpha_i y_i \mathbf{x}_i \mathbf{x}_0 + b = \sum_{i=0}^{n-1} \alpha_i y_i K_{0i} + b \implies$$

$$b_0^{\text{new}} = y_0 - \sum_{i=0}^{n-1} \alpha_i y_i K_{0i} = y_0 - \sum_{i=2}^{n-1} \alpha_i y_i K_{0i} - \alpha_0^{\text{new}} y_0 K_{00} - \alpha_1^{\text{new,clip}} y_1 K_{01}$$

上式的加和项只含有未改变的，可由当前值得出：

$$E_0 = f(\mathbf{x}_0) - y_0 = \sum_{i=2}^{n-1} \alpha_i y_i K_{0i} + \alpha_0^{\text{old}} y_0 K_{00} + \alpha_1^{\text{old}} y_1 K_{01} + b^{\text{old}} - y_0$$

$$b_0^{\text{new}} = y_0 - \sum_{i=0}^{n-1} \alpha_i y_i K_{0i} = y_0 - \sum_{i=2}^{n-1} \alpha_i y_i K_{0i} - \alpha_0^{\text{new}} y_0 K_{00} - \alpha_1^{\text{new,clip}} y_1 K_{01}$$

上面两式相加得到 b_0^{new} 的计算式：

$$b_0^{\text{new}} = -E_0 - y_0 K_{00} (\alpha_0^{\text{new}} - \alpha_0^{\text{old}}) - y_1 K_{01} (\alpha_1^{\text{new,clip}} - \alpha_1^{\text{old}}) + b^{\text{old}}$$

类似地，如果 $0 < \alpha_1^{\text{new,clip}} < C$ ，得到位移项更新公式，记为 b_1^{new} ：

$$b_1^{\text{new}} = -E_1 - y_0 K_{01} (\alpha_0^{\text{new}} - \alpha_0^{\text{old}}) - y_1 K_{11} (\alpha_1^{\text{new,clip}} - \alpha_1^{\text{old}}) + b^{\text{old}}$$

如果 $0 < \alpha_0^{\text{new}}, \alpha_1^{\text{new,clip}} < C$ ，即两个拉格朗日乘子都落在边界上，这时 b_0^{new} 与 b_1^{new} 是一致的，取前者即可。

如果 α_0^{new} 和 $\alpha_1^{\text{new,clip}}$ 都没有落在 $(0,C)$ 区间，仍然按上面两式计算 b_0^{new} 和 b_1^{new} ，取二者的均值作为新的位移项。因此位移项更新规则为：

$$b^{\text{new}} = \begin{cases} b_0^{\text{new}}, \text{ while } 0 < \alpha_0^{\text{new}} < C \\ b_1^{\text{new}}, \text{ while } 0 < \alpha_1^{\text{new,clip}} < C \\ \frac{b_0^{\text{new}} + b_1^{\text{new}}}{2}, \text{ otherwise} \end{cases}$$

最后，还有一个问题需要解决，即如何选择要优化的两个拉格朗日乘子。在SMO算法中利用一个外层循环选择第一个要优化的拉格朗日乘子（比如 α_1 ），在内层循环中选择第二个拉格朗日乘子（比如 α_0 ）并完成优化。外层循环中，先遍历所有数据，如果某一个 α 违反KKT条件则将它选为第一个待优化的拉格朗日乘子。

遍历一次以后，外循环中会继续遍历那些非边界处的 α 即 $0 < \alpha < C$ ，这是由于边界处的 α （ $\alpha=0$ 或 $\alpha=C$ ）在优化中的改动通常不大，选择非边界处 α 则可加速收敛过程。经过一系列的循环，当所有非边界 α 都满足KKT条件后，再对整个数据集进行遍历。这样在单次遍历整个数据集及多次遍历非边界 α 之间进行切换，当所有 α 都满足KKT条件后终止程序。

判断某个 α_i 是否满足KKT条件的标准包括：

（1）如果 $\alpha_i < C$ ，则对应样本为正确分类点且不会落在间隔内部，因而 $y_i f(x_i) - 1 \geq 0$ ，考虑到允许误差 tol ，如果 $y_i E_i = y_i f(x_i) - 1 < -\text{tol}$ 则视为违反KKT条件；

（2）如果 $\alpha_i > 0$ ，则对应样本为支持向量， $y_i f(x_i) - 1 = 0$ ，考虑到允许误差，如果 $y_i E_i = y_i f(x_i) - 1 > \text{tol}$ ，则视为违反KKT条件。

依据外循环中选出的第一个拉格朗日乘子 α_1 ，在内层循环中进行第二个拉格朗日乘子 α_0 的选择，选择原则是尽量让优化中产生的变化更大，从而“大步”走向最优点。根据更新式：

$$\alpha_1^{\text{new}} = \alpha_1^{\text{old}} + \frac{y_1(E_0 - E_1)}{\eta}$$

可以看出，优化后拉格朗日乘子的变化正比于两个样本的误差，即 $|E_0 - E_1|$ ，因此，应选择与 E_1 距离最大的 E 所对应的 α ，如果 $E_1 > 0$ 则选择最小的 E 为 E_0 ，如果 $E_1 < 0$ 则选择最大的 E 为 E_0 ，相应地可以确定 α_0 。

6. SMO算法python实现

下面利用python实现SMO算法，测试数据由scikit-learn中datasets.samples_generator模块中的make_blobs函数产生。

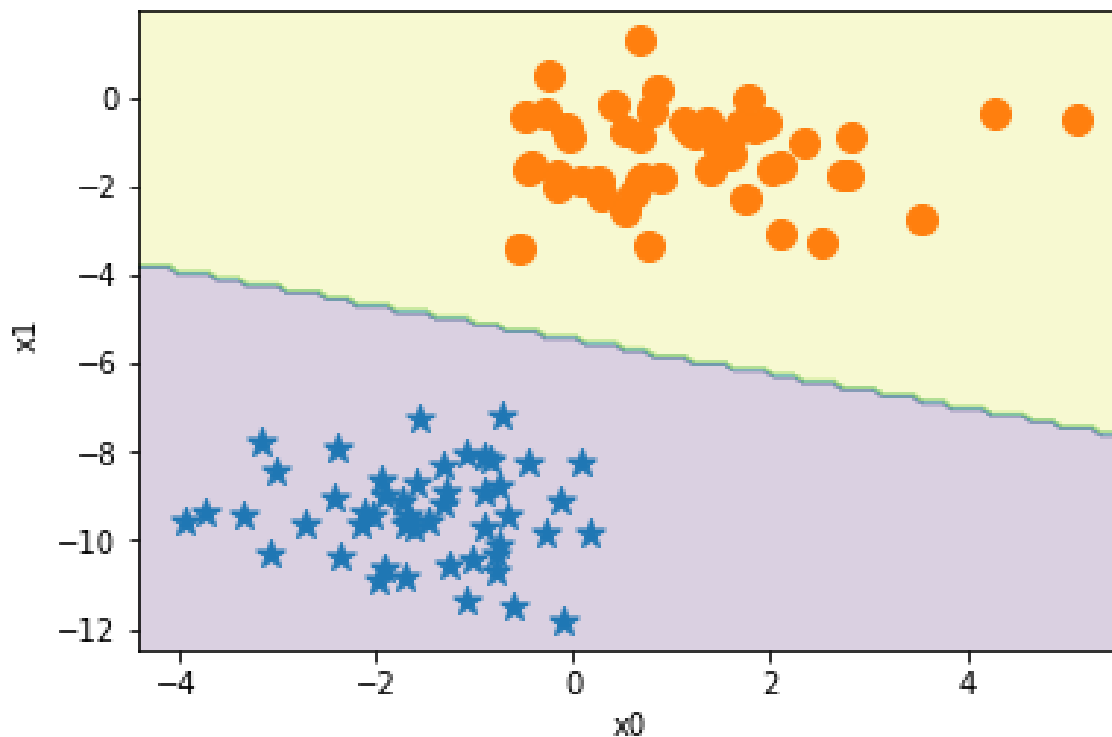


图7-7 SMO算法产生的分隔线

7. 核函数

前面的讨论都假设数据是线性可分的，然而在现实任务中也许并不存在一个能正确划分两类样本的超平面，如图7-8所示的两类数据，该数据由scikit-learn中的make_moons函数产生。

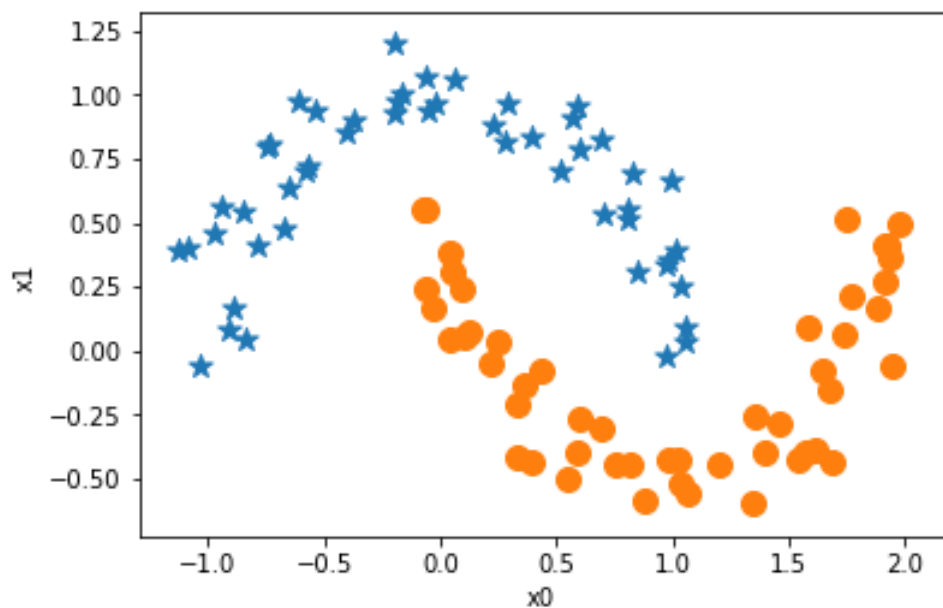


图7-8 用make_moons产生的数据集

图7-8中的数据，人眼很容易进行划分，只需要一条曲线就可以将所有样本正确分割。对这样的问题，可将样本从原始空间映射到一个更高维的空间，使得样本在新空间中线性可分。

幸运的是，如果原始样本空间是有限维，即特征数目有限，那么一定存在一个高维空间使样本可分。

令 $\phi(\mathbf{x})$ 表示将 \mathbf{x} 映射后的特征向量，于是，在新的特征空间中划分超平面所对应的模型可以表示为：

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

与前面的推导过程类似地，也可以得到其对偶问题为：

$$\max_{\alpha} \sum_{i=0}^{n-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

$$\text{s. t. } \sum_{i=0}^{n-1} \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C, i = 0, 1, \dots, n-1$$

上式与原始空间的优化问题几乎相同，唯一的区别是将 \mathbf{x}_i 与 \mathbf{x}_j 在原始空间中的内积 $\mathbf{x}_i^T \mathbf{x}_j$ 改变为映射空间中新特征向量的内积 $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ 。由于新的特征空间维数可能很高，甚至是无穷维的，因此直接计算 $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ 通常是困难的。为了避开这个障碍，人们采用核函数（kernel function）。

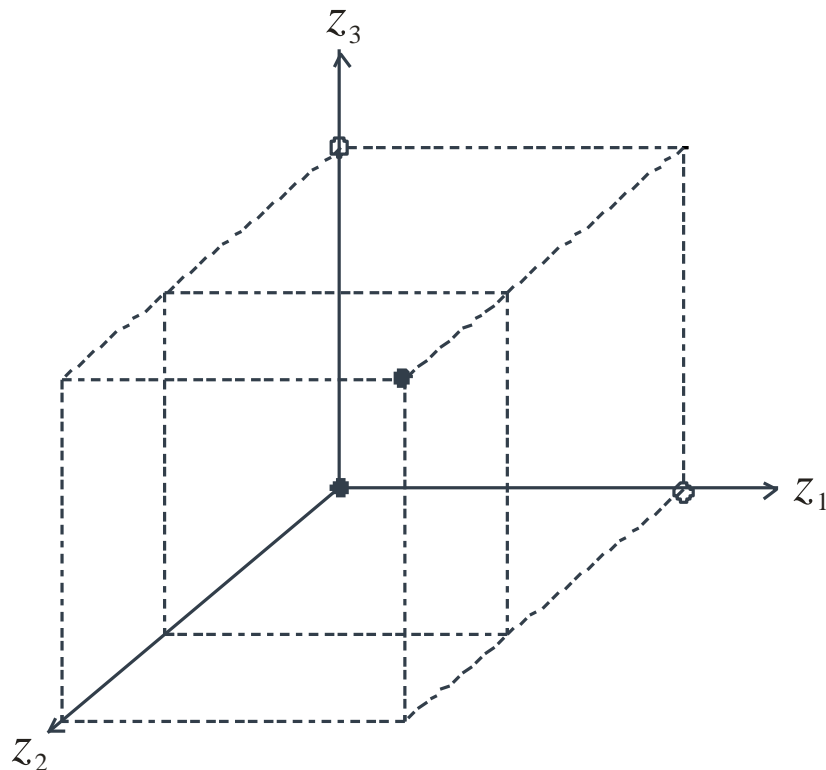
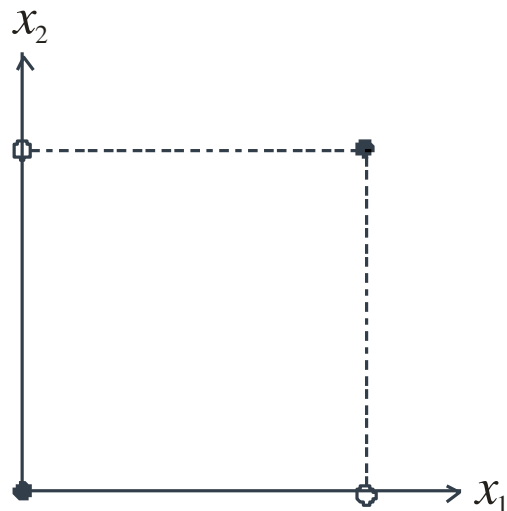
以二维空间中的异或问题为例，由于它在二维空间无法线性分隔，可以将其映射到更高维空间，比如三维空间，作如下映射：

$$\mathbf{x} = (x_1, x_2) \quad \Longrightarrow \quad \varphi(\mathbf{x}) = \mathbf{z} = (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

$$(0, 0) \quad \Longrightarrow \quad (0, 0, 0), \quad (0, 1) \quad \Longrightarrow \quad (0, 0, 1),$$

$$(1, 0) \quad \Longrightarrow \quad (1, 0, 0), \quad (1, 1) \quad \Longrightarrow \quad (1, 1.4, 1)$$

$$(z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



在新特征空间中，易于实现线性分割。比如 $(0.5, 0, 0)$, $(0, 0, 0.5)$, $(1, 1, 1)$ 这三个点构成的平面即可将异或问题正确分类。

在支持向量机的学习过程，会大量计算两个向量的内积，在新的多维特征空间中的内积最终都可以用原始空间的坐标来计算，比如上面的映射：

$$(z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

$$\begin{aligned}\mathbf{z} \cdot \mathbf{z}' &= z_1z_1' + z_2z_2' + z_3z_3' = x_1^2x_1'^2 + 2x_1x_1'x_2x_2' + x_2^2x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (\mathbf{x} \cdot \mathbf{x}')^2 = K(\mathbf{x}, \mathbf{x}')\end{aligned}$$

其中的 K 即为核函数，它是原始样本空间中的一个函数，它以两个样本（即向量）作为输入，其结果就是这两个样本映射后新特征向量的内积，即：

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

这样就把高维空间中两个向量的内积转换为原始样本空间中的一个函数，一旦有了这样的函数，就避免了在高维空间中进行计算。

利用核函数，优化问题转化为：

$$\max_{\alpha} \sum_{i=0}^{n-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

求解后得到：

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=0}^{n-1} \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b = \sum_{i=0}^{n-1} \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

显然，如果已知映射 ϕ 则容易写出其对应的核函数，但在现实任务中，构造合适的映射 ϕ 往往是很困难的，那如何得到其核函数呢？

我们有下面的定理：

核函数定理：令 χ 为输入空间， $K(\cdot, \cdot)$ 是定义在 $\chi \times \chi$ 上的对称函数， K 是核函数的充分必要条件是对于任意数据集 $D=\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{n-1}\}$ ，核矩阵（kernel matrix） \mathbf{K} 总是半正定的：

$$\mathbf{K} = \begin{bmatrix} K(\mathbf{x}_0, \mathbf{x}_0) & \cdots & K(\mathbf{x}_0, \mathbf{x}_{n-1}) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_{n-1}, \mathbf{x}_0) & \cdots & K(\mathbf{x}_{n-1}, \mathbf{x}_{n-1}) \end{bmatrix}$$

核函数定理表明，只要一个对称函数所对应的核矩阵是半正定的，它就能作为核函数使用，或者说，就一定存在一个与它对应的映射 ϕ ，因此，每一个核函数都隐式地定义了一个映射。既然如此，只要找到合适的核函数，也就不必关心相应的映射是什么样子的了。但是，我们建立映射的目的是为了让样本在新的特征空间中线性可分，既然不知道映射的具体形式，也就无法确切知道映射后是否是线性可分的。

通常的做法是可以比较几个核函数，看哪一个性能最佳。因而，核函数的选择成为支持向量机的一大变数，若核函数选择不合适，必然导致分类器的性能不佳。

表7-1 常用核函数

名称	表达式	参数
线性核	$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$	
多项式核	$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^d$	$d \geq 1$ 为多项式次数, $d=1$ 即为线性核
高斯核	$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}\right)$	$\sigma > 0$ 高斯核带宽
拉普拉斯核	$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ }{\sigma}\right)$	$\sigma > 0$
Sigmoid核	$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta \mathbf{x}_i^T \mathbf{x}_j + \theta)$	$\beta > 0, \theta < 0$

在非线性核中，应用最多的是高斯核函数，也称为径向基函数（Radial basis function，简称RBF）：

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right) = \exp\left(-\gamma\|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$$

其中 $\gamma=1/(2\sigma^2)$ ，是一个由用户设定的参数，它与方差成反比。

利用 python 实现 RBF 支持向量机分类图 7-8 的数据集。程序输出结果如图 7-9 所示。

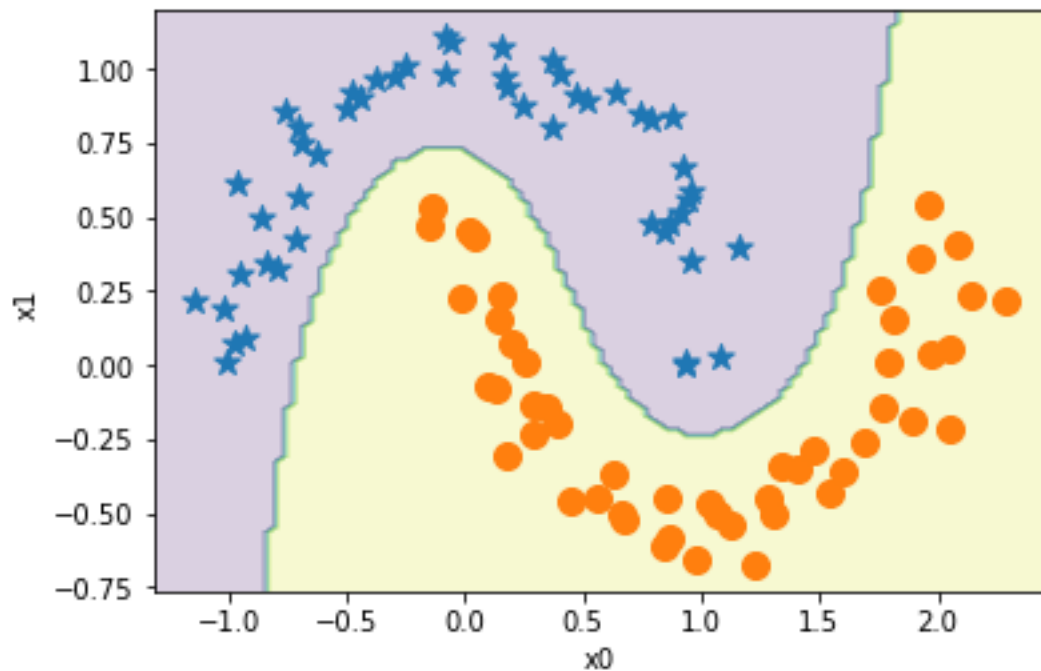


图7-9 RBF支持向量机产生的分隔线

8. 利用scikit-learn实现支持向量机

在scikit-learn中的svm模块实现了多个支持向量机相关的类：

- LinearSVC用于线性支持向量分类，
- LinearSVR用于线性支持向量回归，
- NuSVC用于核函数支持向量分类，
- NuSVR用于核函数支持向量回归。

LinearSVC的用法如下：

```
LinearSVC(penalty='l2', loss='squared_hinge', dual=True,  
tol=0.0001, C=1.0, multi_class='ovr', fit_intercept=True,  
intercept_scaling=1, class_weight=None, verbose=0,  
random_state=None, max_iter=1000)
```

利用LinearSVC类对图7-7中的线性分类问题进行求解。

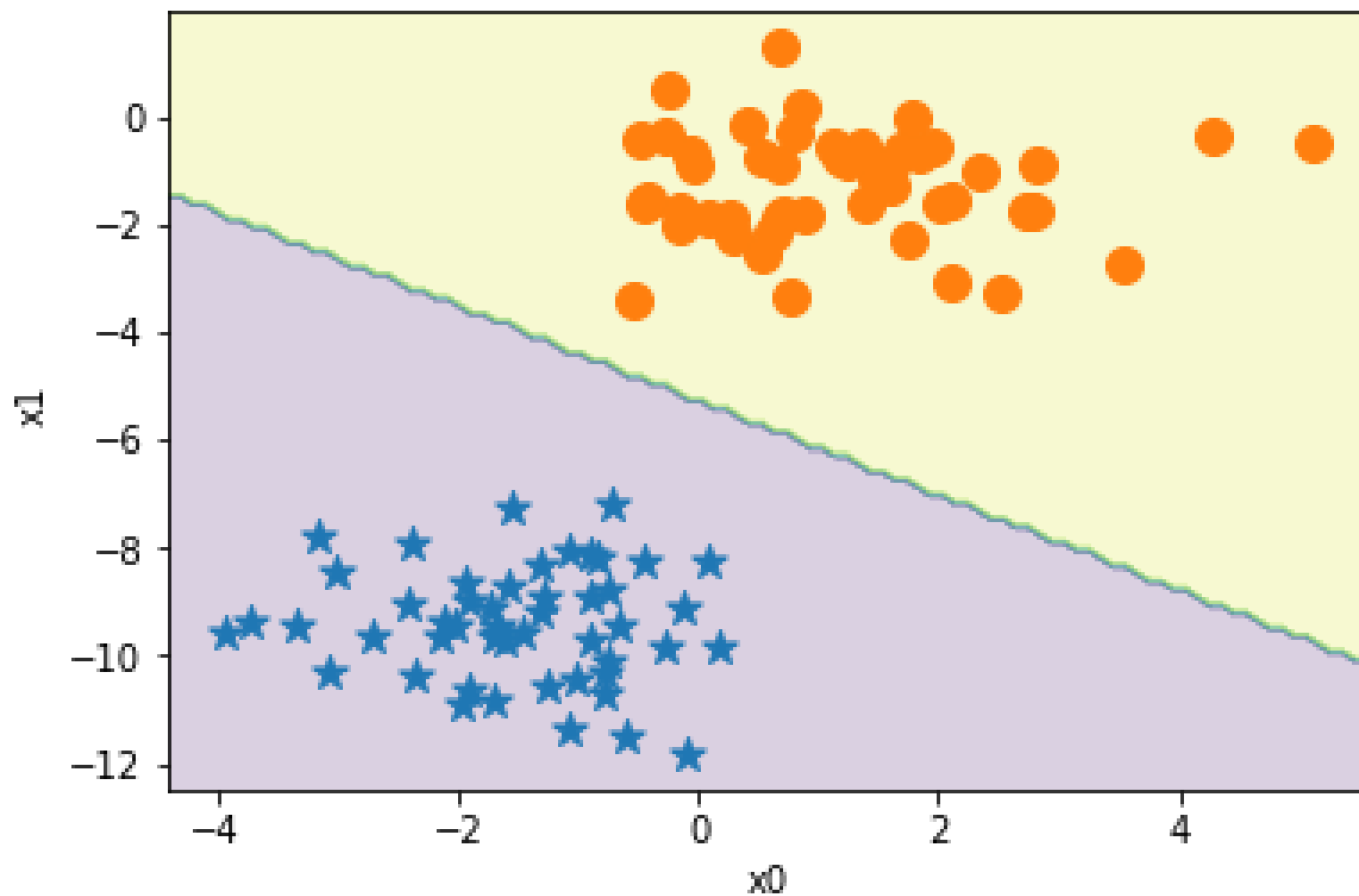


图7-10 利用LinearSVC分类的分隔线

NuSVC类用法如下：

```
NuSVC(nu=0.5, kernel='rbf', degree=3,  
gamma='auto_deprecated', coef0=0.0, shrinking=True,  
probability=False, tol=0.001, cache_size=200,  
class_weight=None, verbose=False, max_iter=-1,  
decision_function_shape='ovr', random_state=None)
```

利用NuSVC类来对make_moons产生的数据进行分类。

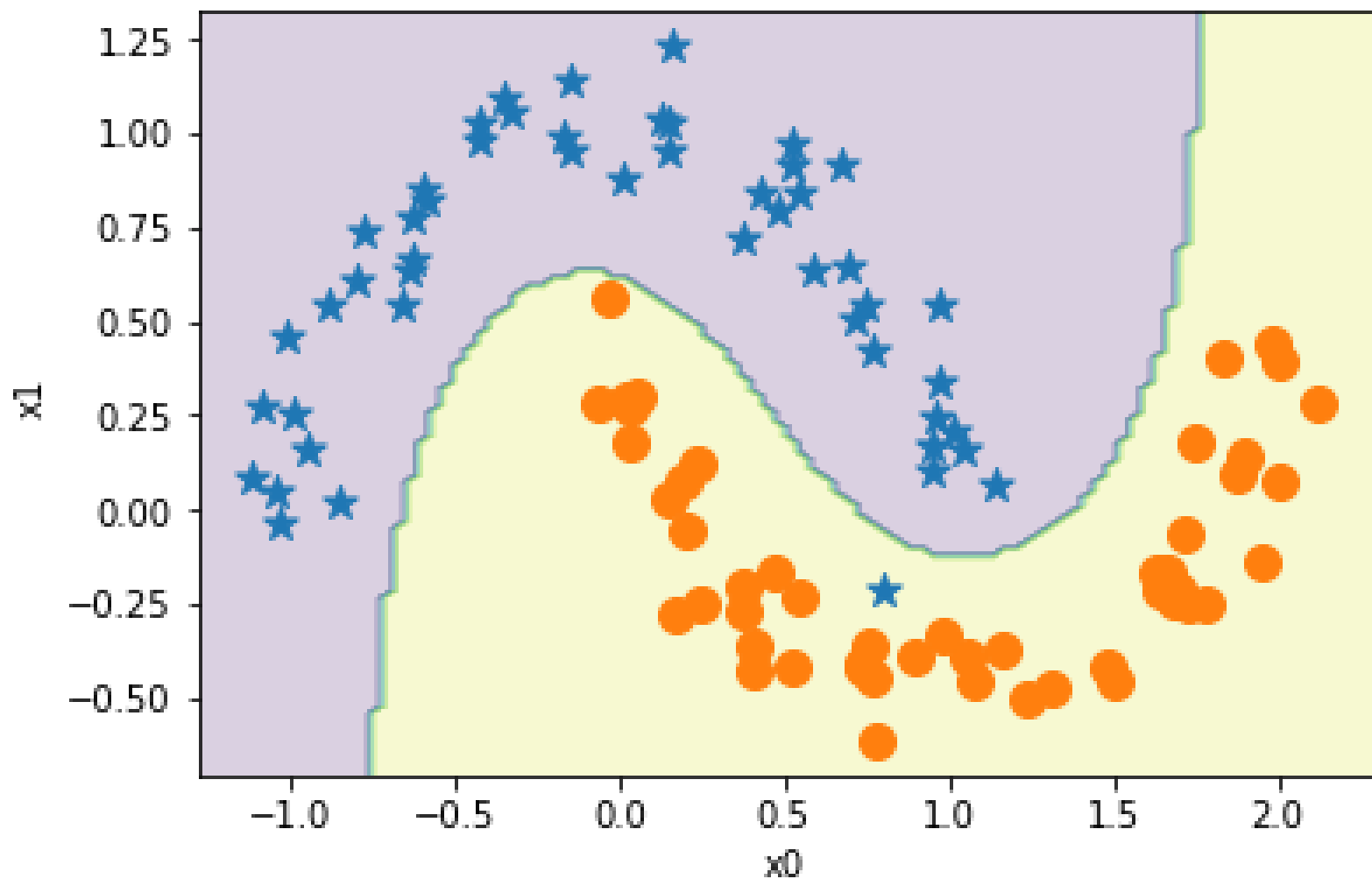


图7-11 利用NuSVC类对make_moons数据分类结果