

给定一个训练集，我们的目标是找到一个决策边界，使我们能够对训练示例做出所有正确和自信（意味着远离决策边界）的预测。

- Given a point (x_0, y_0) , the distance from the point to the line $Ax + By + C = 0$:

$$distance = \frac{|Ax_0 + By_0 + C|}{\sqrt{A^2 + B^2}}$$

- Given a point \mathbf{x}_i , the distance from the point to the hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$:

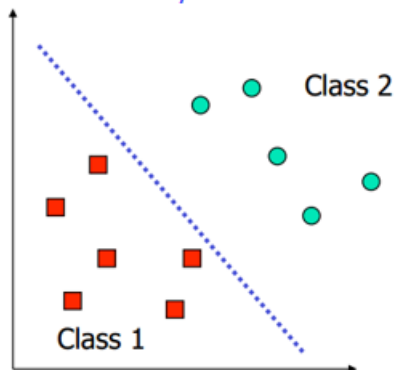
$$distance = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|}$$

\mathbf{w} 是和直线垂直的向量

What Is a Good Decision Boundary?

- We aim to find the hyperplane (i.e., decision boundary) linearly separating our classes. 我们的目标是找到线性分离类的超平面（即决策边界）。
- Our boundary will have equation: $\mathbf{w}^T \mathbf{x} + b = 0$

Decision boundary



- Above the decision boundary should have label 1, i.e., for any \mathbf{x}_i s.t. $\mathbf{w}^T \mathbf{x} + b > 0$, then $y_i = 1$.

- Below the decision boundary should have label -1, i.e., for any \mathbf{x}_i s.t. $\mathbf{w}^T \mathbf{x} + b < 0$, then $y_i = -1$.

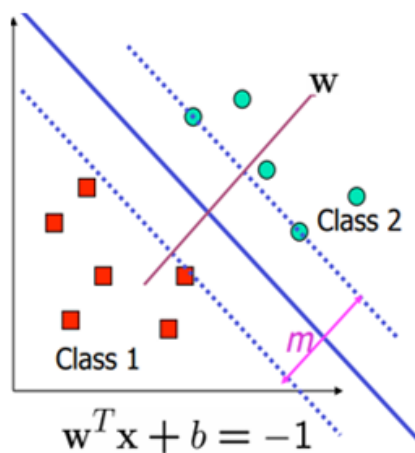
$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

- Moreover, we hope the hyperplane lies in the middle

$$\begin{cases} (\mathbf{w}^T \mathbf{x}_i + b) / \|\mathbf{w}\| \geq \frac{m}{2} & \forall y_i = 1 \\ (\mathbf{w}^T \mathbf{x}_i + b) / \|\mathbf{w}\| \leq -\frac{m}{2} & \forall y_i = -1 \end{cases}$$

$$\text{distance} = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|}$$

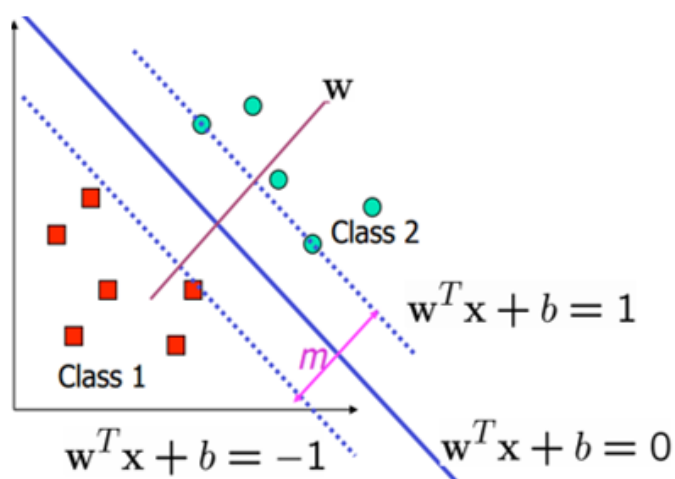
m is the margin



m is the margin, 两类之间最近的距离

- Therefore,

$$\begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 1 & \forall y_i = 1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq -1 & \forall y_i = -1 \end{cases} \quad \Rightarrow \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$



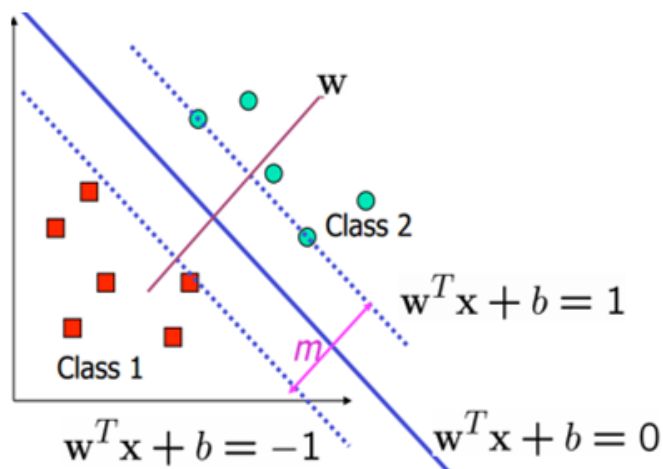
决策边缘应尽可能远离这两个类的数据，应该取 m 的最大值

- For the **support vectors** (data points nearest to the hyperplane)

$$\text{Distance} = |w^T x_i + b| / \|w\|$$

$$= 1 / \|w\|$$

$$m = 2 / \|w\|$$



以上是一个具有凸二次目标和仅线性约束的优化问题。

Exercise

$$\min_w \frac{1}{2} \|w\|^2 = \frac{1}{2} w^T w$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1, i = 1, 2, \dots, n$$

- Given the dataset consist of two positive samples $x_1 = (3, 3)^T$, $x_2 = (4, 3)^T$, and one negative sample $x_3 = (1, 1)^T$. Please write the objective function with SVM.

Answer

$$\min_{w,b} \frac{1}{2} \|w\|^2 = \frac{1}{2} w_1^2 + \frac{1}{2} w_2^2$$

$$\text{s.t. } 3w_1 + 3w_2 + b \geq 1,$$

$$4w_1 + 3w_2 + b \geq 1,$$

$$-w_1 - w_2 - b \geq 1,$$

拉格朗日对偶性

对偶形式将使我们能够推导出一种有效的算法来解决优化问题。

对偶形式将允许我们使用核来获得最优边缘分类器，以便在非常高维的空间中有效地工作。

Constrained Optimization

Consider a problem of the following form:

$$\begin{aligned} \min_{\mathbf{w}} f(\mathbf{w}) \\ \text{s.t. } h_i(\mathbf{w}) = 0, i = 1, \dots, l. \end{aligned}$$

Lagrange multiplier method:

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

β_i 's are the Lagrange multipliers.
No constraint now.

Set the partial derivatives to zero:

$$\frac{\partial \mathcal{L}(\mathbf{w}, \boldsymbol{\beta})}{\partial w_j} = 0 \quad \frac{\partial \mathcal{L}(\mathbf{w}, \boldsymbol{\beta})}{\partial \beta_i} = 0$$

不等式约束优化

将其推广到约束优化问题，在这些问题中，我们可能存在不等式和等式约束。

Consider the following **primal** optimization problem:

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

Generalized Lagrangian

α_i 's and β_i 's are the Lagrange multipliers.

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

$$\alpha_i \geq 0$$

Why?

Optimization with Inequality Constraints

Consider the following primal optimization problem:

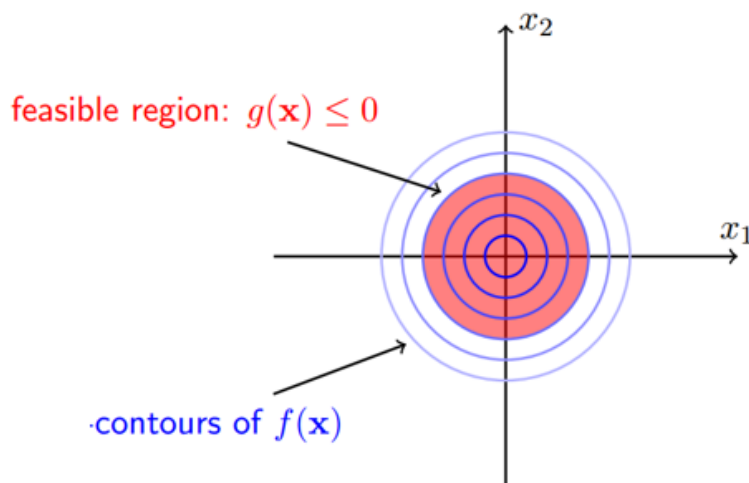
$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) \\ \text{s.t. } g(\mathbf{x}) \leq 0 \end{aligned}$$

Example1:

$$f(\mathbf{x}) = x_1^2 + x_2^2 \text{ and } g(\mathbf{x}) = x_1^2 + x_2^2 - 1$$

Optimization with Inequality Constraints

$$f(\mathbf{x}) = x_1^2 + x_2^2 \text{ and } g(\mathbf{x}) = x_1^2 + x_2^2 - 1$$



$$g(\mathbf{x}) = x_1^2 + x_2^2 - 1$$

Optimization with Inequality Constraints

Problem:

Our constrained optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) \text{ subject to } g(\mathbf{x}) \leq 0$$

where

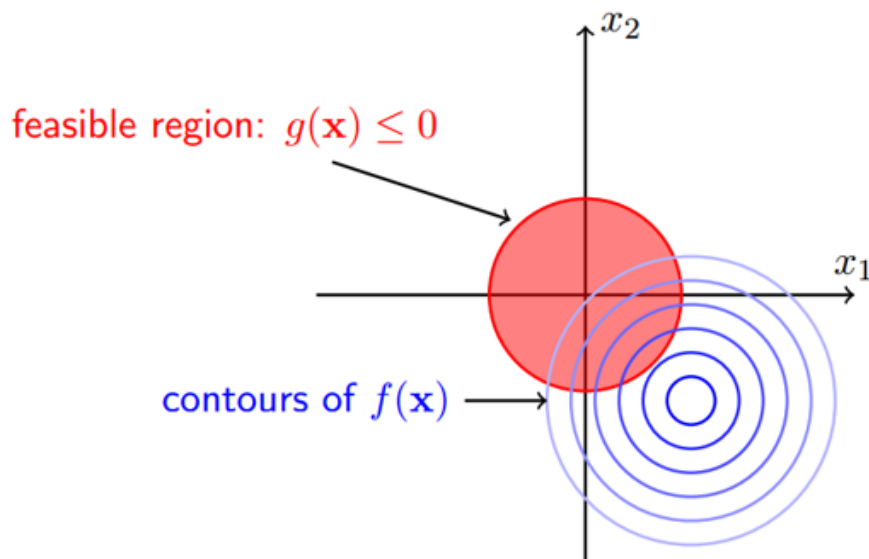
$$f(\mathbf{x}) = x_1^2 + x_2^2 \text{ and } g(\mathbf{x}) = x_1^2 + x_2^2 - 1$$

Constraint is not active at the local minimum ($g(\mathbf{x}^*) < 0$):

Therefore the local minimum is identified by the same conditions as in the unconstrained case. 有约束和无约束条件下的解一样

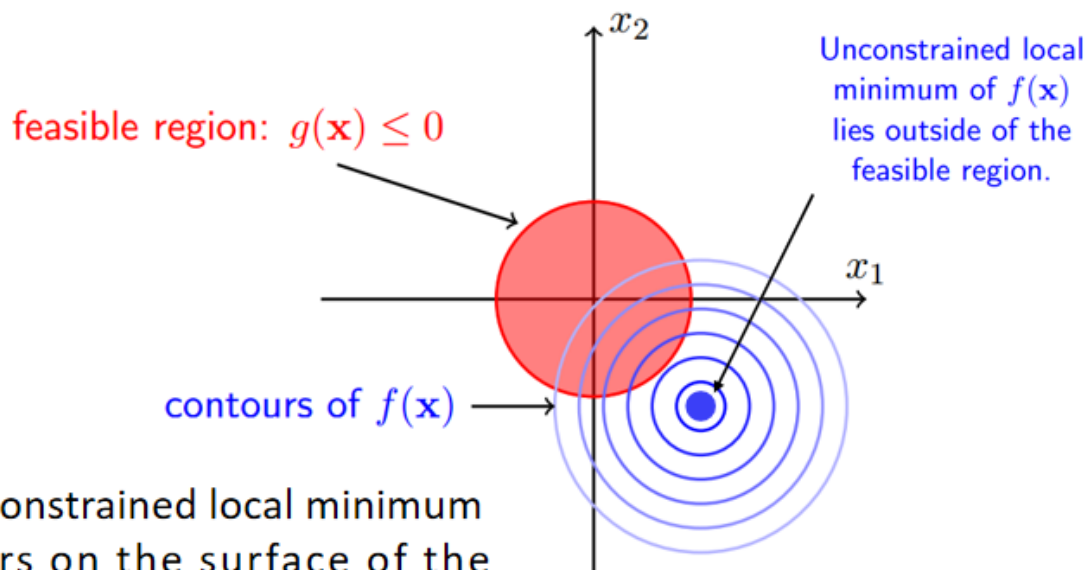
Optimization with inequality constraints

$$f(\mathbf{x}) = (x_1 - 1.1)^2 + (x_2 + 1.1)^2 \text{ and } g(\mathbf{x}) = x_1^2 + x_2^2 - 1$$



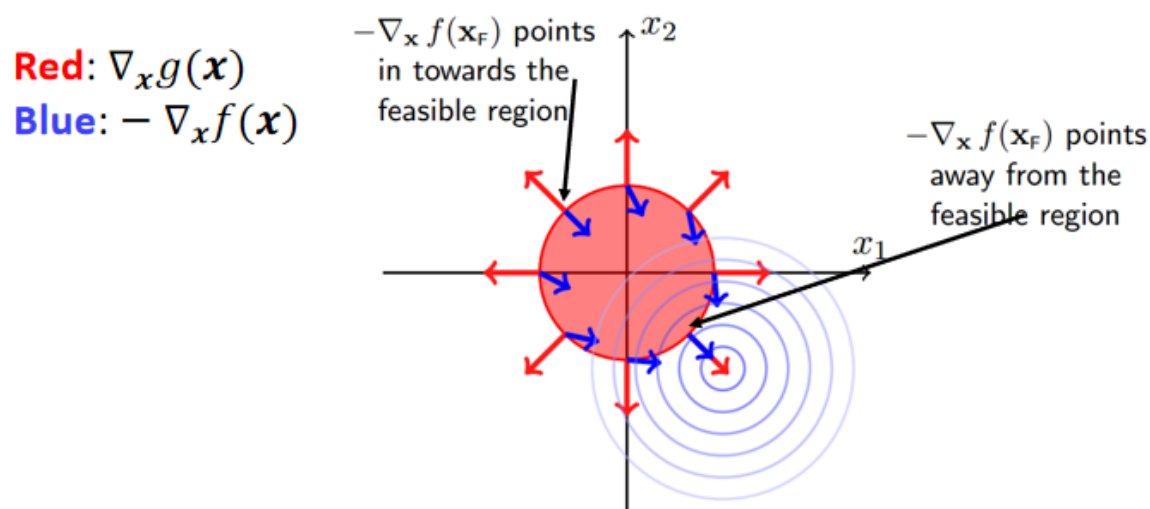
$$g(\mathbf{x}) = x_1^2 + x_2^2 - 1$$

约束条件在 $g(\mathbf{x}) = 0$ 时起作用



The constrained local minimum occurs on the surface of the constraint surface. $g(\mathbf{x}) = 0$

$$g(\mathbf{x}) = x_1^2 + x_2^2 - 1$$



\therefore Constrained local minimum occurs when $-\nabla_{\mathbf{x}} f(\mathbf{x})$ and $\nabla_{\mathbf{x}} g(\mathbf{x})$ point in the same direction:

$$-\nabla_{\mathbf{x}} f(\mathbf{x}) = \lambda \nabla_{\mathbf{x}} g(\mathbf{x}) \quad \text{and} \quad \lambda > 0$$

当 $f(x)$ 与 $g(x)$ 的导数不同号时，取得局部最小值

If \mathbf{x}^* corresponds to a constrained local minimum then

Case 1:

Unconstrained local minimum occurs **in** the feasible region.

- ① $g(\mathbf{x}^*) < 0$ 不起约束作用
- ② $\nabla_{\mathbf{x}} f(\mathbf{x}^*) = \mathbf{0}$

Case 2:

Unconstrained local minimum lies **outside** the feasible region.

- ① $g(\mathbf{x}^*) = 0$
- ② $-\nabla_{\mathbf{x}} f(\mathbf{x}^*) = \lambda \nabla_{\mathbf{x}} g(\mathbf{x}^*)$
with $\lambda > 0$

起约束作用，在边界

Dual optimization problem

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

坐标上升方法

Coordinate Ascent 坐标上升法

- Consider trying to solve the **unconstrained** optimization problem

$$\max_{\alpha} L(\alpha_1, \alpha_2, \dots, \alpha_l)$$

- Coordinate Ascent

Loop until convergence:{

For $i = 1, \dots, l$ {

$$\alpha_i := \operatorname{argmax}_{\hat{\alpha}_i} L(\alpha_1, \dots, \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, \dots, \alpha_l)$$

}

}

在该算法的最内循环中，我们将保持除某些 α_i 之外的所有变量固定，并仅针对参数 α_i 重新优化L。

SMO最小序列方法

每次只优化一个参数，其他参数先固定住，仅求当前这个优化参数的极值。我们来看一下 SMO 算法在 SVM 中的应用。

假设我们有一组满足约束的 α_i

假设我们固定 $\alpha_2, \dots, \alpha_n$ ，我们可以采取坐标上升步骤并优化关于 α_1 的函数吗？

• **NO!!!**

$$\sum_{i=1}^n \alpha_i y_i = 0 \qquad \alpha_1 = -y_1 \sum_{i=2}^n \alpha_i y_i$$

我们必须**同时更新至少两个 α_i** 。

重复，直到收敛{

选择一些对 α_i 和 α_j 进行下一次更新（使用启发式方式，尝试选择两个，使我们能够朝着全局最大值取得最大进展）。

关于 α_i 和 α_j 重新优化L (α)，同时保持所有其他 α_k ($k \neq i, j$) 固定

}

SMO是有效的，因为可以非常有效地计算对 α_i 和 α_j 的更新。

支持向量及其性质

$$\alpha_1 y_1 + \alpha_2 y_2 = - \sum_{i=3}^n \alpha_i y_i = \zeta \quad \text{Constant, 常数}$$

$$L(\alpha_1, \alpha_2, \dots, \alpha_n) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$L(\alpha_1, \alpha_2, \dots, \alpha_n) = L(y_1(\zeta - \alpha_2 y_2), \alpha_2, \dots, \alpha_n)$$

- 这是一个 α_2 的二次函数
- Once we have α_2^{new} , we can obtain α_1^{new} with $\alpha_1 y_1 + \alpha_2 y_2 = \zeta$

不能没有支持向量

- Question: can we have no support vector? $\alpha^* = \mathbf{0}$

- Answer: No.
- If $\alpha^* = \mathbf{0}$, then $\mathbf{w}^* = \mathbf{0}$. (This is not the optimal solution for the primal optimization problem)这不是原始优化问题的最优解

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

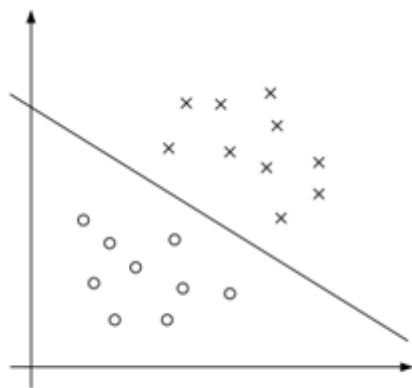
$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, n$$

大多数的 α 是0.因为支持向量是少数

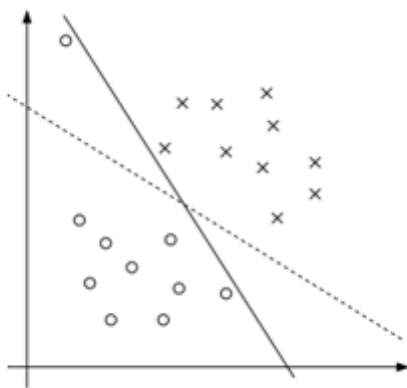
正则化与不可分情形

在某些情况下（由于异常值），我们不清楚找到一个分离超平面是否正是我们想要做的

图（a）显示了一个最优裕度分类器，当在左上角区域添加单个异常值时（图b），它会导致决策边界发生剧烈波动，结果分类器的裕度要小得多（对异常值敏感）。

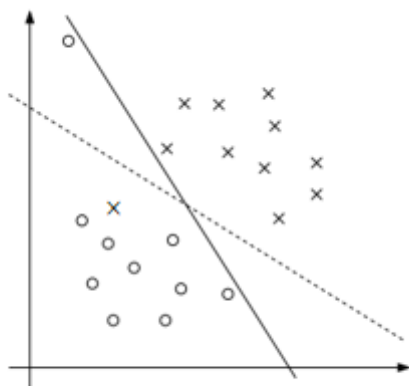


(a) Linearly separable



(b) Linearly separable
with outliers

在某些情况下（图c），数据不能完全线性分离。



(c) Non-linearly separable

软间隔

正松弛变量

如果一个例子的裕度为 $1 - \xi_i$ ($\xi_i > 0$)，我们将付出目标函数增加 $C\xi_i$ 的代价。

C 控制两个目标之间的相对权重

使“ w ”变小（使裕度变大）；

确保大多数示例的裕度至少为1

非线性

我们不想使用原始输入空间 x 来应用SVM，而是想使用一些特征空间 $\phi(x)$ 来学习

要做到这一点，我们只需要回顾我们以前的SVM算法，并将其中所有的 x 都替换为 $\phi(x)$ 。

Train SVM:

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

s.t. $0 \leq \alpha_i \leq C, i = 1, \dots, n$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

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

Test SVM:

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = \sum_{i \in S} \alpha_i y_i \mathbf{x}_i$$

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

$$= \text{sign}\left(\left(\sum_{i \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}\right) + b\right)$$

$\phi(\mathbf{x}_i)^T \phi(\mathbf{x})$

能计算在特征空间下的 $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$, 就不需要求显式的 $\phi(\mathbf{x}_i)$

核函数

- Define the **kernel** function K by

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

Train SVM:

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

s.t. $0 \leq \alpha_i \leq C, i = 1, \dots, n$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

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

Test SVM:

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = \sum_{i \in S} \alpha_i y_i \mathbf{x}_i$$

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

$$= \text{sgn}\left(\left(\sum_{i \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}\right) + b\right)$$

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

Intuition:

- If $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$ are close, we might want $K(\mathbf{x}_i, \mathbf{x}_j)$ to be large.
- If $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$ are far apart, we might want $K(\mathbf{x}_i, \mathbf{x}_j)$ to be small.

我们可以把 $K(x_i, x_j)$ 看作是对 $\xi(x_i)$ 和 $\xi(x_j)$ 相似程度的测量。

高斯核

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) \Leftrightarrow K(x_i, x_j) = \exp\left(-\gamma\|x_i - x_j\|^2\right), \gamma > 0$$

γ 越大，越容易过拟合

核函数的判定

Given the training dataset $\{x_1, x_2, \dots, x_n\}$.

Let $K_{ij} = K(x_i, x_j)$ be the (i, j) -entry of $K \in \mathbb{R}^{n \times n}$.

K is called the **Kernel matrix**.

If K is a valid kernel, then $K_{ij} = \phi(x_i)^T \phi(x_j) = \phi(x_j)^T \phi(x_i) = K_{ji}$.

K is **symmetric**. K is **positive semi-definite**.

要满足 $K_{ij}=K_{ji}$, K 是对称矩阵, 是半正定的 (所有的特征值都是非负的)

SVM的评价

优点:

没有局部解

它对高维数据的扩展性相对较好

分类器复杂度和误差之间的权衡可以明确控制

字符串和树等非传统数据可以用作SVM的输入, 而不是特征向量

缺点: 需要选择一个好的核函数

SVM回归

特征空间中的线性回归

与最小二乘回归不同, 误差函数为 ϵ -不敏感损失函数, 直觉上, 小于 ϵ 的错误被忽略

