

# 第4章：支持向量机学习

蒋良孝



中国地质大学（武汉）



CUG-Miner

机器学习与数据挖掘团队

**ljiaang@cug.edu.cn**

<http://www.escience.cn/people/jlx/>



# 本章内容

一、最大边缘超平面

二、线性支持向量机

三、非线性支持向量机

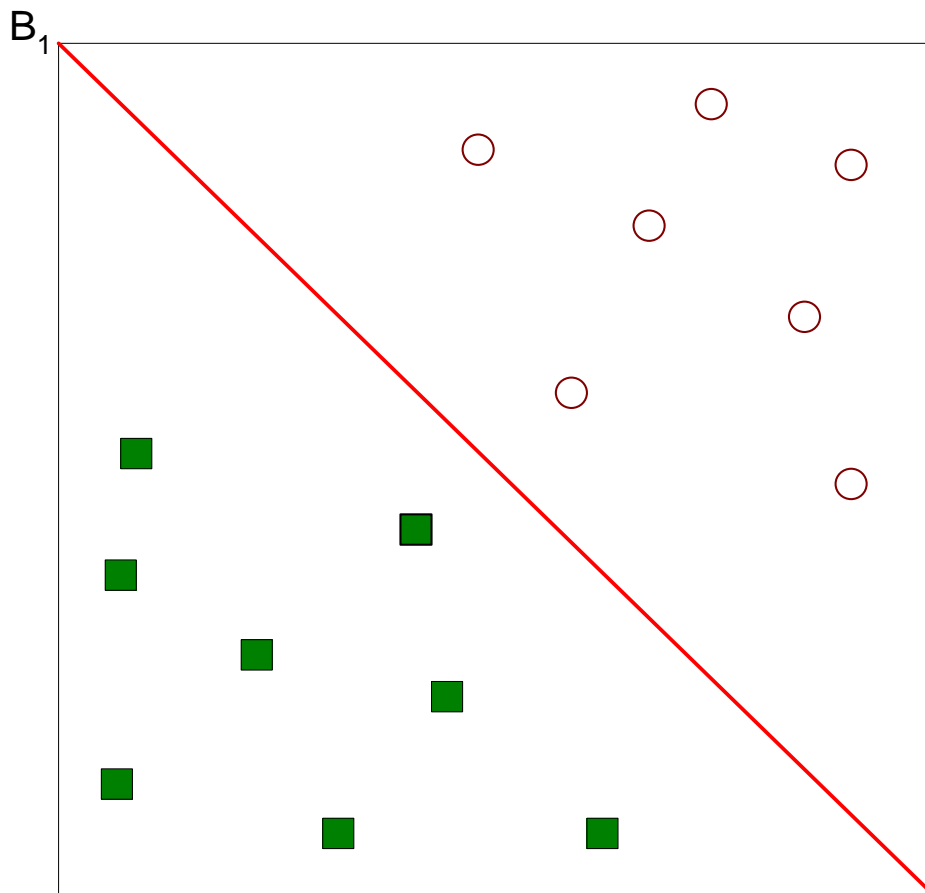
# 一、最大边缘超平面

---

- 支持向量机(Support Vector Machine, SVM)不仅具有坚实的统计学理论基础，还可以很好地应用于高维数据、避免维度灾难问题，已经成为一种倍受关注的机器学习分类技术。
- 为了解释SVM的基本思想，我们首先介绍一下最大边缘超平面(Maximal Margin Hyperplane)。

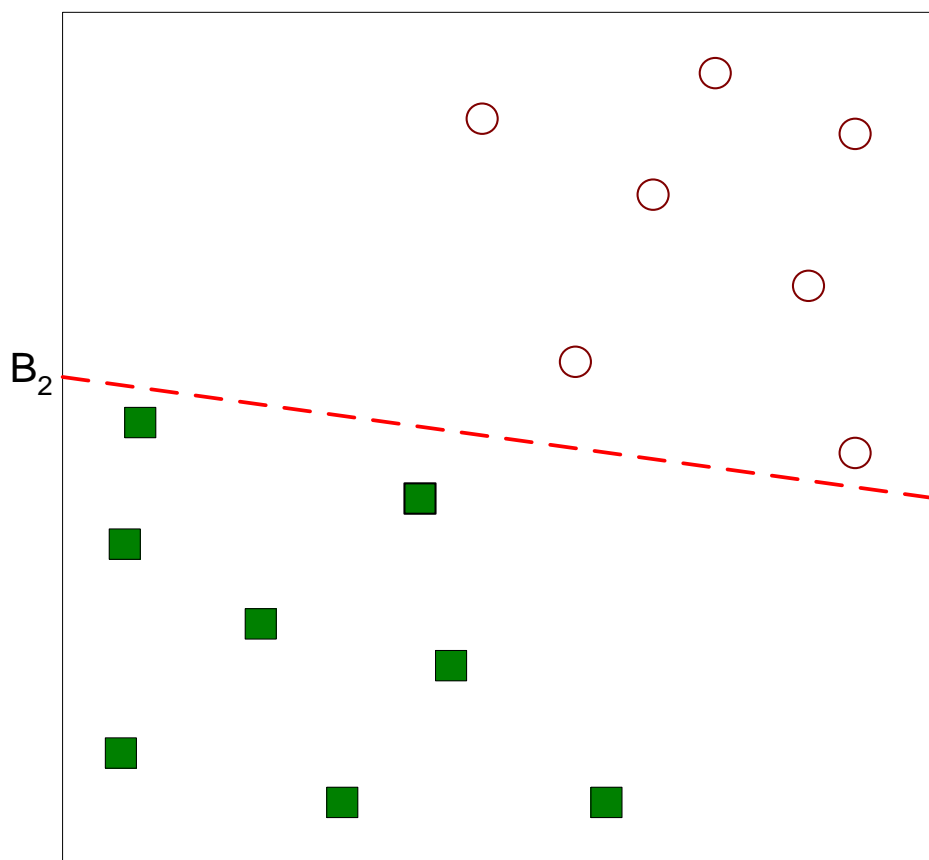
# 一、最大边缘超平面

给定训练数据集，线性分类器最基本的想法是：在样本空间中寻找一个超平面，将不同类别的样本分开。



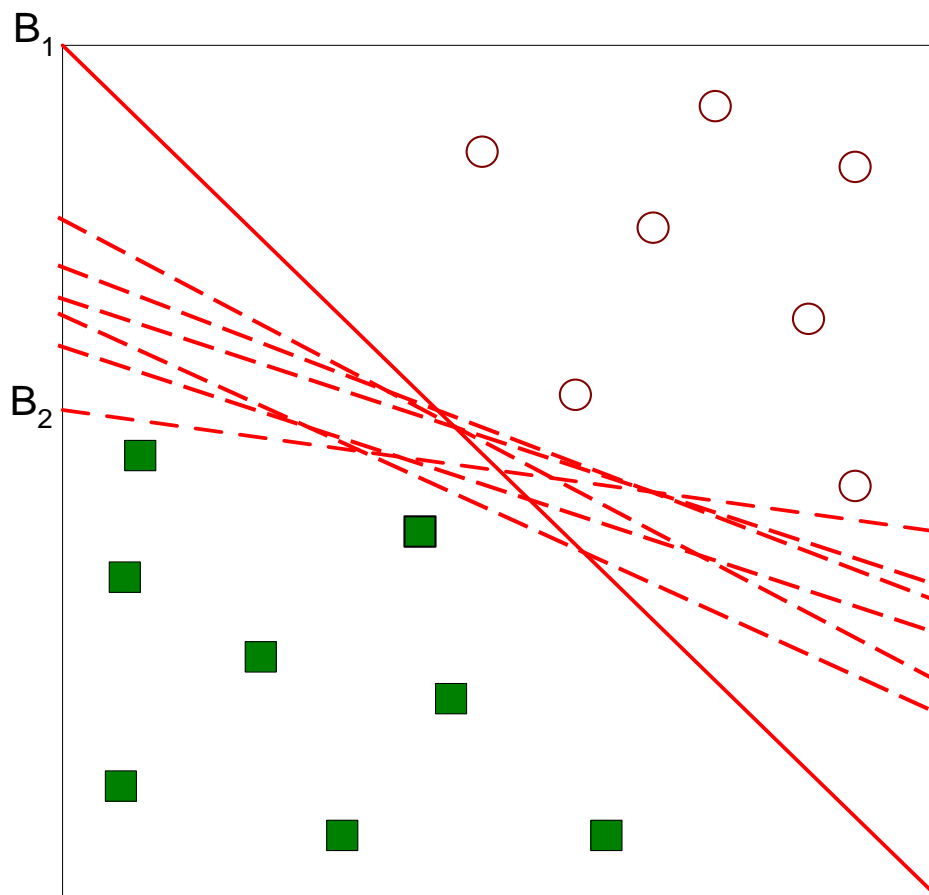
# 一、最大边缘超平面

给定训练数据集，线性分类器最基本的想法是：在样本空间中寻找一个超平面，将不同类别的样本分开。



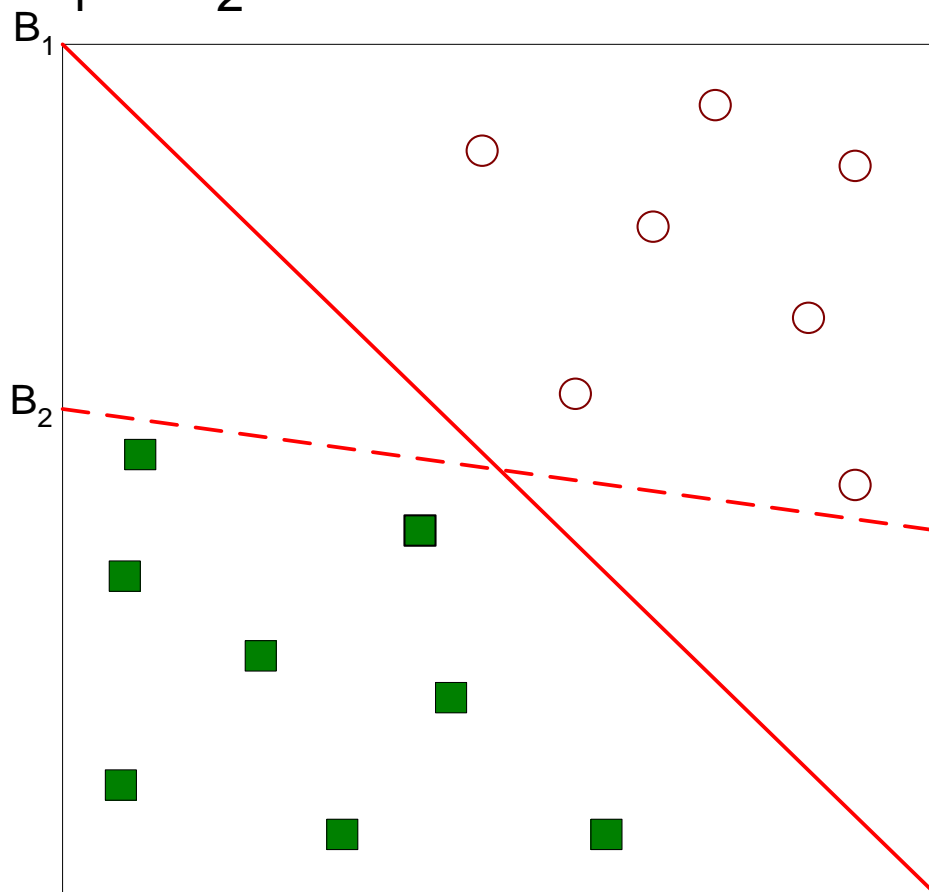
# 一、最大边缘超平面

其实能将训练样本分开的超平面可能有很多，分类器必须从这些超平面中选择哪个来表示它的决策边界？



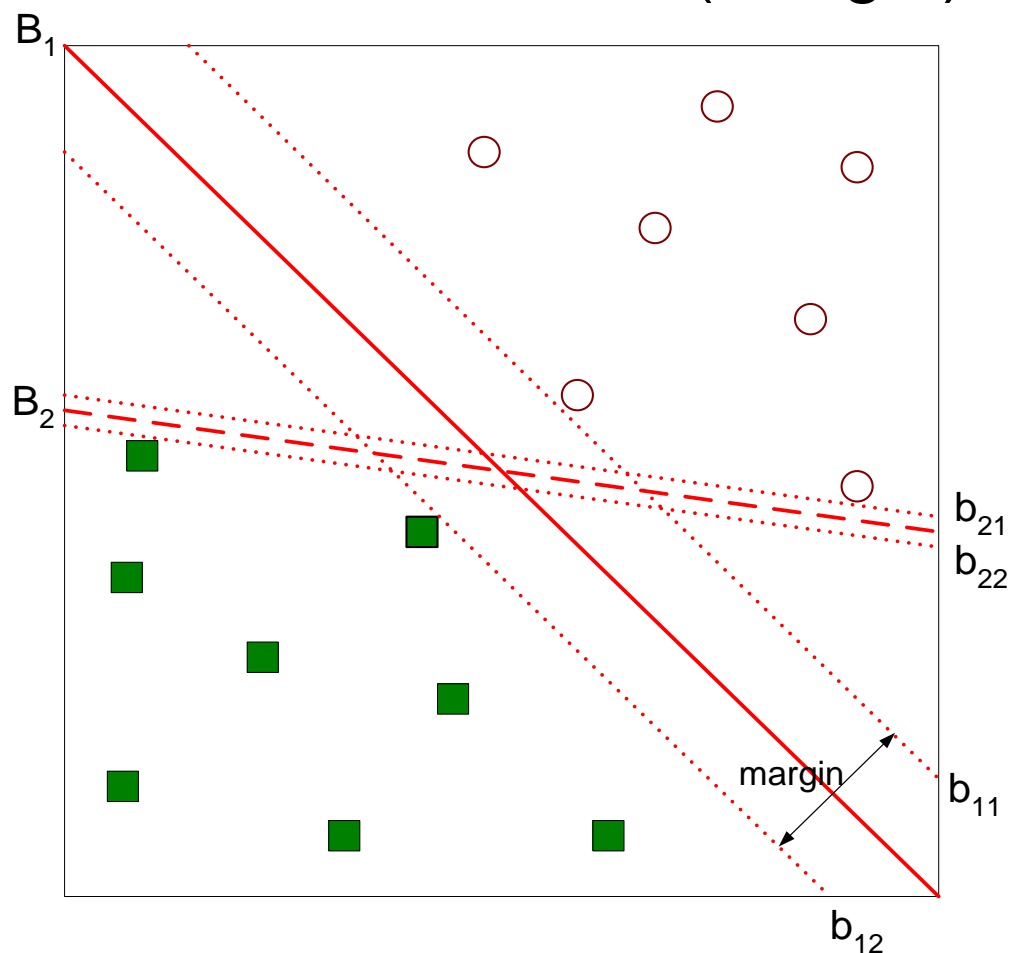
# 一、最大边缘超平面

为了更好地解释不同的超平面对泛化误差的影响，考虑两个超平面 $B_1$ 和 $B_2$ 。



# 一、最大边缘超平面

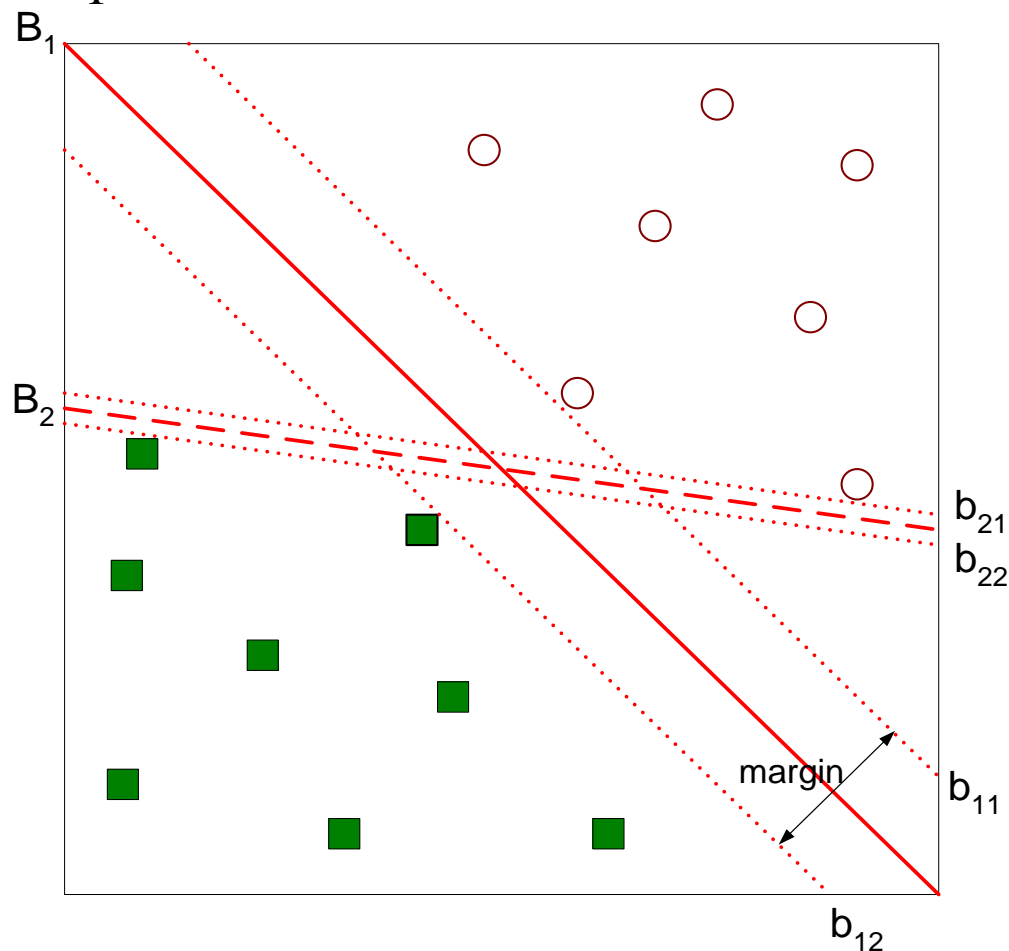
每个超平面 $B_i$ 都对应着一对平行的超平面： $b_{i1}$ 和 $b_{i2}$ ，它们之间的间隔称为超平面的边缘(margin)。





# 一、最大边缘超平面

图中，超平面 $B_1$ 的边缘明显大于 $B_2$ 的边缘。因此，在这个例子中， $B_1$ 就是训练样本的最大边缘超平面。



## 二、线性支持向量机

---

- 直觉上，决策边界的边缘较小，决策边界任何轻微的扰动都可能对分类结果产生较大的影响。也就是说，具有较大边缘的决策边界比那些具有较小边缘的决策边界具有更好的泛化误差。
- 因此，根据结构风险最小化原理，需要设计最大化决策边界的边缘的线性分类器，以确保最坏情况下的泛化误差最小。线性支持向量机(linear SVM) 就是这样的分类器。

## 二、线性支持向量机

- 给定训练数据集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}, y \in \{-1, +1\}$   
线性分类器的决策边界可以写成如下线性方程：

$$\mathbf{w}^\top \mathbf{x} + b = 0$$

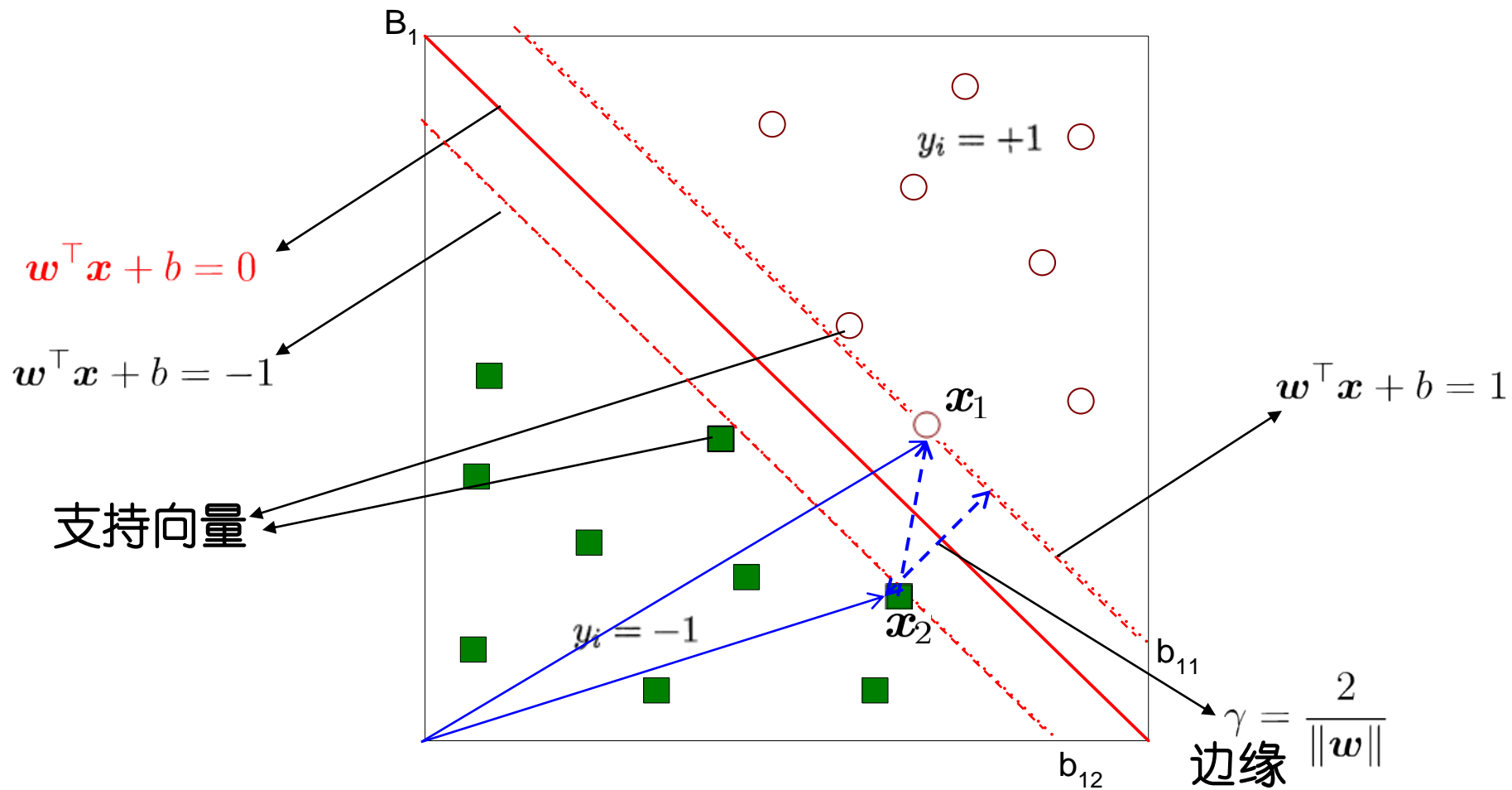
其中  $\mathbf{w} = (w_1; w_2; \dots; w_m)$  为法向量，决定了决策边界的方向； $b$  为位移量，决定了决策边界与原点之间的距离。显然，决策边界由参数  $\mathbf{w}$  和  $b$  确定。

- 假设决策边界能将训练样本正确分类，即对于任意样本点  $(\mathbf{x}_i, y_i) \in D$ ：若  $y_i = +1$ ，则有  $\mathbf{w}^\top \mathbf{x}_i + b > 0$ ；若  $y_i = -1$  则有  $\mathbf{w}^\top \mathbf{x}_i + b < 0$ 。那么通过调整决策边界的参数  $\mathbf{w}$  和  $b$  总可以得到：
$$\begin{cases} \mathbf{w}^\top \mathbf{x}_i + b \geq +1, & y_i = +1; \\ \mathbf{w}^\top \mathbf{x}_i + b \leq -1, & y_i = -1. \end{cases}$$

## 二、线性支持向量机

- 距离决策边界最近的训练样本点使得上式中的等号成立，因此被称为“支持向量” (support vector)。
- 两个异类支持向量到决策边界的距离之和被称为决策边界的“边缘” (margin)，刚好等于超平面  $b_{i1}$  和  $b_{i2}$  之间的间隔  $\gamma$ 。
- 为了计算边缘  $\gamma$ ，令  $x_1$  是  $b_{i1}$  上的一个样本点， $x_2$  是  $b_{i2}$  上的一个样本点，将  $x_1$  和  $x_2$  分别代入上式，得到  $b_{i1} : w^T x_1 + b = 1$   $b_{i2} : w^T x_2 + b = -1$   
令两式相减得到： $w^T (x_1 - x_2) = 2$  即  $w^T \overrightarrow{x_2 x_1} = 2$   
向量  $\overrightarrow{x_2 x_1}$  在  $w$  上的投影乘以  $w$  的模： $\|w\| \times \gamma = 2$   
即：
$$\gamma = \frac{2}{\|w\|}$$

## 二、线性支持向量机



## 二、线性支持向量机

- 因此，线性支持向量机的学习就是要寻找满足约束条件的参数 $\mathbf{w}$ 和 $b$ ，使得 $\gamma$ 最大，即：

$$\arg \max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|}$$
$$s.t. \begin{cases} \mathbf{w}^\top \mathbf{x}_i + b \geq +1, & y_i = +1; \\ \mathbf{w}^\top \mathbf{x}_i + b \leq -1, & y_i = -1. \end{cases}$$



等价于

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$
$$s.t. \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, i = 1, 2, \dots, n.$$

## 二、线性支持向量机

---

- 由于目标函数是二次的，并且约束条件在参数  $w$  和  $b$  上是线性的，因此线性支持向量机的学习问题是一个凸二次优化问题，可以直接用现成的优化计算包求解，或者用更高效的拉格朗日乘子法求解。

## 二、线性支持向量机

### ● 拉格朗日乘子法

- ✓ 第一步：引入拉格朗日乘子 $\alpha_i \geq 0$ 得到拉格朗日函数：

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

- ✓ 第二步：令 $L(\mathbf{w}, b, \boldsymbol{\alpha})$ 对 $\mathbf{w}$ 和 $b$ 的偏导为零可得：

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

- ✓ 第三步：回代拉格朗日函数得到对偶优化问题：

$$L(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$$



## 二、线性支持向量机

- 因为对偶优化问题的二次项前面有个负号，而原始线性支持向量机学习的优化问题的二次项前面没有负号，这说明原来要优化的最小化问题已经转化成了最大化对偶优化问题。即：

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

- 不过，如果再对上述最大问题稍作变形，就可等价于下面的最小化问题：

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

## 二、线性支持向量机

- 再结合拉格朗日函数的约束条件，就可得到原始问题的最终优化问题：

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

$$s.t. \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad i = 1, 2, \dots, n.$$

## 二、线性支持向量机

- 解出  $\alpha$  后, 求出  $\mathbf{w}$  和  $b$  即可得到线性SVM的最终模型:  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^\top \mathbf{x} + b$
- 因为线性SVM还需满足不等式约束  $y_i f(\mathbf{x}_i) \geq 1$ , 同时  $\alpha_i \geq 0$ , 所以可以把不等式约束变换成等式约束, 这种变换得到拉格朗日乘子约束, 也称为Karush-Kuhn-Tucker (KKT)  $\alpha_i (y_i f(\mathbf{x}_i) - 1) = 0$
- KKT条件表明, 除非训练样本满足  $y_i f(\mathbf{x}_i) = 1$ , 否则必有拉格朗日乘子  $\alpha_i = 0$ 。那些  $\alpha_i > 0$  的训练样本, 即  $y_i f(\mathbf{x}_i) = 1$  的训练样本都落在最大边缘的边界  $b_{i1}$  和  $b_{i2}$  上, 都是支持向量。可见, 最终支持向量机模型的参数  $\mathbf{w}$  和  $b$  仅仅依赖于这些支持向量。

## 二、线性支持向量机

---

- 由于对偶优化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；然而，该问题的规模正比于训练样本的数目，对于大型数据集，会造成很大的开销。为了避开这个障碍，人们通过利用问题本身的特性，提出了很多高效的算法，比如：**SMO** 序列最小化优化算法 (Sequential Minimal Optimization) 就是其中一个著名的代表，具体算法可参阅[Platt, 1998]:

J. Platt: **Fast Training of Support Vector Machines using Sequential Minimal Optimization**. In B. Schoelkopf and C. Burges and A. Smola, editors, *Advances in Kernel Methods - Support Vector Learning*, 1998.

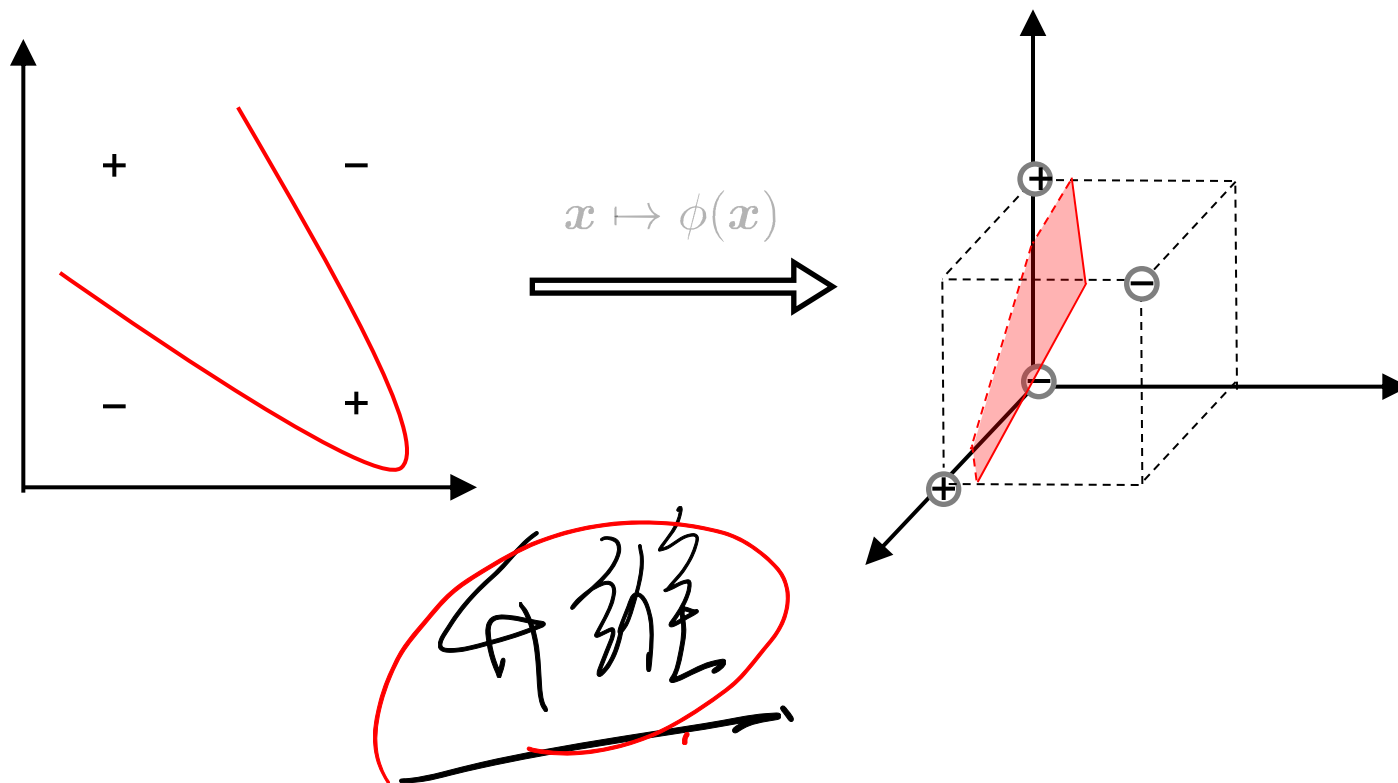
## 三、非线性支持向量机

---

- 线性SVM假定训练样本是线性可分的，即存在一个线性的决策边界能将所有的训练样本正确分类。
- 然而在实际应用中，在原始的样本空间内也许并不存在这样的决策边界。
- 对于这样的问题，可将样本从原始空间映射到一个更高维的特征空间，使得样本在映射后的特征空间内线性可分。例如在下图中，如果将原始的二维空间映射到一个合适的三维空间，就能找到一个合适的划分超平面。幸运的是，如果原始空间是有限维，即属性数目有限，那么一定存在一个更高维的特征空间使得样本线性可分。

### 三、非线性支持向量机

#### 异或问题的非线性映射



# 三、非线性支持向量机

- 设样本  $\mathbf{x}$  映射后的向量为  $\phi(\mathbf{x})$ , 决策边界为  $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$

原始问题

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

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

对偶问题

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

$$s.t. \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad i = 1, 2, \dots, n.$$

预测模型

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_{i=1}^n \alpha_i y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) + b$$

### 三、非线性支持向量机

- 求解上面的优化问题，涉及到计算  $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ ，这是样本  $\mathbf{x}_i$  与  $\mathbf{x}_j$  映射到特征空间之后的内积。由于特征空间维数可能很高，甚至可能是无穷维，因此直接计算  $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$  通常是非常困难的。为了避开这个障碍，一个基本的想法是：不显式地设计映射  $\phi(\cdot)$ ，而是设计一个核函数：

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

- $\mathbf{x}_i$  和  $\mathbf{x}_j$  在高维特征空间的内积等于它们在原始样本空间中通过核函数  $k(\cdot, \cdot)$  计算的结果。有了这样的核函数，我们就不必直接去计算高维甚至无穷维特征空间中的内积了。



## 三、非线性支持向量机

- 那么，合适的核函数是否一定存在呢？什么样的函数才能做核函数呢？
- 根据**Mercer**定理，只要一个对称函数所对应的核矩阵半正定，那么它就能作为核函数来使用。
- 常用核函数包括：

名称	表达式	参数
线性核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$	
多项式核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j)^d$	$d \geq 1$ 为多项式的次数
高斯核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\delta^2}\right)$	$\delta > 0$ 为高斯核的带宽(width)
拉普拉斯核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ }{\delta}\right)$	$\delta > 0$
Sigmoid核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta \mathbf{x}_i^\top \mathbf{x}_j + \theta)$	$\tanh$ 为双曲正切函数, $\beta > 0, \theta < 0$

- 于是，核函数的选择成为非线性**SVM**的最大变数，若核函数选择不合适，就意味着将样本映射到了一个不合适的特征空间，从而很可能导致非线性**SVM**性能不佳。