



机器学习与数据挖掘

支持向量机

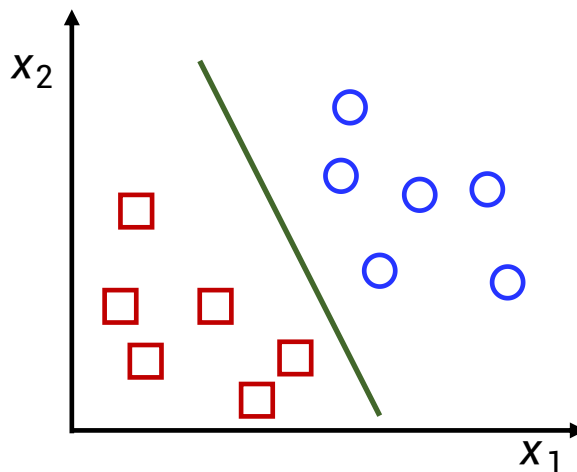


课程大纲

- 线性分类器的决策边界
- 线性最大间隔分类器
- 软线性最大间隔分类器
- 支持向量机
- 与逻辑回归的关系

线性分类器中的决策边界

- 在线性分类器中，决策边界总是一个超平面。目标是找到能够分离不同类型样本的超平面

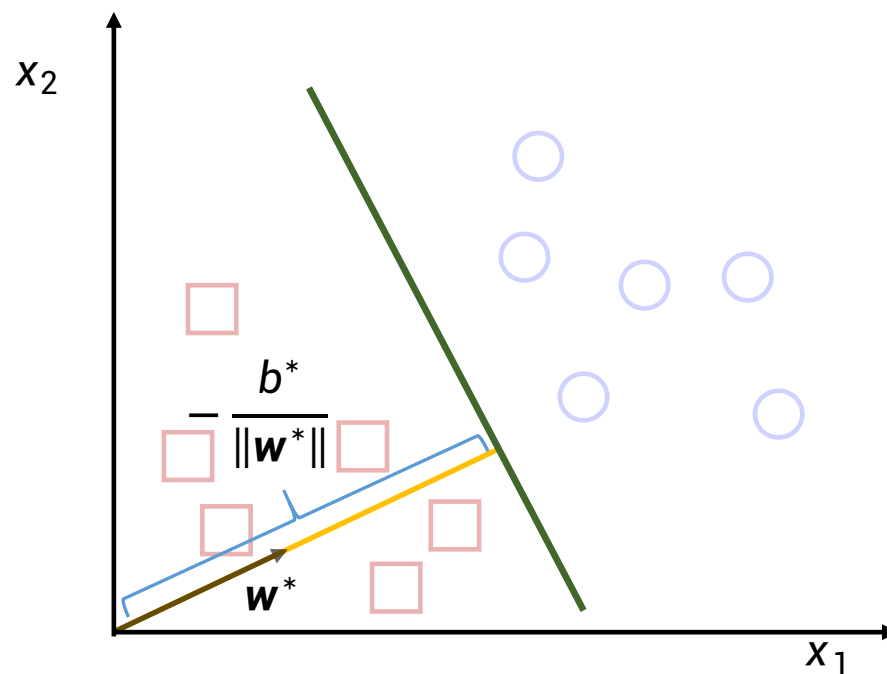


- 逻辑回归
 - 决策边界超平面是通过最小化交叉熵损失来找到

$$L(\mathbf{w}, b) = -y \log(\sigma(\mathbf{w}^T \mathbf{x} + b)) - (1 - y) \log(1 - \sigma(\mathbf{w}^T \mathbf{x} + b))$$

➤ 当得到最优的 \mathbf{w}^* 和 b^* 时，超平面由下式的 \mathbf{x} 组成：

$$\{\mathbf{x} | \mathbf{w}^{*T} \mathbf{x} + b^* = 0\}$$



1) 超平面与向量 \mathbf{w}^* 垂直

2) 从原点到平面的距离是 $-\frac{b^*}{\|\mathbf{w}^*\|}$

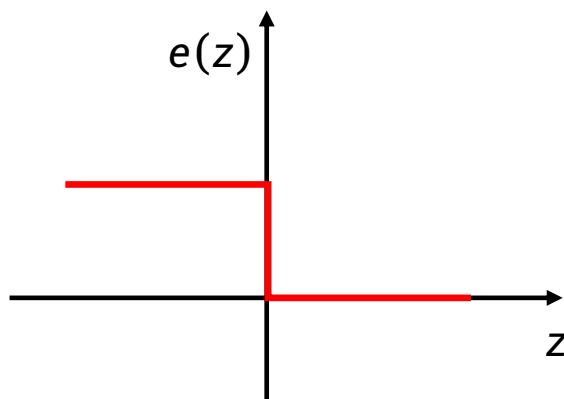
- 理想分类器

➤ 超平面是通过最小化损失来确定的

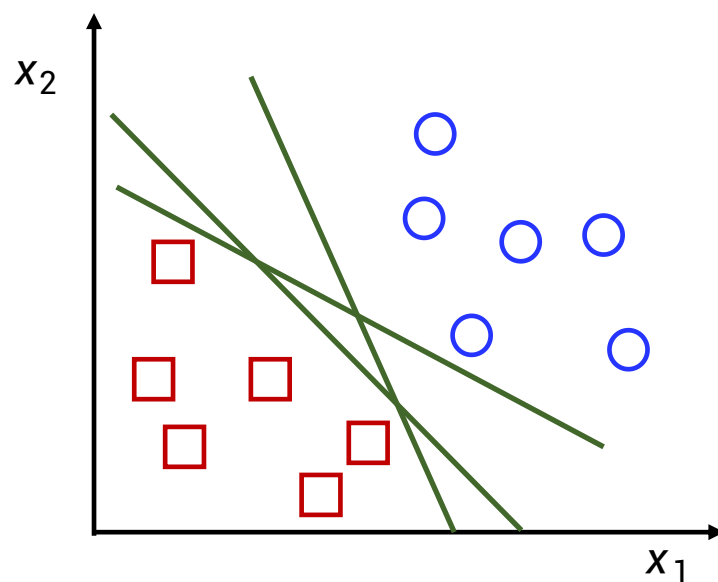
$$L(\mathbf{w}, b) = \sum_{\ell=1}^N e\left(y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + b)\right)$$

$L(\mathbf{w}, b)$ 代表被错误分类的样本数量

- $y \in \{-1, 1\}$
- 如果 $z \geq 0$, 则 $e(z) = 0$; 否则 $e(z) = 1$



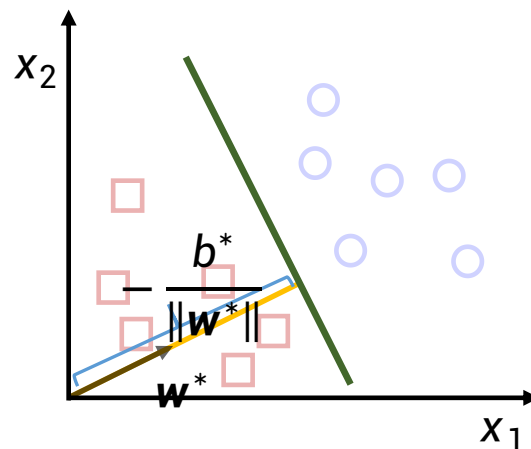
- 如果样本是线性可分的，将会有无数个理想分类器，它们由 \mathbf{w}^* 和 b^* 决定
- 每一对 \mathbf{w}^* 和 b^* 都对应一个超平面



上述所有超平面都能使损失降为零

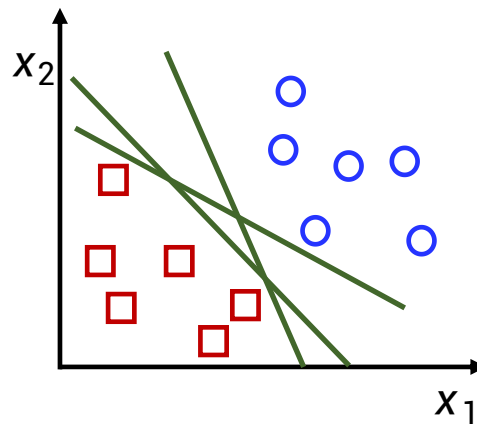
哪个超平面是最好的？

- 下面的超平面从最小化交叉熵损失的角度来看是最优的



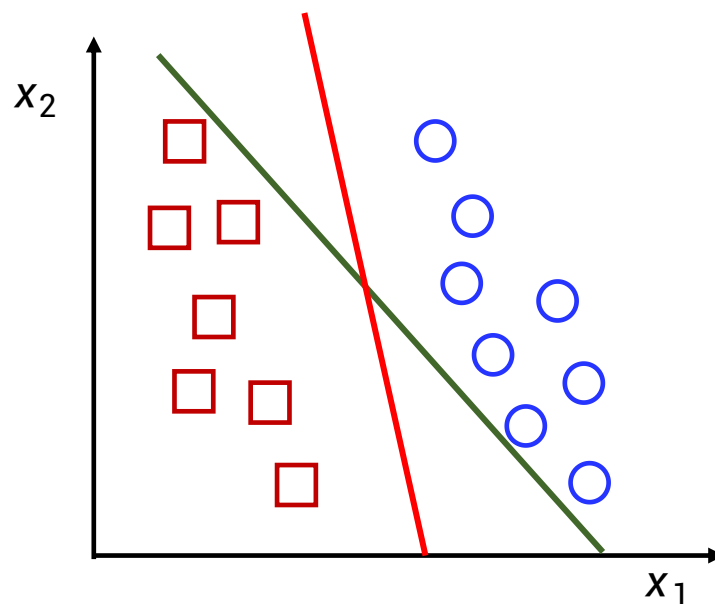
哪个超平面是最好的？

- 下面的超平面从**最小化误分类样本数量**的角度来看是最优的



- 这些超平面都是在 *训练样本* 上进行评估的
- 但我们真正需要的是在 *(未见过的) 测试数据* 上的表现

根据我们的直觉，下面的哪个超平面更有可能在测试数据上产生更好的结果？

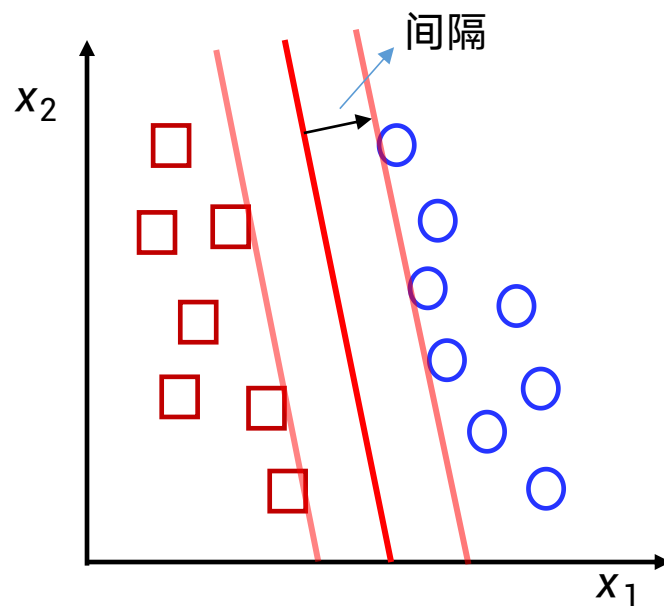


课程大纲

- 线性分类器的决策边界
- 线性最大间隔分类器
- 软线性最大间隔分类器
- 支持向量机
- 与逻辑回归的关系

最大间隔目标

- 为了在未见过的数据上表现良好，直觉是找到一个能尽可能扩大间隔的超平面



当间隔较大时，我们可以预期 *未见过的样本有更高的几率被正确分类*

如何表示间隔？

- 样本 \mathbf{x} 到超平面 \mathcal{H} 的距离。记 $\mathbf{w}^T \mathbf{x} + b = h(\mathbf{x})$

➤ 每个 \mathbf{x} 都可以被分解为：

$$\mathbf{x} = \mathbf{m}_1 + \mathbf{m}_2$$

- \mathbf{m}_1 在超平面 \mathcal{H} 上，即 $\mathbf{w}^T \mathbf{m}_1 + b = 0$

- $\mathbf{m}_2 \perp \mathcal{H}$ 及 $\mathbf{m}_2 \parallel \mathbf{w}$

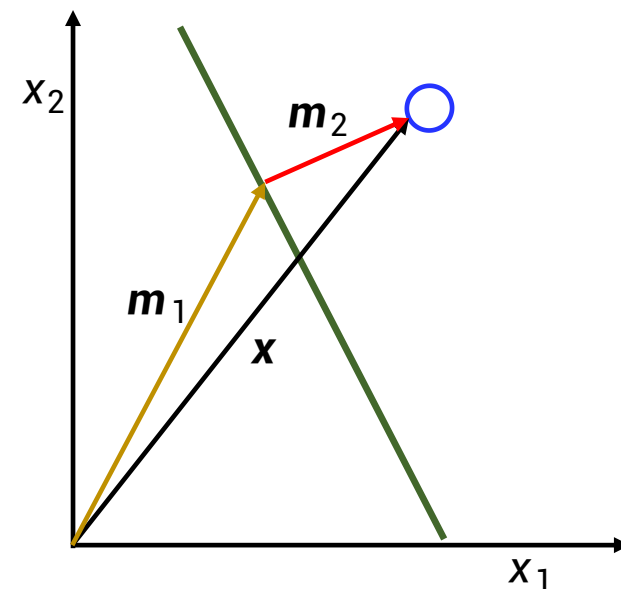
➤ 因此，我们有：

$$\mathbf{w}^T \mathbf{x} + b = \mathbf{w}^T (\mathbf{m}_1 + \mathbf{m}_2) + b = \mathbf{w}^T \mathbf{m}_2 = h(\mathbf{x})$$

➤ 由于 $\mathbf{m}_2 \parallel \mathbf{w}$ ，我们可以写成：

$$\mathbf{m}_2 = \gamma \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|},$$

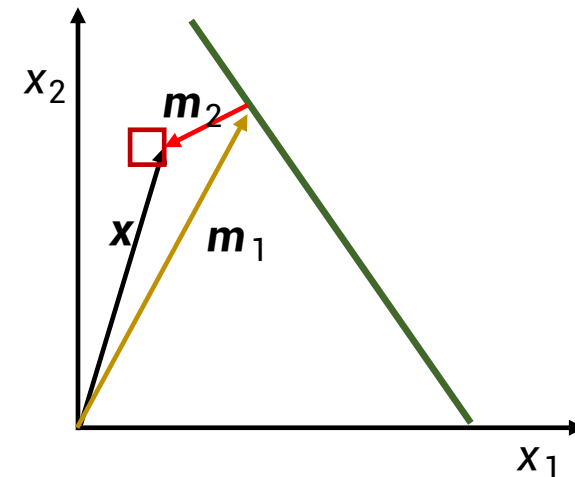
其中 $|\gamma|$ 代表 \mathbf{m}_2 的长度



➤ 将 $\mathbf{m}_2 = y \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|}$ 代入 $h(\mathbf{x}) = \mathbf{w}^T \mathbf{m}_2$ 得到:

$$h(\mathbf{x}) = y \cdot \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \Rightarrow y = \frac{h(\mathbf{x})}{\|\mathbf{w}\|}$$

- 位于超平面上方的样本 \mathbf{x} 到超平面的距离为 $\frac{h(\mathbf{x})}{\|\mathbf{w}\|}$
- 位于超平面下方的样本 \mathbf{x} 到超平面的距离为 $-\frac{h(\mathbf{x})}{\|\mathbf{w}\|}$



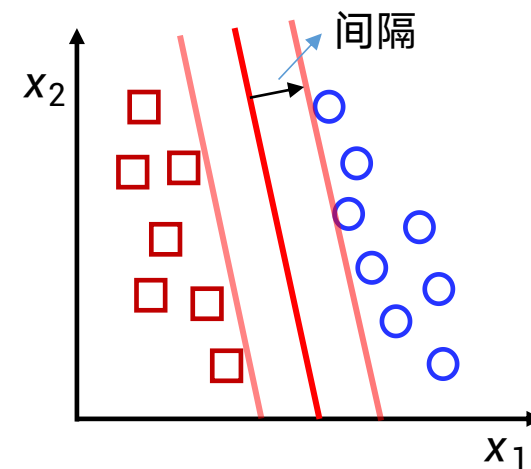
➤ 样本 (\mathbf{x}, y) 到超平面的距离由下式给出:

$$\frac{y \cdot h(\mathbf{x})}{\|\mathbf{w}\|} = \frac{y \cdot (\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|}$$

其中 $y \in \{-1, 1\}$

- 一个超平面在一个数据集上的间隔由最小距离给出，即：

$$\text{Margin} = \min_{\ell} \frac{y^{(\ell)} \cdot (\mathbf{w}^T \mathbf{x}^{(\ell)} + b)}{\|\mathbf{w}\|}$$



- 因此，最大间隔分类器的目标是找到最优的 \mathbf{w}^* 和 b^* 来使间隔最大化，即：

$$\mathbf{w}^*, b^* = \arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_{\ell} [y^{(\ell)} \cdot (\mathbf{w}^T \mathbf{x}^{(\ell)} + b)] \right\}$$

但是，如何优化是未知的

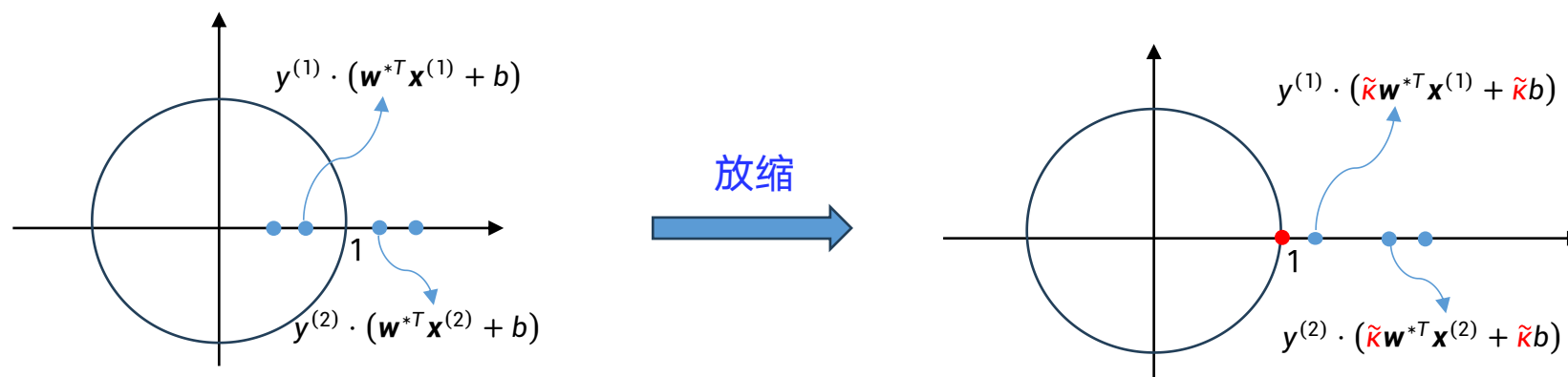
变换后的目标函数

- 基本思想：优化一个与原问题具有相同最优解的目标函数

- 假设 \mathbf{w}^* 和 b^* 是 $\frac{1}{\|\mathbf{w}\|} \min_{\ell} [y^{(\ell)} \cdot (\mathbf{w}^T \mathbf{x}^{(\ell)} + b)]$ 的最优解。那么，对于所有 $\kappa \neq 0$ ， $\kappa \mathbf{w}^*$ 和 κb^* 也必定是最优解
- 此外，总存在一个特定的 $\tilde{\kappa}$ 使得

$$\text{对于所有 } \ell = 1, 2, \dots, n, \quad y^{(\ell)} \cdot (\tilde{\kappa} \mathbf{w}^{*T} \mathbf{x}^{(\ell)} + \tilde{\kappa} b^*) \geq 1$$

并且至少有一个等号成立



➤ 因此，可以通过求解无约束优化问题来找到 $\frac{1}{\|\mathbf{w}\|} \min_{\ell} [y^{(\ell)} \cdot (\mathbf{w}^T \mathbf{x}^{(\ell)} + b)]$ 的最大值：

$$\max_{\mathbf{w}, b} \left[\frac{1}{\|\mathbf{w}\|} \min_{\ell} [y^{(\ell)} \cdot (\mathbf{w}^T \mathbf{x}^{(\ell)} + b)] \right]$$

或者通过求解带约束的优化问题来找到：

$$\max_{\tilde{\mathbf{w}}, \tilde{b}} \left[\frac{1}{\|\tilde{\mathbf{w}}\|} \min_{\ell} [y^{(\ell)} \cdot (\tilde{\mathbf{w}}^T \mathbf{x}^{(\ell)} + \tilde{b})] \right]$$

$$\text{s. t.: 对于所有 } \ell = 1, 2, \dots, n, y^{(\ell)} \cdot (\tilde{\mathbf{w}}^T \mathbf{x}^{(\ell)} + \tilde{b}) \geq 1$$

并且至少有一个等号成立

- 最优解 \mathbf{w}^* 和 $\tilde{\mathbf{w}}^*$ 可能不相同，但它们所得到的 $\frac{1}{\|\mathbf{w}^*\|} \min_{\ell} [y^{(\ell)} \cdot (\mathbf{w}^{*T} \mathbf{x}^{(\ell)} + b^*)]$ 和 $\frac{1}{\|\tilde{\mathbf{w}}^*\|} \min_{\ell} [y^{(\ell)} \cdot (\tilde{\mathbf{w}}^{*T} \mathbf{x}^{(\ell)} + \tilde{b}^*)]$ 的值必定相等

- 在第二个优化问题中，由于对于所有 $\ell = 1, 2, \dots, n$ ， $y^{(\ell)} \cdot (\tilde{\mathbf{w}}^T \mathbf{x}^{(\ell)} + \tilde{b}) \geq 1$ 并且至少有一个等号成立，我们可以很容易地得到：

$$\min_{\ell} [y^{(\ell)} \cdot (\tilde{\mathbf{w}}^T \mathbf{x}^{(\ell)} + \tilde{b})] = 1$$

因此，优化目标可以简化为 $\frac{1}{\|\tilde{\mathbf{w}}\|}$

- 最大化 $\frac{1}{\|\tilde{\mathbf{w}}\|}$ 可以等价于最小化 $\|\tilde{\mathbf{w}}\|^2$ 。因此，该问题可以等价地写为：

$$\begin{aligned} \min_{\tilde{\mathbf{w}}, \tilde{b}} \quad & \|\tilde{\mathbf{w}}\|^2 \\ \text{s. t.} \quad & \text{对于所有 } \ell = 1, 2, \dots, n, \quad y^{(\ell)} \cdot (\tilde{\mathbf{w}}^T \mathbf{x}^{(\ell)} + \tilde{b}) \geq 1 \end{aligned}$$

至少有一个等号成立

- 当最小化 $\|\tilde{\mathbf{w}}\|^2$ 时，“至少有一个等号成立”的约束条件将自动满足（为什么？）。因此，可以在不影响结果的情况下将其移除


- 因此，最大间隔超平面可以通过求解下面的优化问题来找到：

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t. : } & y^{(\ell)} \cdot (\mathbf{w}^T \mathbf{x}^{(\ell)} + b) \geq 1, \quad \text{for } \ell = 1, 2, \dots, N \end{aligned}$$

- 这是一个二次规划问题。它的最优解可以高效地通过数值方法找到

- 这是一个二次规划问题。它的最优解可以高效地通过数值方法找到
- 二次规划问题：

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T Q x + c^T x \text{ subject to } Ax \leq b, \quad Ex = d, \quad l \leq x \leq u$$


$$\frac{1}{2} \|w\|^2$$

x 是决策变量的向量

Q 是一个 $n \times n$ 的对称正定矩阵，定义了二次项

c 是线性项的系数向量

A 和 b 定义了不等式约束

E 和 d 定义了等式约束

l 和 u 定义了变量的下界和上界

- 这是一个二次规划问题。它的最优解可以高效地通过数值方法找到

数值方法：牛顿法，
梯度下降法（神经网络学习），
共轭梯度下降法，
.....

- 当得到最优的 \mathbf{w}^* 和 b^* 时，一个未见过的数据 \mathbf{x} 可以被分类为：

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

等价对偶形式

- 每个凸优化问题都对应一个等价的对偶形式

本节内容摘自凸优化这一学科

- 原始优化问题的拉格朗日函数

$$\mathcal{L}(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{\ell=1}^N \mathbf{a}_{\ell} (y^{(\ell)} (\mathbf{w}^T \mathbf{x}^{(\ell)} + b) - 1),$$

其中，拉格朗日乘子 \mathbf{a}_{ℓ} 必须满足 $\mathbf{a}_{\ell} \geq 0$

- 拉格朗日对偶函数

$$g(\mathbf{a}) = \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \mathbf{a})$$

- 原始优化问题的对偶形式

$$\begin{array}{ll} \max_{\mathbf{a}} & g(\mathbf{a}) \\ \text{s. t.} & \mathbf{a} \geq \mathbf{0} \end{array}$$

- 推导函数 $g(\mathbf{a})$ 的闭式表达式

➤ 将梯度 $\frac{\partial \mathcal{L}}{\partial \mathbf{w}}$ 和 $\frac{\partial \mathcal{L}}{\partial b}$ 设置为 0, 得到:

$$\mathbf{w} = \sum_{\ell=1}^N a_{\ell} \mathbf{y}^{(\ell)} \mathbf{x}^{(\ell)} \quad \sum_{\ell=1}^N a_{\ell} \mathbf{y}^{(\ell)} = 0$$

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{\ell=1}^N a_{\ell} (y^{(\ell)} (\mathbf{w}^T \mathbf{x}^{(\ell)} + b) - 1)$$

➤ 将它们代入 $\mathcal{L}(\mathbf{w}, b, \mathbf{a})$, 得到 $g(\mathbf{a}) = \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \mathbf{a})$ 为:

$$g(\mathbf{a}) = \underbrace{\sum_{\ell=1}^N a_{\ell} - \frac{1}{2} \sum_{\ell=1}^N \sum_{j=1}^N a_{\ell} a_j y^{(\ell)} y^{(j)} \mathbf{x}^{(\ell)T} \mathbf{x}^{(j)}}_{= \mathbf{1}^T \mathbf{a} - \frac{1}{2} \mathbf{a}^T \mathbf{M} \mathbf{a}}$$

其中 $[\mathbf{M}]_{\ell j} \triangleq y^{(\ell)} y^{(j)} \mathbf{x}^{(\ell)T} \mathbf{x}^{(j)}$

- 因此，对偶优化问题变为：

$$\begin{array}{ll} \max_{\mathbf{a}} & g(\mathbf{a}) \\ \text{s. t. :} & \mathbf{a} \geq \mathbf{0} \text{ 及 } \sum_{\ell=1}^N a_{\ell} y^{(\ell)} = 0 \end{array}$$

$$\begin{aligned} \text{其中 } g(\mathbf{a}) &= \underbrace{\sum_{\ell=1}^N a_{\ell} - \frac{1}{2} \sum_{\ell=1}^N \sum_{j=1}^N a_{\ell} a_j y^{(\ell)} y^{(j)} \mathbf{x}^{(\ell)T} \mathbf{x}^{(j)}}_{= \mathbf{1}^T \mathbf{a} - \frac{1}{2} \mathbf{a}^T \mathbf{M} \mathbf{a}} \end{aligned}$$

它也是一个二次规划问题，可以通过数值方法高效地求解

- 最优解 \mathbf{w}^* 、 b^* 与最优解 \mathbf{a}^* 之间的关系

➤ 给定最优解 \mathbf{a}^* ，根据 $\mathbf{w} = \sum_{\ell=1}^N a_{\ell} y^{(\ell)} \mathbf{x}^{(\ell)}$ ，最优的 \mathbf{w}^* 可以等价地表示为：

$$\mathbf{w}^* = \sum_{\ell=1}^N a_{\ell}^* y^{(\ell)} \mathbf{x}^{(\ell)}$$

➤ 由于对于所有位于间隔边界上的样本 $(\mathbf{x}^{(\ell)}, y^{(\ell)})$ ，都有 $y^{(\ell)} (\mathbf{w}^{*T} \mathbf{x}^{(\ell)} + b) = 1$ ，我们可以推导出：

$$b^* = \frac{1}{N_{\mathcal{S}}} \sum_{n \in \mathcal{S}} \left(y^{(n)} - \sum_m a_m^* y^{(m)} \mathbf{x}^{(n)T} \mathbf{x}^{(m)} \right)$$

其中 \mathcal{S} 表示位于间隔边界上的样本集合

- 最大间隔分类器

- 原问题形式

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

- 对偶问题形式

将 $\mathbf{w}^* = \sum_{n=1}^N a_n^* y^{(n)} \mathbf{x}^{(n)}$ 代入原问题形式，得到：

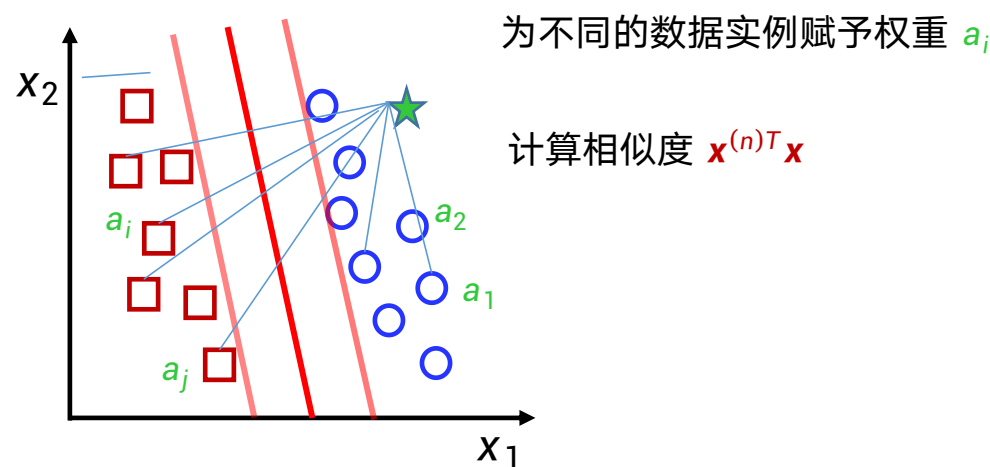
$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N a_n^* y^{(n)} \mathbf{x}^{(n)T} \mathbf{x} + b^* \right)$$

这两个分类器是等价的

$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N a_n^* \cdot (\mathbf{x}^{(n)T} \mathbf{x}) \cdot y^{(n)} + b^* \right)$$

- 如何理解对偶最大间隔分类器？

- 对于一个测试样本 \mathbf{x} ，通过计算它与所有训练样本 $\mathbf{x}^{(n)}$ ($n = 1, \dots, N$) 的相似度 $\mathbf{x}^{(n)T} \mathbf{x}$
- 将所有标签 $y^{(n)}$ 按样本相似度 $\mathbf{x}^{(n)T} \mathbf{x}$ 和乘子 a_n^* 进行加权求和



原问题与对偶问题比较

- 优化复杂度

原问题

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t.} \quad & y^{(\ell)} \cdot (\mathbf{w}^T \mathbf{x}^{(\ell)} + b) \geq 1, \\ & \text{for } \ell = 1, 2, \dots, N \end{aligned}$$

要优化的参数数量：特征维度

对偶问题

$$\begin{aligned} \max_{\mathbf{a}} \quad & g(\mathbf{a}) \\ \text{s. t.} \quad & \mathbf{a} \geq \mathbf{0} \\ & \sum_{\ell=1}^N a_{\ell} y^{(\ell)} = 0 \end{aligned}$$

要优化的参数数量：训练样本数量

在高维特征情况下，求解对偶问题更高效

- 测试复杂度

原问题

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

只需要一次内积运算
 $\mathbf{w}^{*T} \mathbf{x}$

对偶问题

$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N a_n^* (\mathbf{x}^{(n)T} \mathbf{x}) y^{(n)} + b^* \right)$$

需要 N 次内积运算 $\mathbf{x}^{(n)T} \mathbf{x}$, 其中
 $n = 1, 2, \dots, N$

乍一看, 对偶分类器看起来比原分类器昂贵得多

- 幸运的是, 可以证明绝大多数 a_n^* 都为 0

拉格朗日乘子 a^* 的稀疏性

- 对于任何凸优化问题，最优解都满足 **KKT 条件**，对于我们这个问题，KKT 条件为：

KKT条件：

- 1.原始可行性：约束被满足
- 2.对偶可行性：拉格朗日乘子不小于零
- 3.互补松弛性：拉格朗日乘子与约束的乘积为零
- 4.梯度条件：最优点的梯度为零

$$a_n^* \geq 0$$

$$y^{(n)} (\mathbf{w}^{*T} \mathbf{x}^{(n)} + b^*) - 1 \geq 0$$

$$a_n^* [y^{(n)} (\mathbf{w}^{*T} \mathbf{x}^{(n)} + b^*) - 1] = 0$$

前两个条件来自原始的原问题和对偶问题

拉格朗日乘子 a^* 的稀疏性

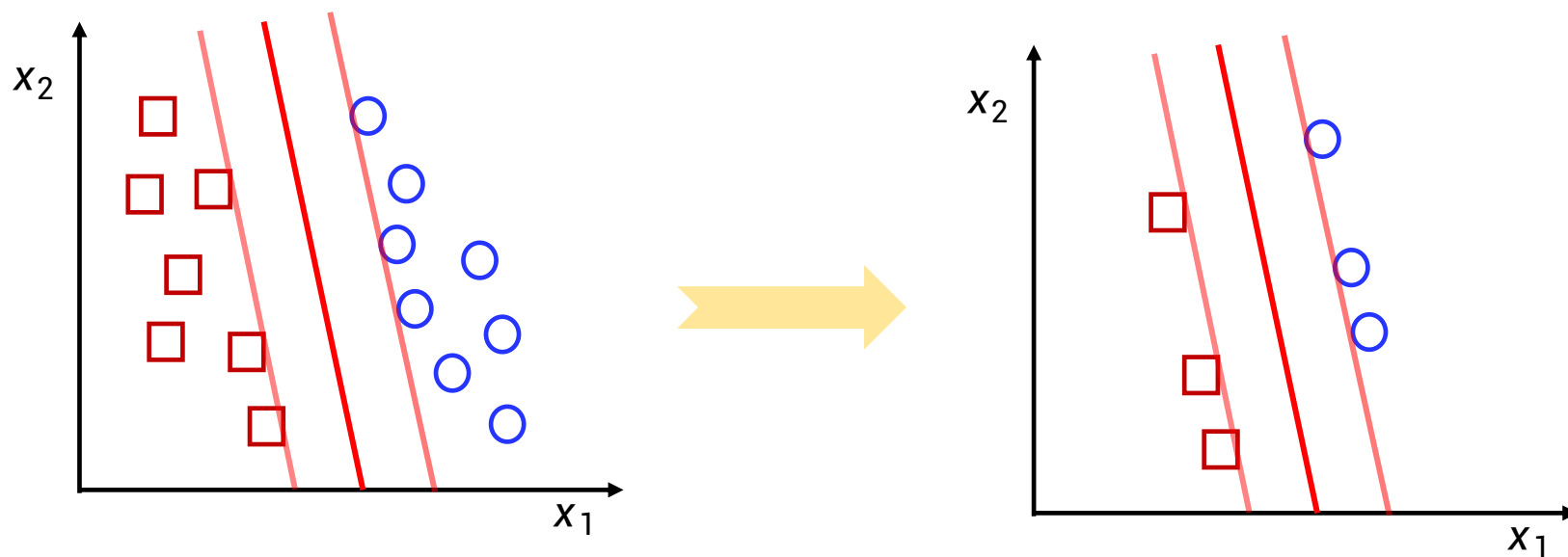
- 从最后一个条件可以看出，只有当 $y^{(n)}(\mathbf{w}^{*T} \mathbf{x}^{(n)} + b^*) = 1$ 时， $a_n^* \neq 0$
- 如果 $\mathbf{x}^{(n)}$ 满足 $y^{(n)}(\mathbf{w}^{*T} \mathbf{x}^{(n)} + b^*) = 1$ ，这意味着它位于间隔边界上

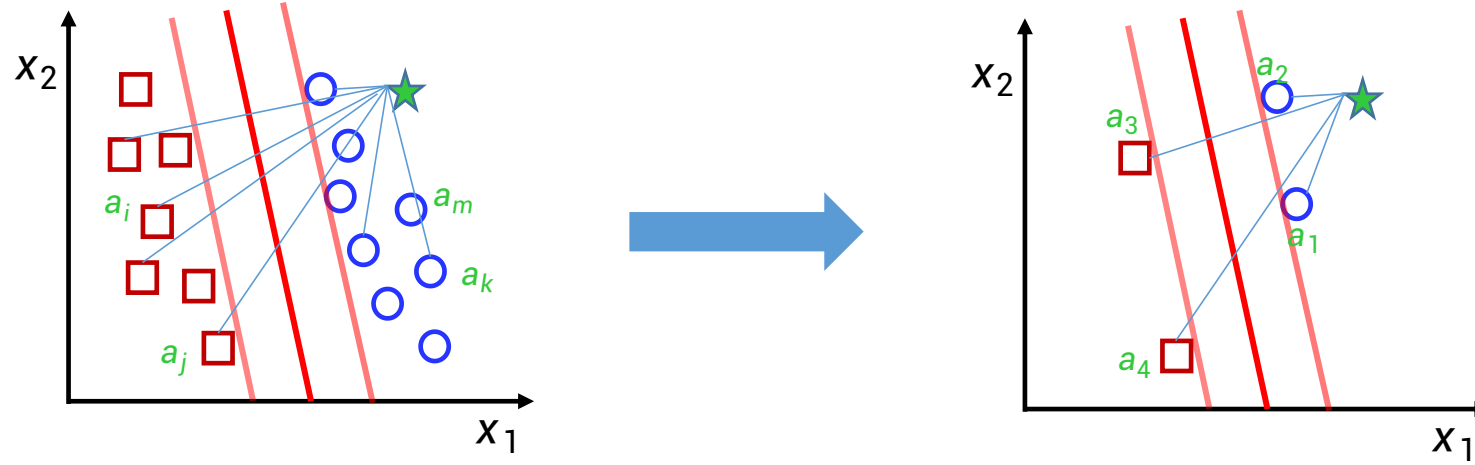
这类样本被称为支持向量

- 因此，当我们对一个未见过的样本 \mathbf{x} 进行分类时：

$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_n a_n^* y^{(n)} \mathbf{x}^{(n)T} \mathbf{x} + b^* \right),$$

我们只需要评估 \mathbf{x} 与支持向量（样本）之间的相似度 $\mathbf{x}^{(n)T} \mathbf{x}$





$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N \left(a_n^* (\mathbf{x}^{(n)T} \mathbf{x}) \right) \cdot y^{(n)} + b^* \right)$$

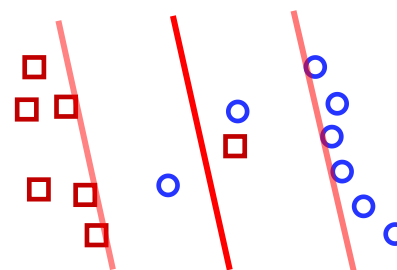
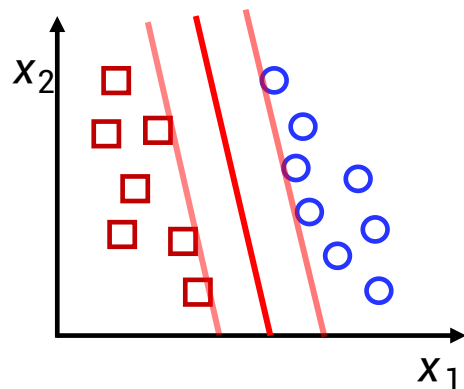
课程大纲

- 线性分类器的决策边界
- 线性最大间隔分类器
- 软线性最大间隔分类器
- 支持向量机
- 与逻辑回归的关系

非线性可分类别

- 之前最大间隔分类器中隐含的假设

训练样本是线性可分的!!!



非线性可分类别

- 在这种情况下，优化问题会发生什么？

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t.} \quad & y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b) \geq 1, \text{ for } n = 1, 2, \dots, N \end{aligned}$$

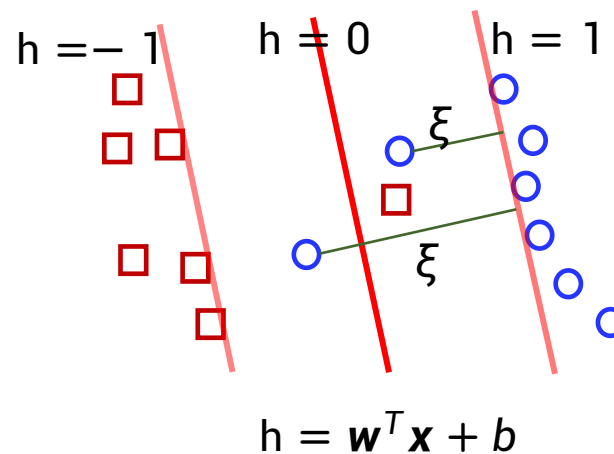
- 该优化问题没有可行解。也就是说，*不存在这样的超平面*

软最大间隔

- 为了解决这个问题，我们不再要求对于所有 $n = 1, \dots, N$ 都必须满足 $y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b) \geq 1$ ，而是只要求：

$$y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b) \geq 1 - \xi_n$$

其中 ξ_n 是一个 **松弛变量**，并且 $\xi_n \geq 0$



- 目标不再仅仅是最小化 $\frac{1}{2} \|\mathbf{w}\|^2$ ，还需要最小化 ξ_n 的总和，从而得到以下目标函数：

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

其中 C 用于控制相对重要性

- 现在的优化问题变为：

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

$$\text{s. t.: } y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b) \geq 1 - \xi_n,$$

$$\xi_n \geq 0, \quad \text{for } n = 1, 2, \dots, N$$

- 使用之前相同的方法，可以推导出对偶形式为：

$$\max_{\mathbf{a}} g(\mathbf{a})$$

$$\text{s. t.: } a_n \geq 0, \text{ } a_n \leq C$$

当 $a_n > C$ 时，可以证明 $g(\mathbf{a}) = -\infty$

$$\sum_{n=1}^N a_n y^{(n)} = 0$$

其中 $g(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m y^{(n)} y^{(m)} \mathbf{x}^{(n)T} \mathbf{x}^{(m)}$

- 当得到最优的 \mathbf{w}^* 和 b^* 时, 样本 \mathbf{x} 被分类为:

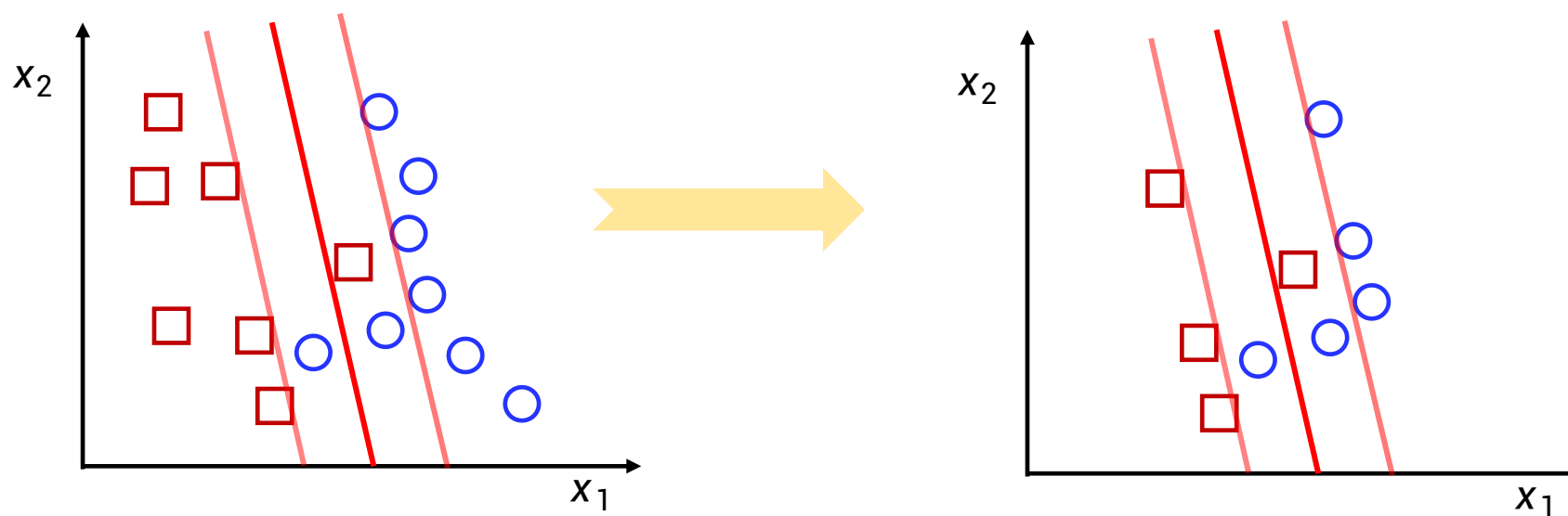
$$\hat{y}(\mathbf{x}) = \text{sign}(\mathbf{w}^{*T} \mathbf{x} + b^*)$$

- 当得到最优的 \mathbf{a}^* 时, 样本 \mathbf{x} 被分类为:

$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N a_n^* y^{(n)} \mathbf{x}^{(n)T} \mathbf{x} + b^* \right)$$

此外, 上面的两个分类器是 **等价的**

- 最优解 \mathbf{a}^* 是稀疏的，只有间隔边界内的元素是非零的



课程大纲

- 线性分类器的决策边界
- 线性最大间隔分类器
- 软线性最大间隔分类器
- 支持向量机
- 与逻辑回归的关系

非线性化

- 到目前为止，最大间隔分类器仍然是线性的
- 为了将模型非线性化，我们可以通过基函数将原始数据 \mathbf{x} 变换到特征空间

$$\phi: \mathbf{x} \rightarrow \phi(\mathbf{x})$$

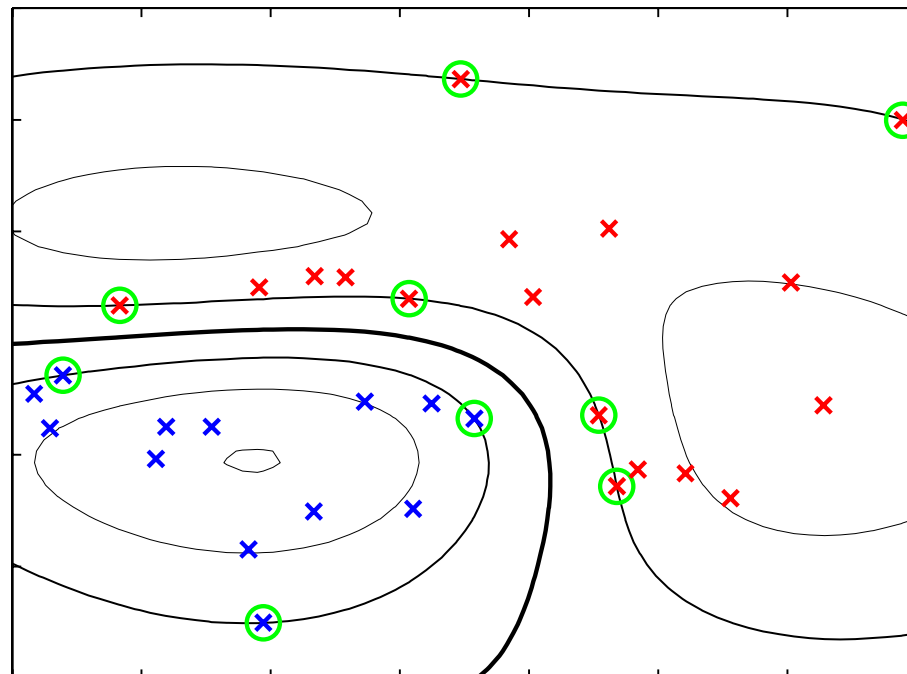
- 那么，原最大间隔优化问题变为：

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

$$\text{s. t. : } y^{(n)} \cdot (\mathbf{w}^T \phi(\mathbf{x}^{(n)}) + b) \geq 1 - \xi_n,$$

$$\xi_n \geq 0, \quad \text{for } n = 1, 2, \dots, N$$

分类器: $\hat{y}(\mathbf{x}) = \text{sign}(\mathbf{w}^{*T} \boldsymbol{\phi}(\mathbf{x}^{(n)}) + b^*)$



- 直观上，数据在高维空间中更容易被分离
- 为了获得更好的性能，我们应该将变换后特征空间 $\phi(\mathbf{x}^{(n)})$ 的维度设置得尽可能高
- 然而，基函数 $\phi(\mathbf{x})$ 的维度不能设置得太高，因为原问题的计算成本会变得非常昂贵

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{s. t.} \quad & y^{(n)} \cdot (\mathbf{w}^T \phi(\mathbf{x}^{(n)}) + b) \geq 1 - \xi_n, \\ & \xi_n \geq 0, \quad \text{for } n = 1, 2, \dots, N \end{aligned}$$

- 该问题可以通过其对偶形式求解

$$\max_{\mathbf{a}} g(\mathbf{a})$$

$$\text{s.t.: } a_n \geq 0, \quad a_n \leq C$$

$$\sum_{n=1}^N a_n y^{(n)} = 0$$

$$\text{其中 } g(\mathbf{a}) = \underbrace{\sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m y^{(n)} y^{(m)} \boldsymbol{\phi}(\mathbf{x}^{(n)})^T \boldsymbol{\phi}(\mathbf{x}^{(m)})}_{= \mathbf{1}^T \mathbf{a} - \frac{1}{2} \mathbf{a}^T \mathbf{M} \mathbf{a}}$$

$$\text{分类器: } \hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N a_n^* y^{(n)} \boldsymbol{\phi}(\mathbf{x}^{(n)})^T \boldsymbol{\phi}(\mathbf{x}) + b^* \right)$$

- \mathbf{a} 的 **维度与 $\boldsymbol{\phi}(\cdot)$ 的维度无关**，因此对偶形式能够在非常大的特征空间 $\boldsymbol{\phi}(\cdot)$ 中工作

对偶形式需要计算内积

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

这在高维情况下计算成本很高

这个问题可以通过使用 **核技巧** 来解决

核函数

- 核函数是一个双变量函数 $k(\mathbf{x}, \mathbf{x}')$ ，可以表示为某个函数 $\phi(\cdot)$ 的内积

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$$

显然， $\mathbf{x}^T \mathbf{x}'$ 和 $\phi(\mathbf{x})^T \phi(\mathbf{x}')$ 是核函数

- 默瑟定理：** 如果一个函数 $k(\mathbf{x}, \mathbf{x}')$ 是对称正定的，即

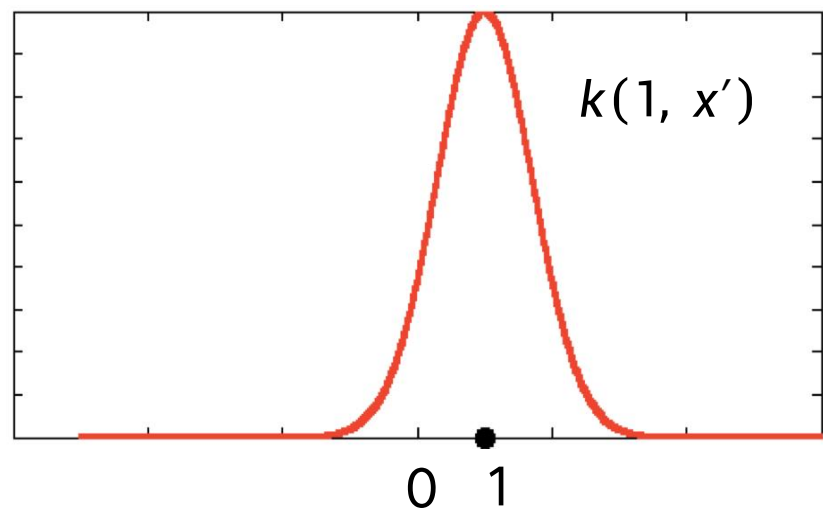
$$\int \int g(\mathbf{x}) k(\mathbf{x}, \mathbf{y}) g(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0, \quad \forall g(\cdot) \in L^2,$$

那么必然存在一个函数 $\phi(\cdot)$ 使得 $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$

如果一个函数 $k(\mathbf{x}, \mathbf{x}')$ 满足对称正定条件，那么它必然是一个核函数

- 最广泛使用的核函数之一是高斯核，其形式为：

$$k(\mathbf{x}, \mathbf{x}') = \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2 \right\}$$



➤ 高斯核的函数 $\phi(\cdot)$ 具有无限维度

$$\phi(x) = e^{-x^2/2\sigma^2} \left[1, \sqrt{\frac{1}{1!\sigma^2}}x, \sqrt{\frac{1}{2!\sigma^4}}x^2, \sqrt{\frac{1}{3!\sigma^6}}x^3, \dots \right]^T$$

- 借助核函数，对偶最大间隔分类器可以等价地重写为：

$$\max_{\mathbf{a}} g(\mathbf{a})$$

$$s. t. : \quad a_n \geq 0, \quad a_n \leq C$$

$$\sum_{n=1}^N a_n y^{(n)} = 0$$

$$\text{其中 } g(\mathbf{a}) = \underbrace{\sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m y^{(n)} y^{(m)} k(\mathbf{x}^{(n)}, \mathbf{x}^{(m)})}_{= \mathbf{1}^T \mathbf{a} - \frac{1}{2} \mathbf{a}^T \mathbf{M} \mathbf{a}}$$

- 导出的分类器

$$\hat{y}(\mathbf{x}) = \text{sign} \left(\sum_{n=1}^N a_n^* y^{(n)} k(\mathbf{x}^{(n)}, \mathbf{x}) + b^* \right)$$

核技巧：用核函数 $k(\mathbf{x}, \mathbf{x}')$ 替代 $\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}')$

- 结论可以总结如下：
 - 如果 $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$ ，这是一个线性最大间隔分类器
 - 如果 $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}')$ ，这是一个基于基函数的 **有限维** 非线性最大间隔分类器
 - 如果 $k(\mathbf{x}, \mathbf{x}') = \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2 \right\}$ ，这是一个 **无限维** 非线性最大间隔分类器

课程大纲

- 线性分类器的决策边界
- 线性最大间隔分类器
- 软线性最大间隔分类器
- 支持向量机
- 与逻辑回归的关系

- 在逻辑回归中，我们最小化损失：

$$L(\mathbf{w}, b) = - \sum_{n=1}^N \left[\tilde{y}^{(n)} \log \sigma(h^{(n)}) + (1 - \tilde{y}^{(n)}) \log (1 - \sigma(h^{(n)})) \right] + \lambda \|\mathbf{w}\|^2$$

$$= \sum_{n=1}^N \log (1 + \exp (-y^{(n)} h^{(n)})) + \lambda \|\mathbf{w}\|^2$$

注： $\tilde{y} \in \{0, 1\}$, $y \in \{-1, 1\}$

$$= \sum_{n=1}^N E_{LR}(y^{(n)} h^{(n)}) + \lambda \|\mathbf{w}\|^2$$

其中， $E_{LR}(z) = \log (1 + \exp (-z))$

- 在理想分类器中，我们最小化损失

$$L(\mathbf{w}, b) = \sum_{n=1}^N E_{Ideal}(y^{(n)}h^{(n)}) + \lambda \|\mathbf{w}\|^2$$

注： $y \in \{-1, 1\}$,
 $h \in \{-\infty, \infty\}$

其中，如果 $z \geq 0$ ， $E_{Ideal}(z) = 0$ ；否则为 1

公式定义： $E_{Ideal}(z) = 0$ if $z \geq 0$ else 1

阶跃函数，不可导

- 回顾“线性最大间隔分类器”：

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

$$\text{s. t. : } y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b) \geq 1, \text{ for } n = 1, 2, \dots, N$$

- 在线性最大间隔分类器中，我们等价地最小化损失：

$$L(\mathbf{w}, b) = \sum_{n=1}^N E_{\infty}(y^{(n)} h^{(n)} - 1) + \frac{1}{2} \|\mathbf{w}\|^2$$

其中，如果 $z \geq 0$ ，则 $E_{\infty}(z) = 0$ ；否则为 $+\infty$

硬约束：所有样本都满足约束，否则损失无穷大

公式定义： $E_{\infty}(z) = 0 \text{ if } z \geq 0 \text{ else } +\infty$

也为阶跃函数，可以发现仍然不可导

- 回顾“软线性最大间隔分类器”：

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

$$\text{s. t.: } y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b) \geq 1 - \xi_n,$$

$$\xi_n \geq 0, \quad \text{for } n = 1, 2, \dots, N$$

- 从约束推导损失函数：

$$\xi_n \geq 1 - y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b),$$

$$\xi_n \geq 0,$$

$$\xi_n = \max(0, 1 - y^{(n)} \cdot (\mathbf{w}^T \mathbf{x}^{(n)} + b))$$

注： $\max(0, 1 - z)$

- 在软线性最大间隔分类器中，我们等价地最小化损失：

$$L(\mathbf{w}, b) = C \sum_{n=1}^N E_{SV}(y^{(n)} h^{(n)}) + \frac{1}{2} \|\mathbf{w}\|^2$$

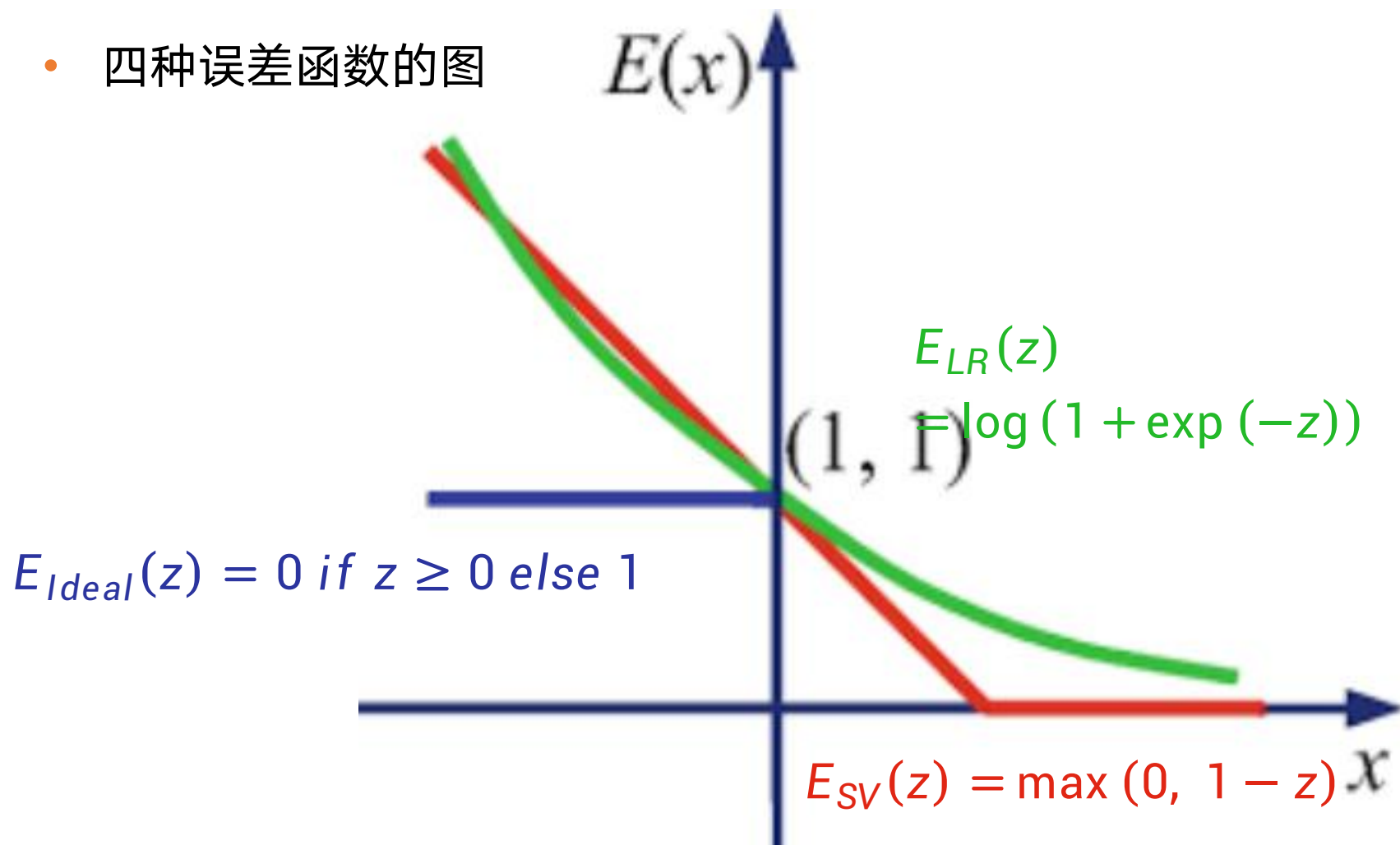
$$= \sum_{n=1}^N E_{SV}(y^{(n)} h^{(n)}) + \lambda \|\mathbf{w}\|^2$$

其中， $E_{SV}(z) = \max(0, 1 - z)$ ，这被称为 **合页损失 (hinge loss)**

理想损失像“台阶”——对了 0，错了 1，不可导；
而在软线性最大间隔分类器里，我们用 hinge 函数把它磨平：
 $\max(0, 1 - z)$ —— **训练时可导，测试时近似理想**

- 我们可以看到，这四种分类器可以在同一个框架下进行建模，**唯一的区别在于所选择的误差函数**

- 四种误差函数的图



注： $E_{SV}(z)$ 类似 ReLU
训练时可导：
可以使用次梯度求导

$$f'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$