

明理精工

笃学致远

第5章 支持向量机

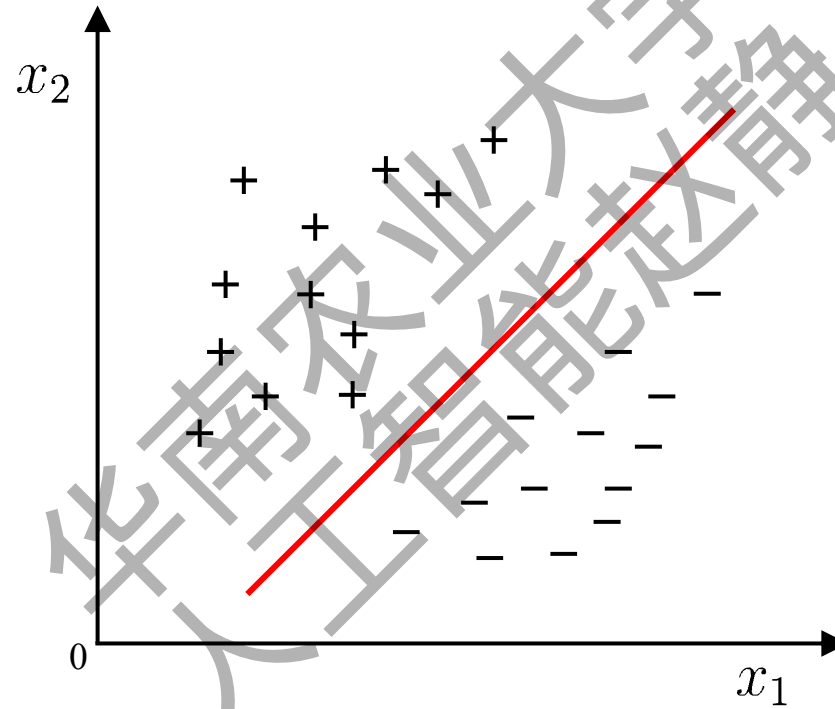
Support Vector Machine



电子工程学院、人工智能学院

college of Electronic Engineering , college of Artificial Intelligence

线性模型：在样本空间中寻找一个超平面，将不同类别的样本分开。

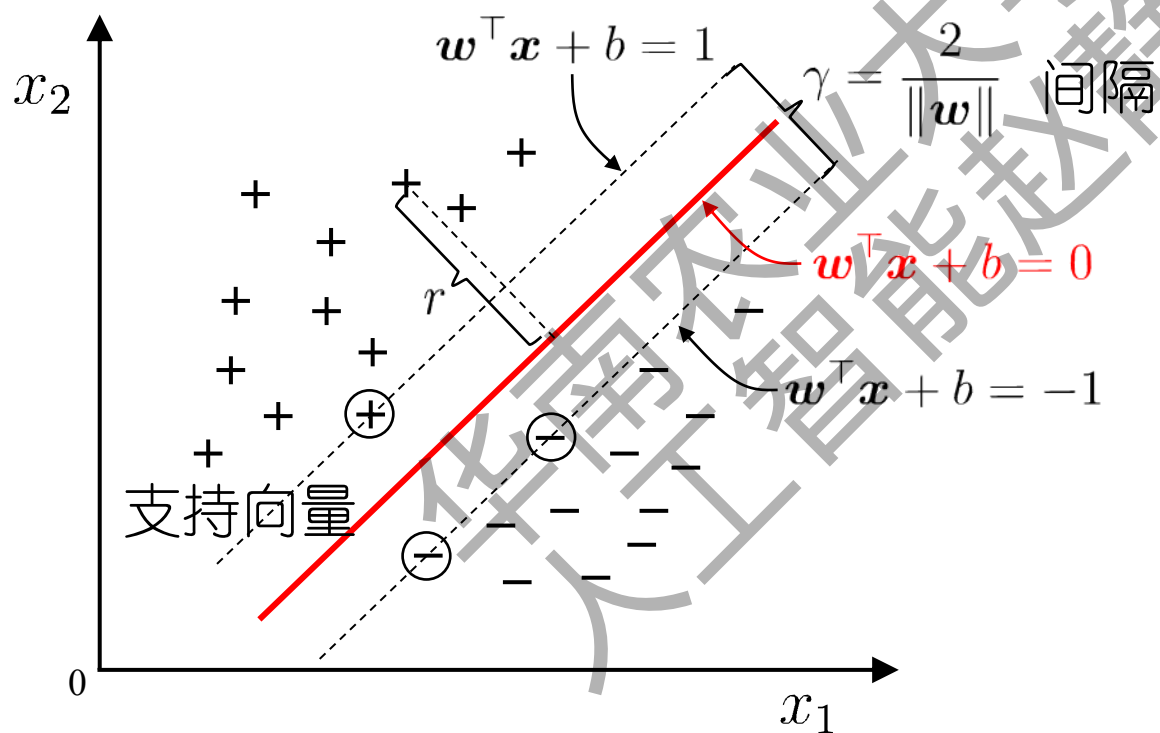




- ◆ 线性SVM
- ◆ 带松弛变量的SVM
- ◆ 合页损失函数
- ◆ SVM的对偶问题
- ◆ 核化SVM模型
- ◆ SVM回归 (SVR)

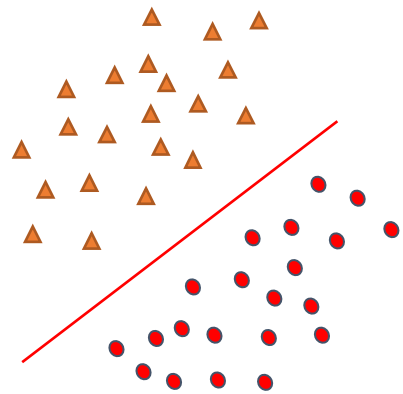
1. 线性SVM

$$w^T x + b = 0$$

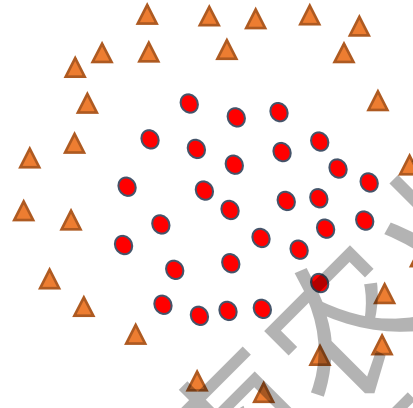


Vladimir Vapnik

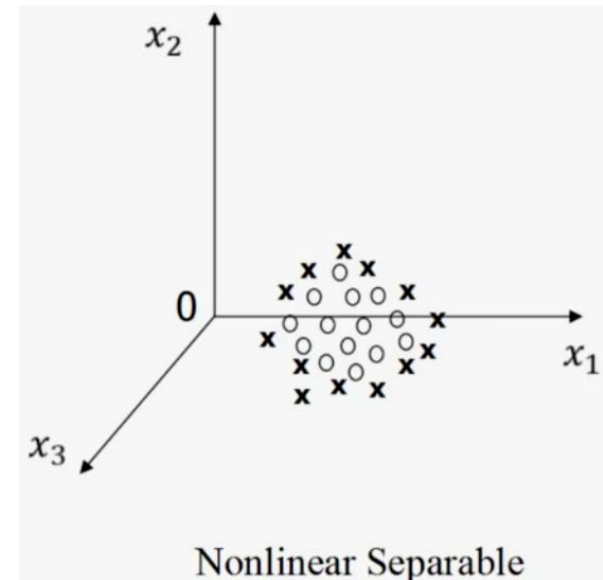
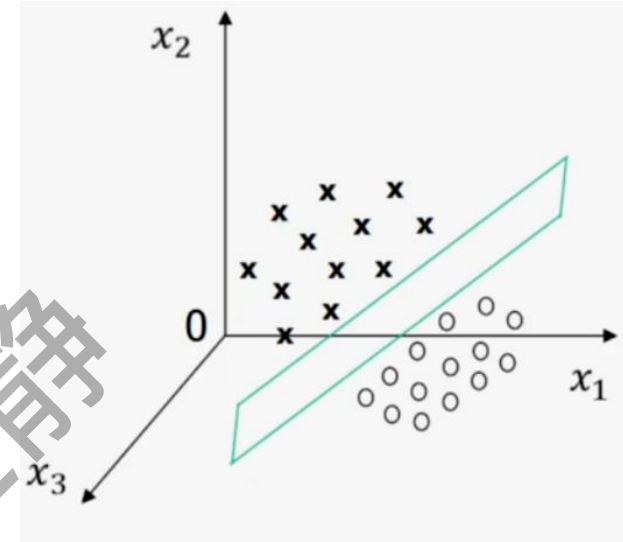
➤ 线性可分与线性不可分



线性可分



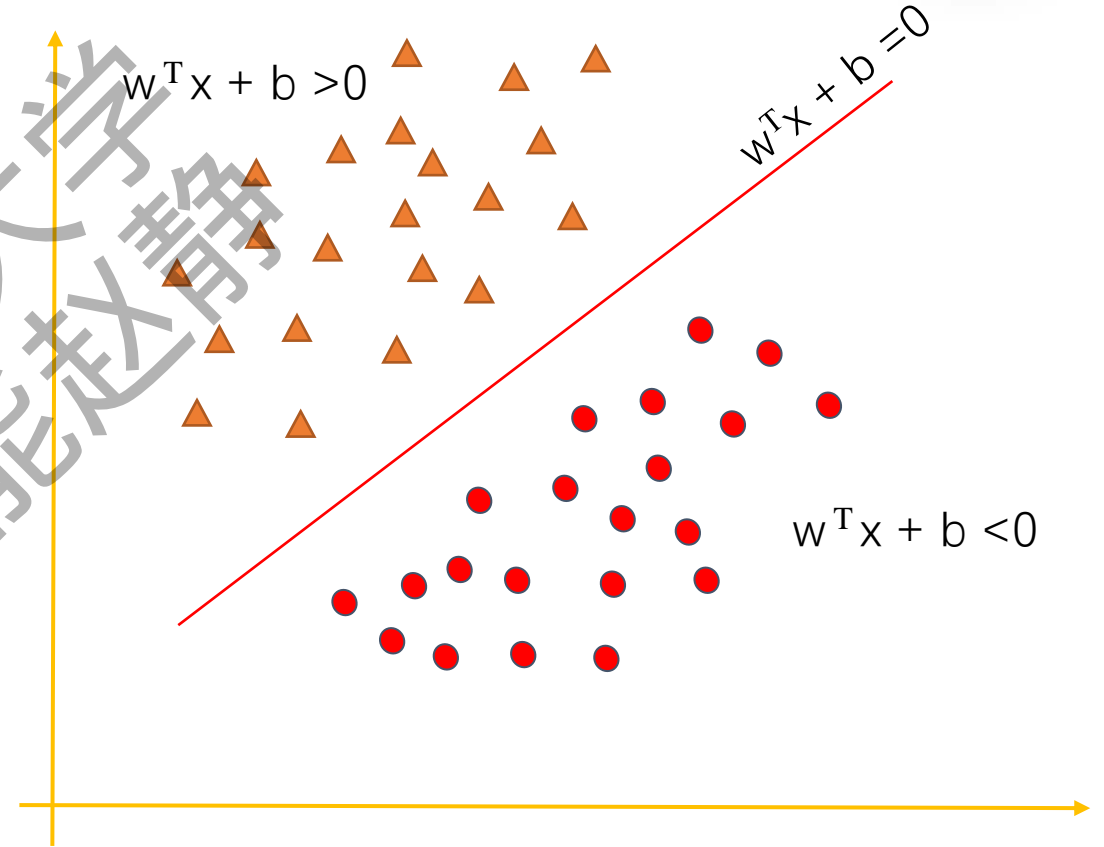
线性不可分



➤ 定义线性可分

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$
$$y \in \{+1, -1\}$$

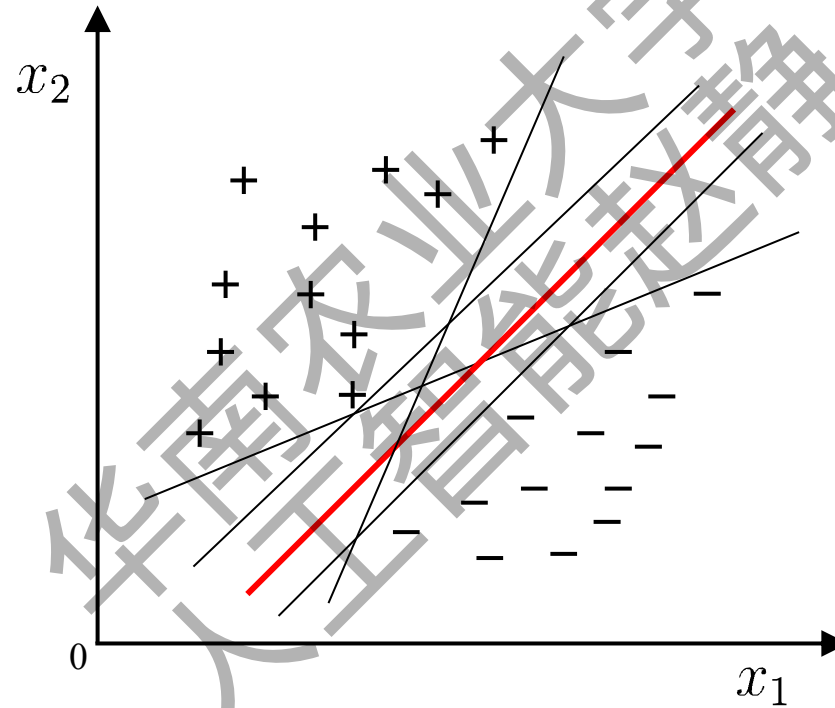
- 线性分类器 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$
- 若存在 w ，使得：
所有满足 $f(\mathbf{x}) < 0$ 的点，其对应的 y 等于 -1
所有满足 $f(\mathbf{x}) > 0$ 的点，其对应的 y 等于 1
则数据线性可分
- 线性判别函数 $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$ ， \mathbf{x} 是位于超平面面上的点



超平面： $\mathbf{w}^T \mathbf{x} + b = 0$

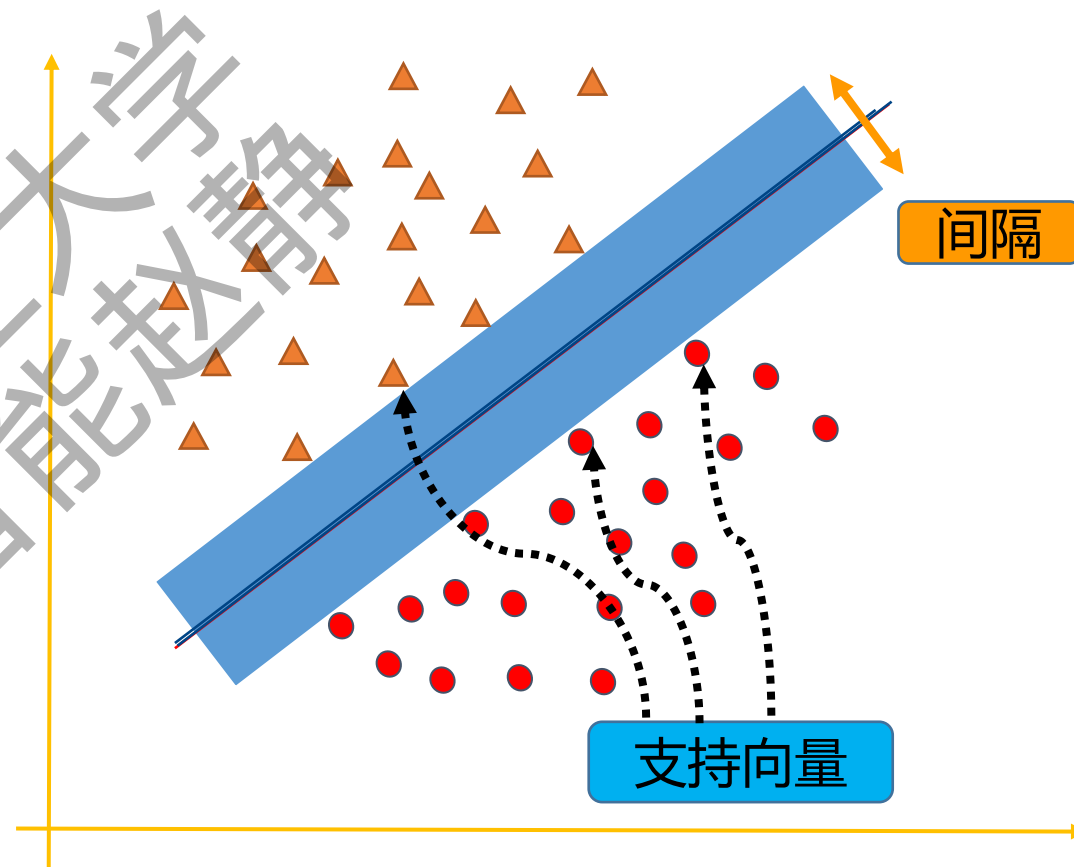
➤ 线性可分支持向量机

-Q: 将训练样本分开的超平面可能有很多, 哪一个好呢?

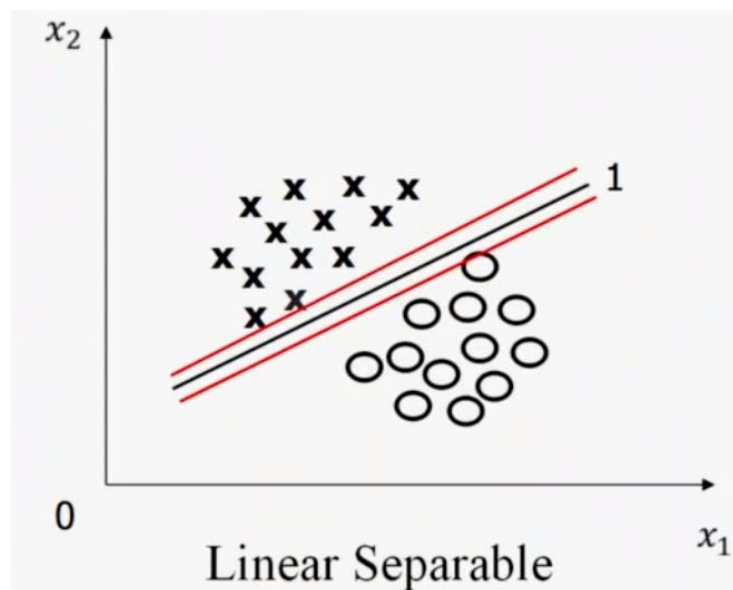


-A: 应选择“**正中间**”, 容忍性好, 鲁棒性高, 泛化能力最强.

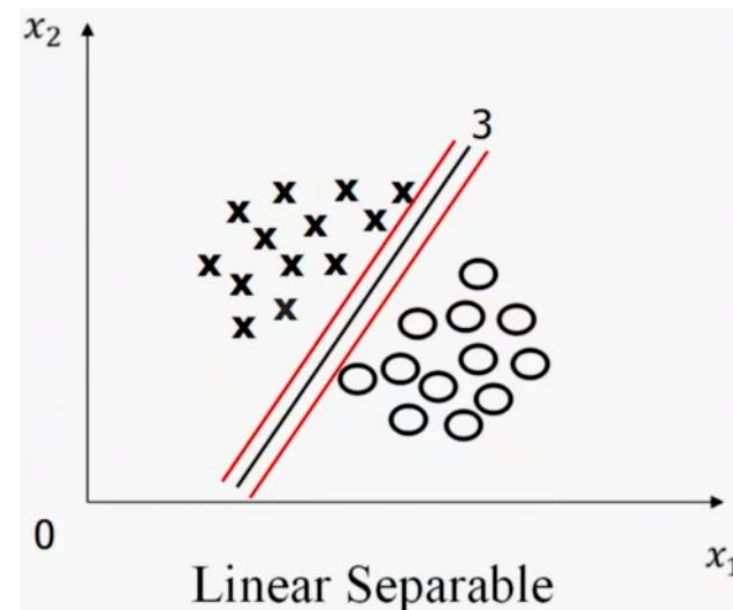
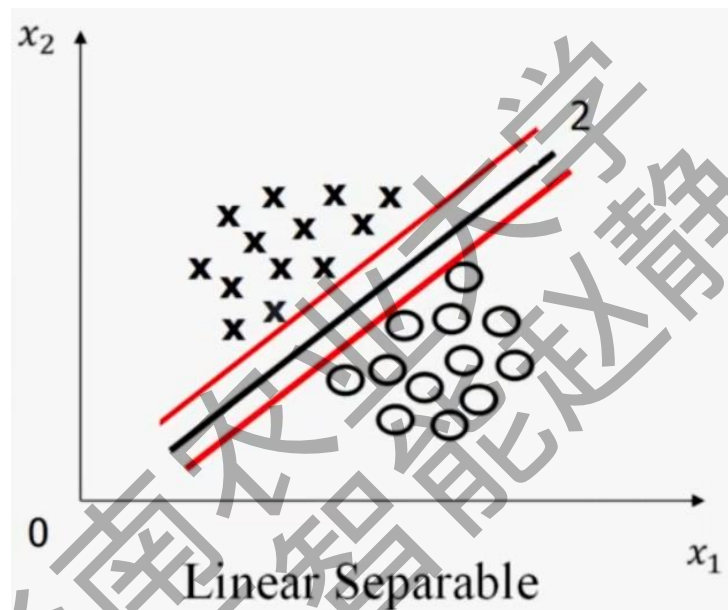
找到集合边缘上的若干数据（称为支持向量（Support Vector）），用这些点找出一个平面（称为决策面），使得支持向量到该平面的距离最大。



➤ 间隔 (Margins) 最大



直线能分开两类
寻找间隔最大的直线
直线在间隔的正中间

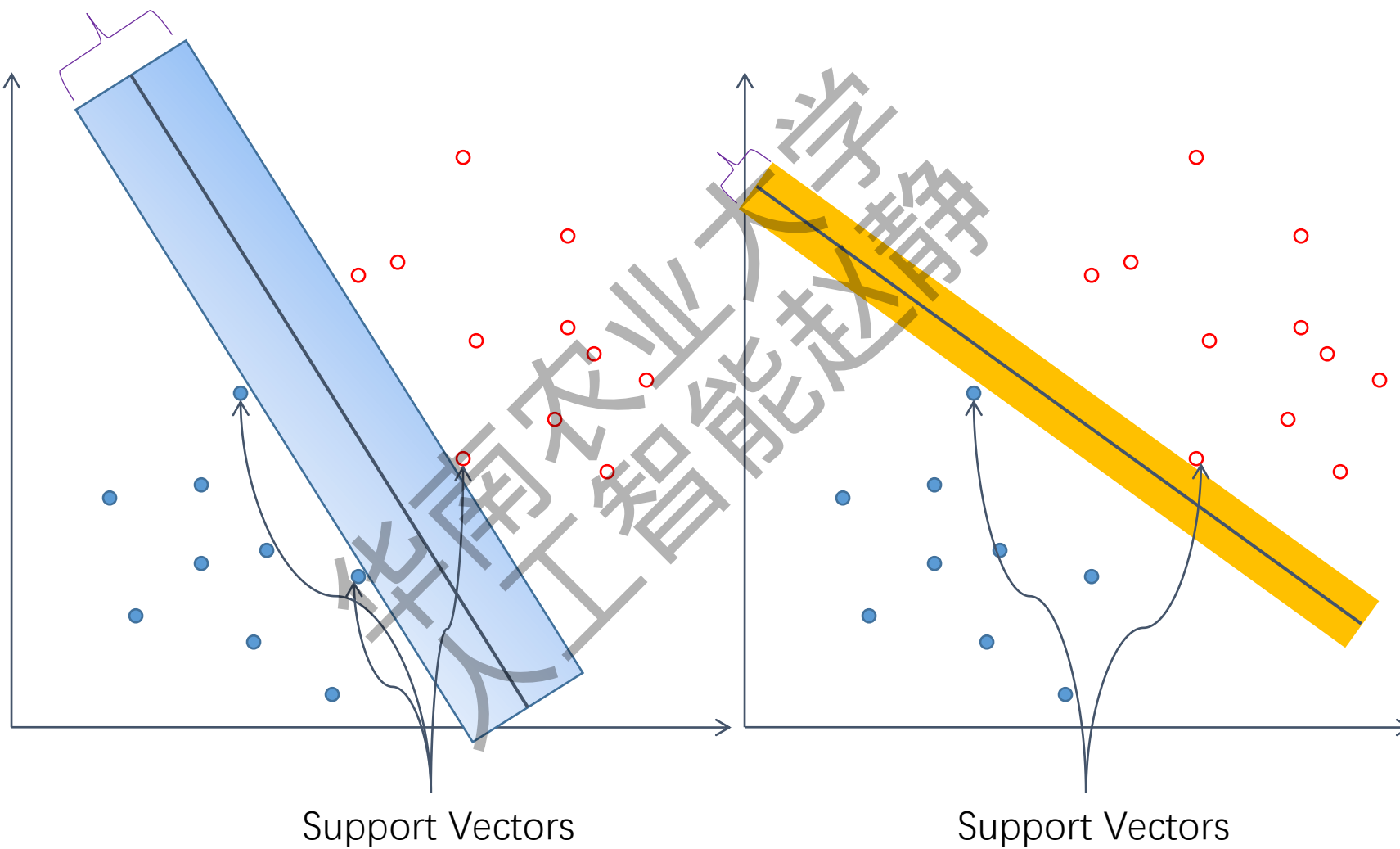


$$\gamma = \frac{2}{\|w\|}$$



超平面能分开两类
寻找间隔最大的超平面
超平面在间隔的正中间

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




➤ 线性可分支持向量机最优化问题

- 最大间隔: 寻找参数 w 和 b , 使得 γ 最大.

$$\operatorname{argmax}_{w,b} \frac{1}{\|w\|}$$

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


$$\operatorname{argmin}_{w,b} \frac{1}{2} \|w\|^2$$

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

$$w = \begin{bmatrix} \omega_1 \\ \omega_2 \\ \vdots \\ \omega_m \end{bmatrix}$$

$$\|w\|^2 = \omega_1^2 + \omega_2^2 \dots + \omega_m^2 = \sum_{i=1}^m \omega_i^2$$

- 凸二次规划(convex quadratic programming)

点到平面的距离：

$$w^T x + b = 0$$

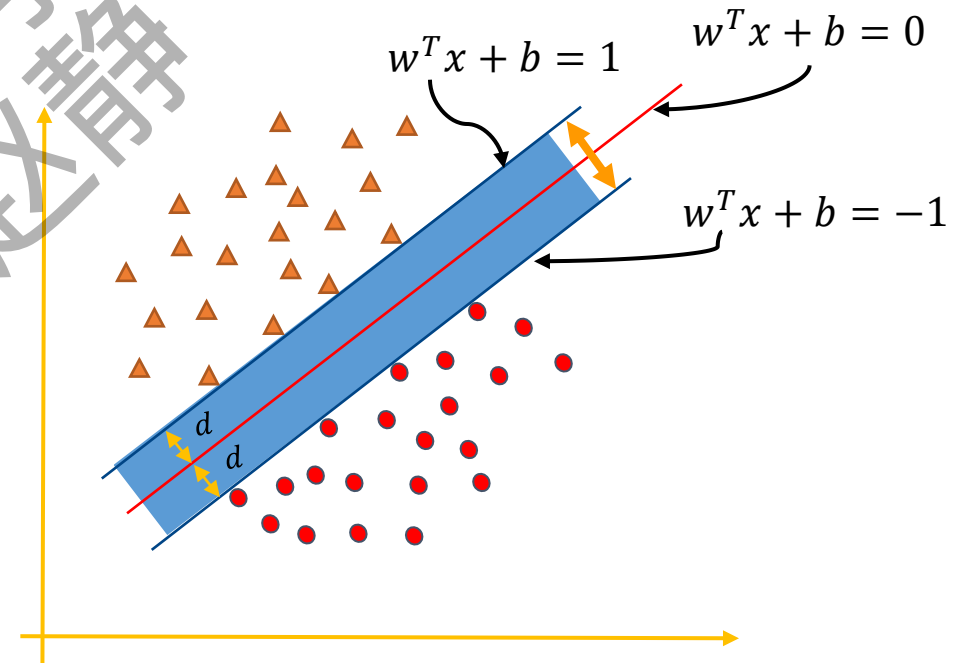
二维空间点 (x, y) 到直线 $Ax + By + C = 0$ 的距离公式是：

$$\frac{|Ax + By + C|}{\sqrt{A^2 + B^2}}$$

扩展到 n 维空间后，点 $x = (x_1, x_2 \dots x_n)$ 到超平面

$$w^T x + b = 0 \text{ 的距离为: } d = \frac{|w^T x + b|}{\|w\|}$$

$$\text{其中 } \|w\| = \sqrt{w_1^2 + \dots w_n^2}$$



决策超平面:

$w^T x + b = 0$ 与 $(aw^T)x + ab = 0$ ($a \neq 0$) 是同一个超平面

$$(w, b) \longrightarrow (aw, ab)$$

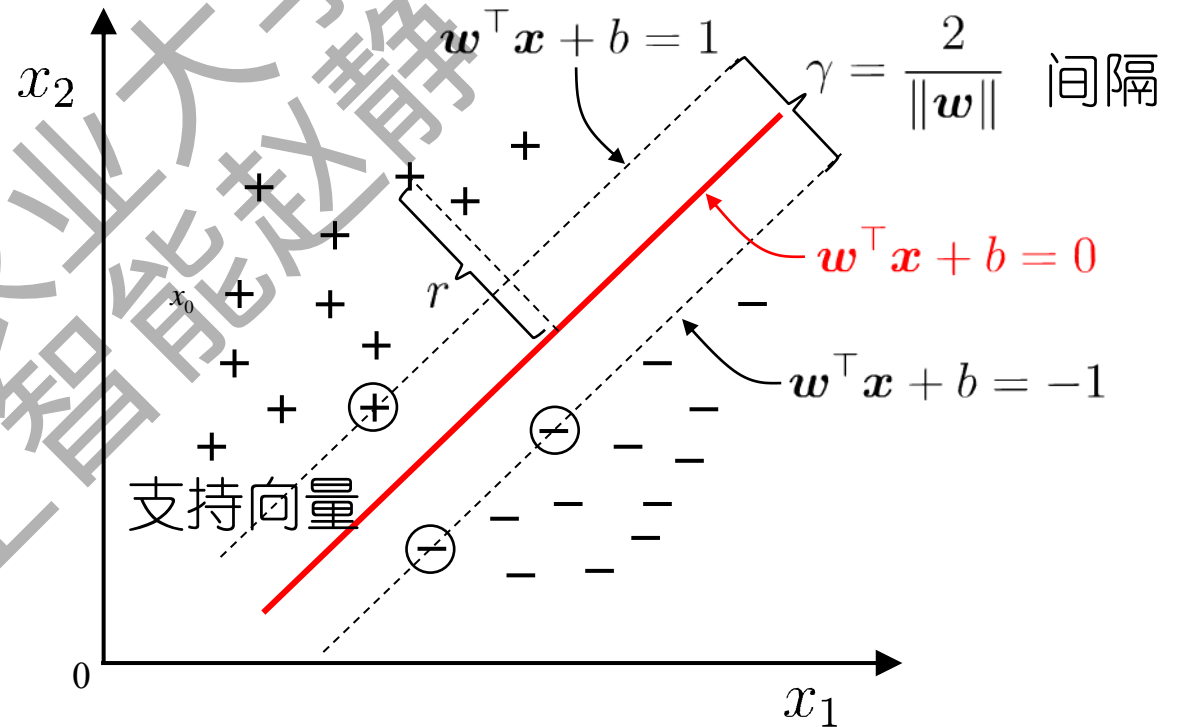
支持向量决定的超平面:

支持向量 x_0 $w^T x_0 + b = 1$

非支持向量 x_0 $w^T x_0 + b > 1$

支持向量 x 到超平面的距离

$$d = \frac{|w^T x + b|}{\|w\|} = \frac{1}{\|w\|}$$



- 支持向量 x 到超平面的距离

$$d = \frac{|w^T x + b|}{\|w\|} = \frac{1}{\|w\|}$$

最大化 $\frac{1}{\|w\|}$, 即最小化 $\|w\|$ \longrightarrow $\underset{w,b}{\operatorname{argmin}} \frac{1}{2} \|w\|^2$

- 非支持向量 x 到超平面的距离

$$\begin{cases} w^T x_i + b \geq +1, & y_i = +1; \\ w^T x_i + b \leq -1, & y_i = -1. \end{cases}$$



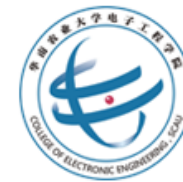
$$s.t. \ y_i(w^T x_i + b) \geq 1, \ i = 1, 2, \dots, m$$

$$\underset{w,b}{\operatorname{argmin}} \frac{1}{2} \|w\|^2$$

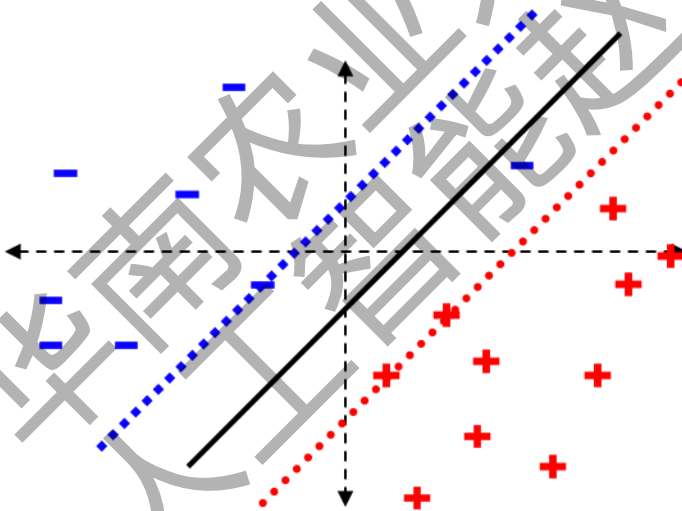
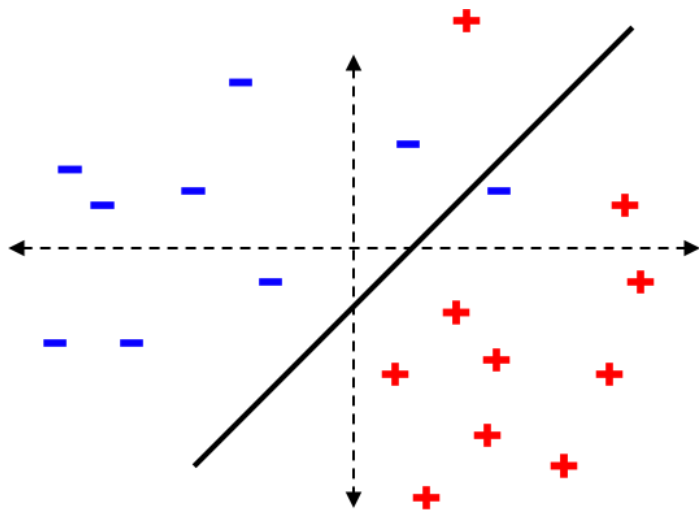
$$s.t. \ y_i(w^T x_i + b) \geq 1, \ i = 1, 2, \dots, m$$



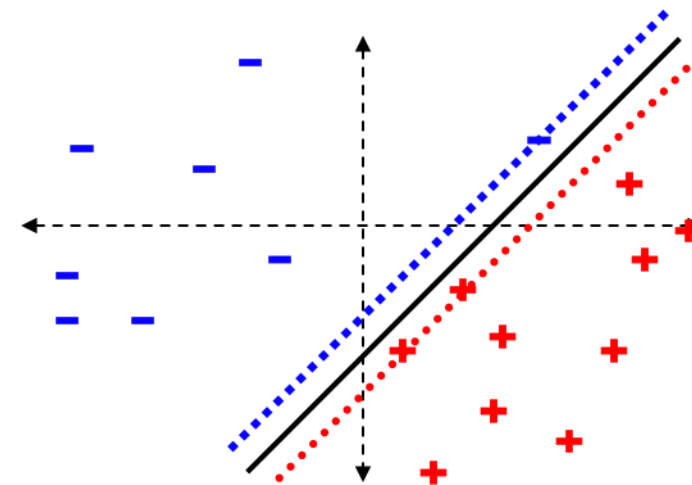
- ◆ 线性SVM
- ◆ 带松弛因子的SVM
- ◆ 合页损失函数
- ◆ SVM的对偶问题
- ◆ 核化SVM模型
- ◆ SVM回归 (SVR)



2.带松弛因子的SVM

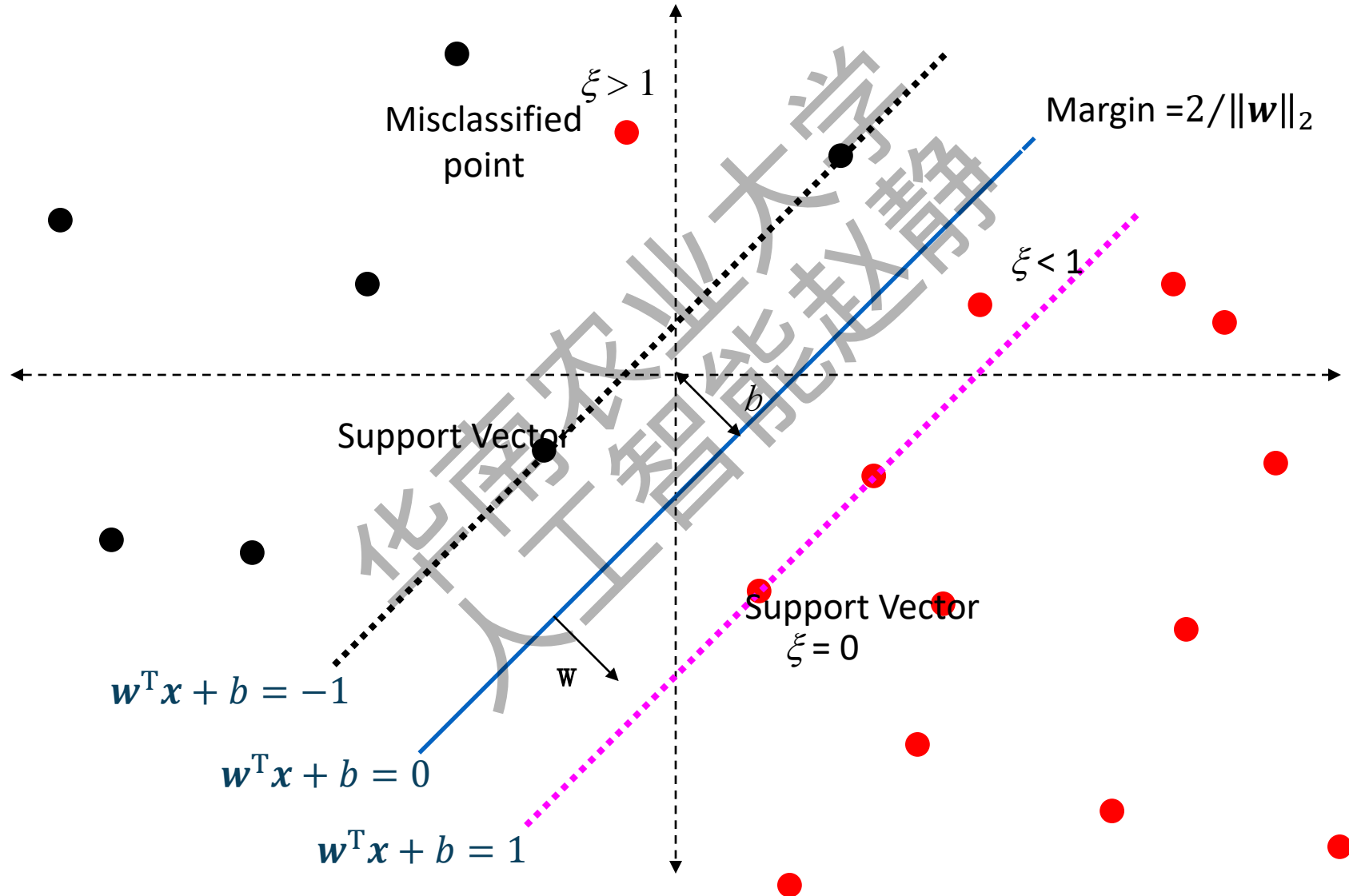


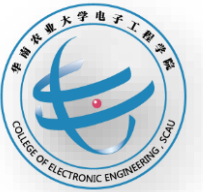
少量样本被错分，但间隔大



样本被完全分对，但间隔小

数据不完全线性可分：松弛变量 ξ





➤ C-SVM

若数据线性不可分，则可以引入松弛变量(slack variable) $\xi \geq 0$ ，使函数间隔加上“**松弛变量**”大于等于1

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

则软间隔最大化SVM (C-SVM) 的**目标函数**

$$J(\mathbf{w}, b, C) = C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$

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

$$\xi_i \geq 0, \quad i = 1, 2, \dots, N$$

➤ C-SVM目标函数

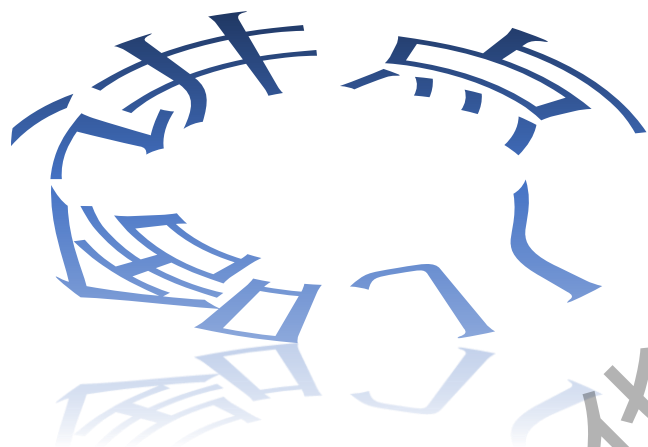
$$J(\mathbf{w}, b, C) = C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$

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

$$\xi_i \geq 0, \quad i = 1, 2, \dots, N$$

形式与带正则的线性回归或Logistic回归的目标函数类似

$$J(\mathbf{w}, \lambda) = C \sum_{i=1}^N L(y_i, f(\mathbf{x}_i, \mathbf{w})) + R(\mathbf{w})$$



- ◆ 线性SVM
- ◆ 带松弛因子的SVM
- ◆ 合页损失函数
- ◆ SVM的对偶问题
- ◆ 核化SVM模型
- ◆ SVM回归 (SVR)

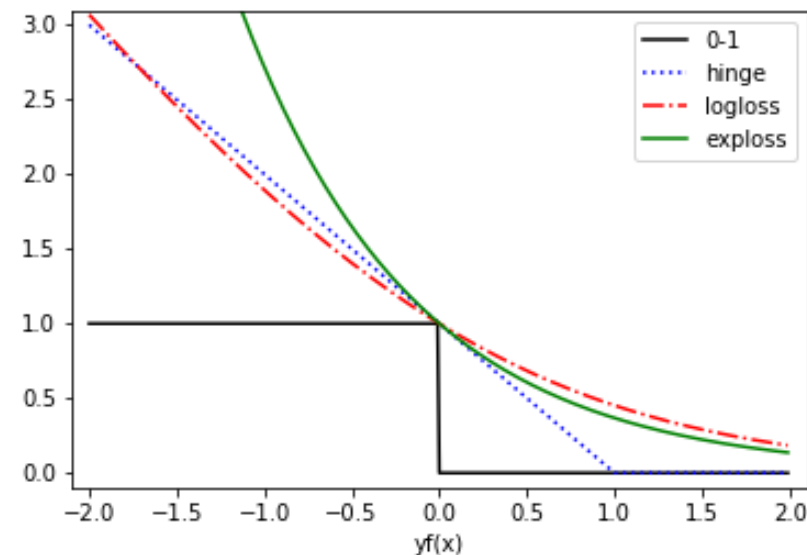
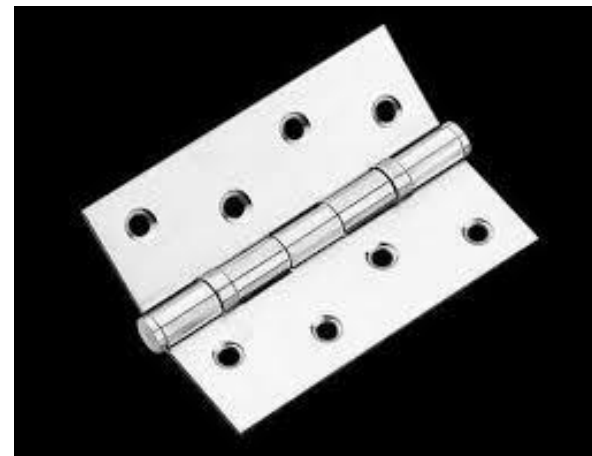
3. 合页损失函数 (Hinge Loss)

在C-SVM中

- 当 $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$, $\xi_i = 0$
- 其他点: $\xi_i = 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)$

合页损失(Hinge Loss)

$$\xi = L_{Hinge}(y, \hat{y}) = \begin{cases} 0 & y\hat{y} \geq 1 \\ 1 - y\hat{y} & \text{otherwise} \end{cases}$$



将合页损失代入C-SVM的目标函数

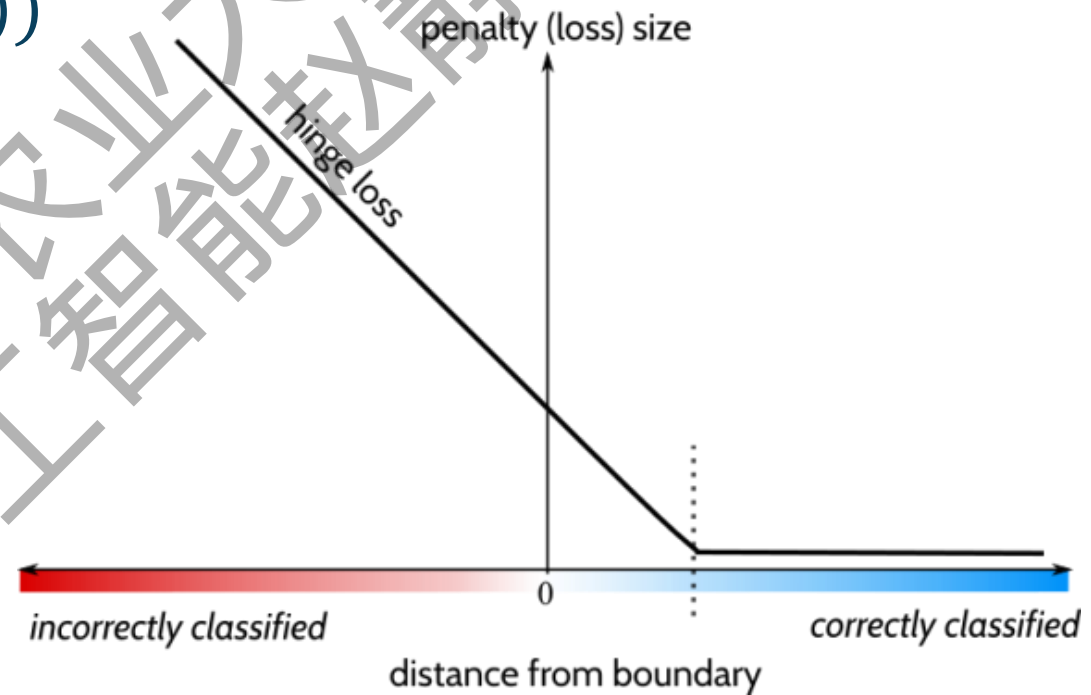
$$\begin{aligned} J(\mathbf{w}, b, C) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i \\ &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N L_{\text{Hinge}}(y_i, f(\mathbf{x}_i, \mathbf{w}, b)) \end{aligned}$$

对比一般机器学习模型的目标函数

$$J(\boldsymbol{\theta}, \lambda) = C \sum_{i=1}^N L(y_i, f(\mathbf{x}_i, \boldsymbol{\theta})) + R(\boldsymbol{\theta})$$

$$\xi = L_{Hinge}(y, \hat{y}) = \begin{cases} 0 & y\hat{y} \geq 1 \\ 1 - y\hat{y} & \text{otherwise} \end{cases}$$

$$\xi_i = \max(0, 1 - y_i(w^T x_i + b))$$

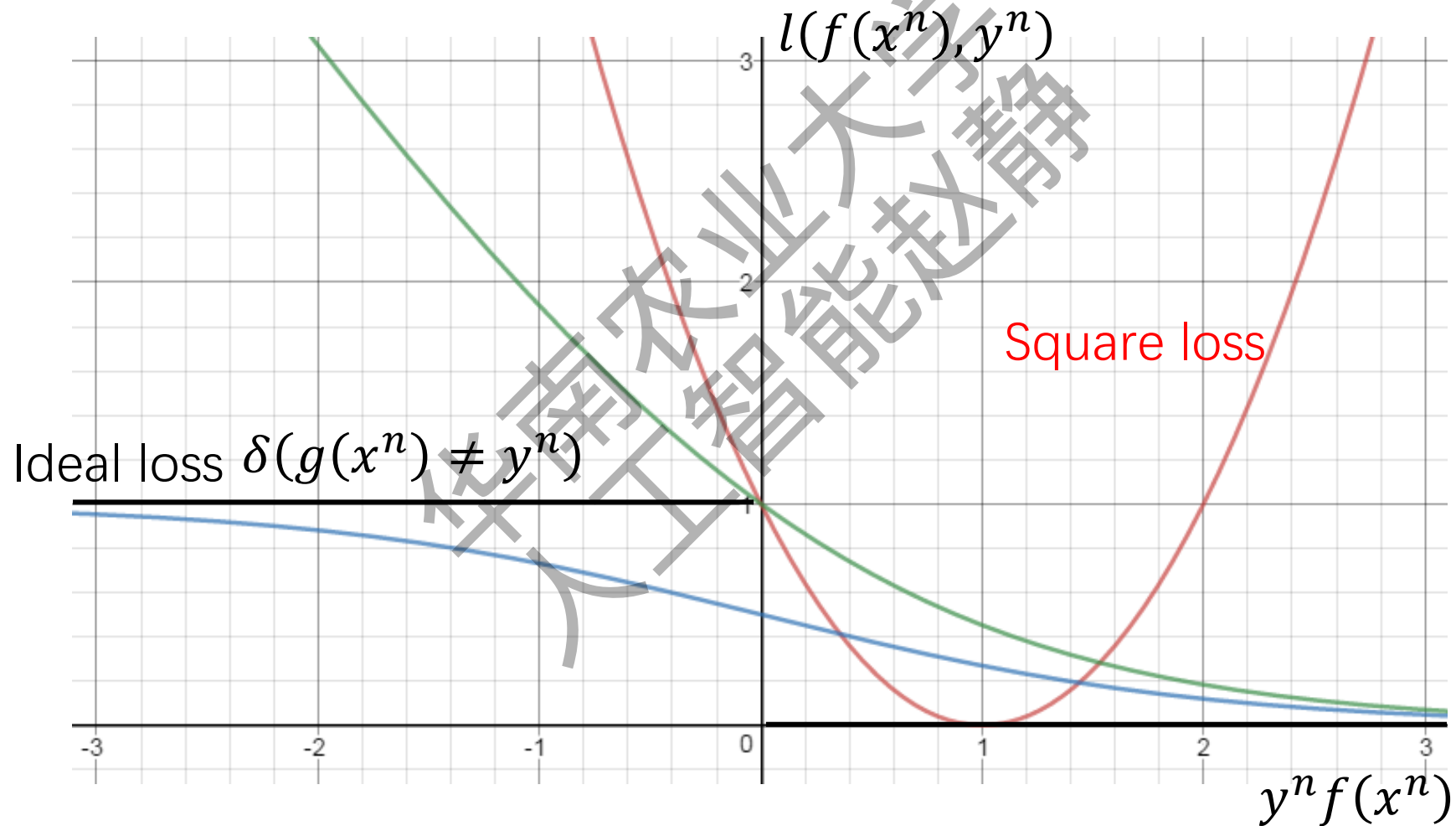


合页函数 **hinge loss**

✓ 总结：损失函数

- 平方差

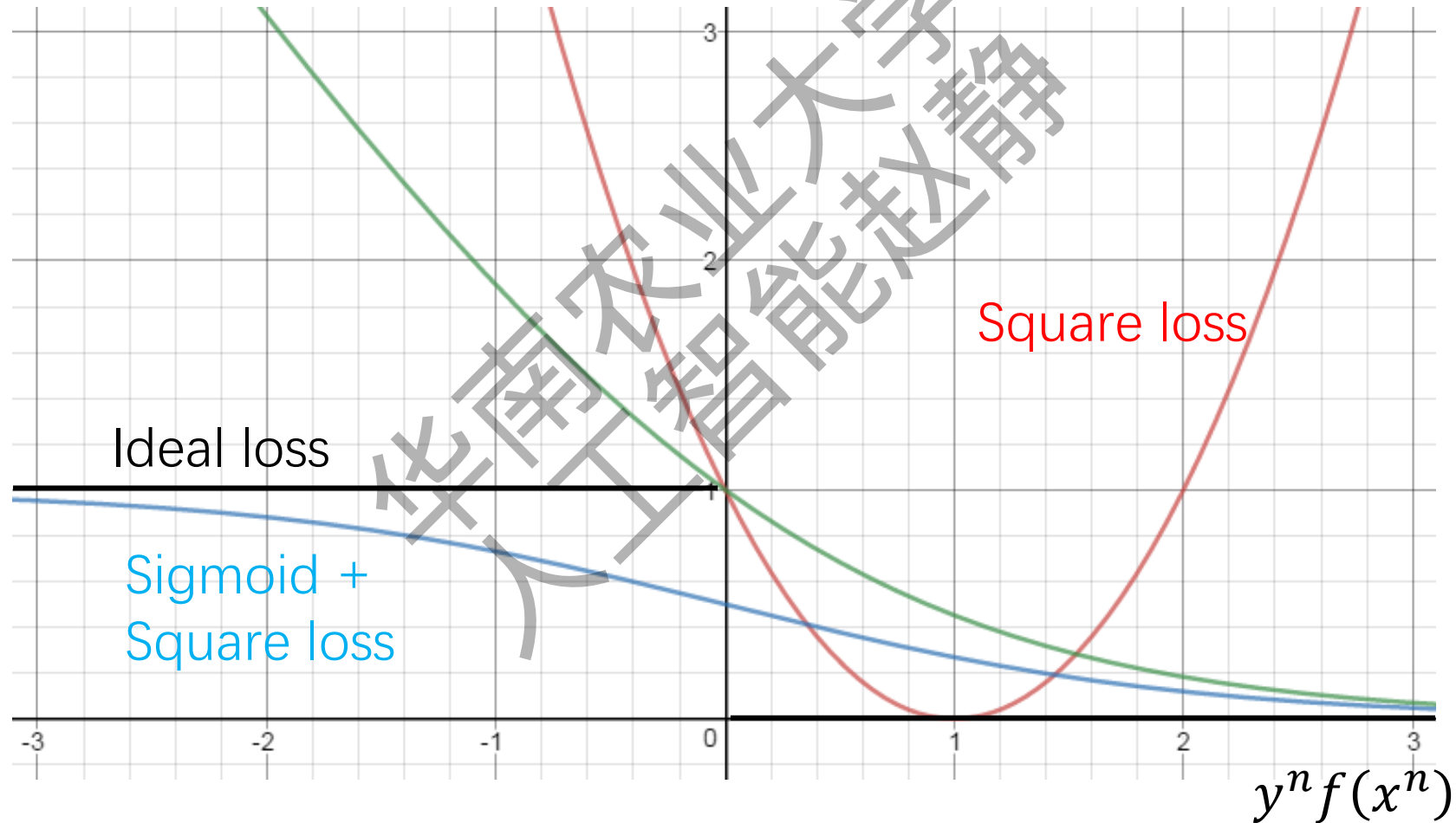
$$\hat{y} = \begin{cases} w^T x_i + b = f(x) > 0 & y = +1 \\ w^T x_i + b = f(x) < 0 & y = -1 \end{cases}$$



- Sigmoid + Square Loss:

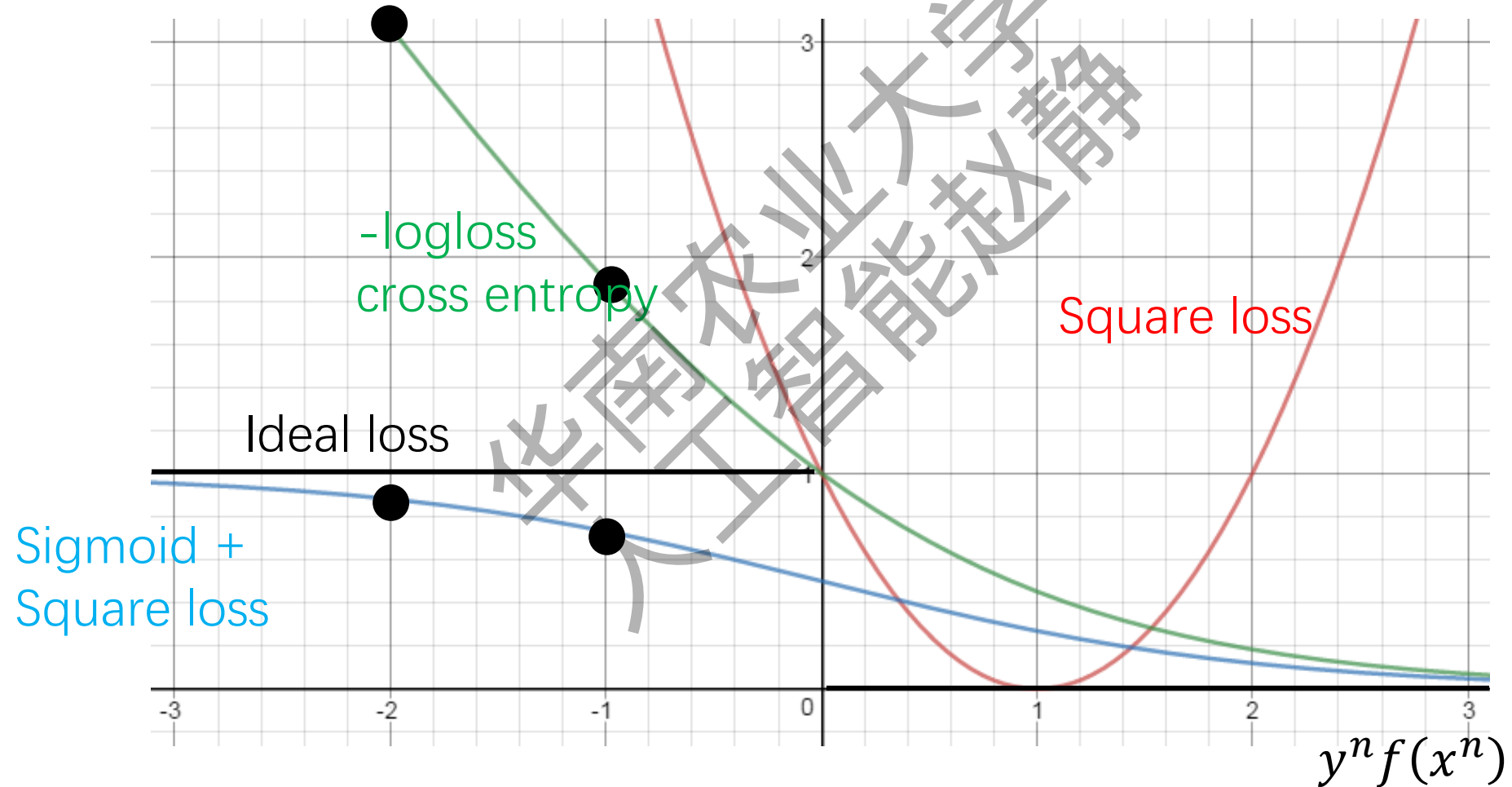
If $y^n = 1$, $\sigma(f(x))$ close to 1
If $y^n = -1$, $\sigma(f(x))$ close to 0

$$l(f(x^n), \hat{y}^n) = (\sigma(\hat{y}^n f(x^n)) - 1)^2$$



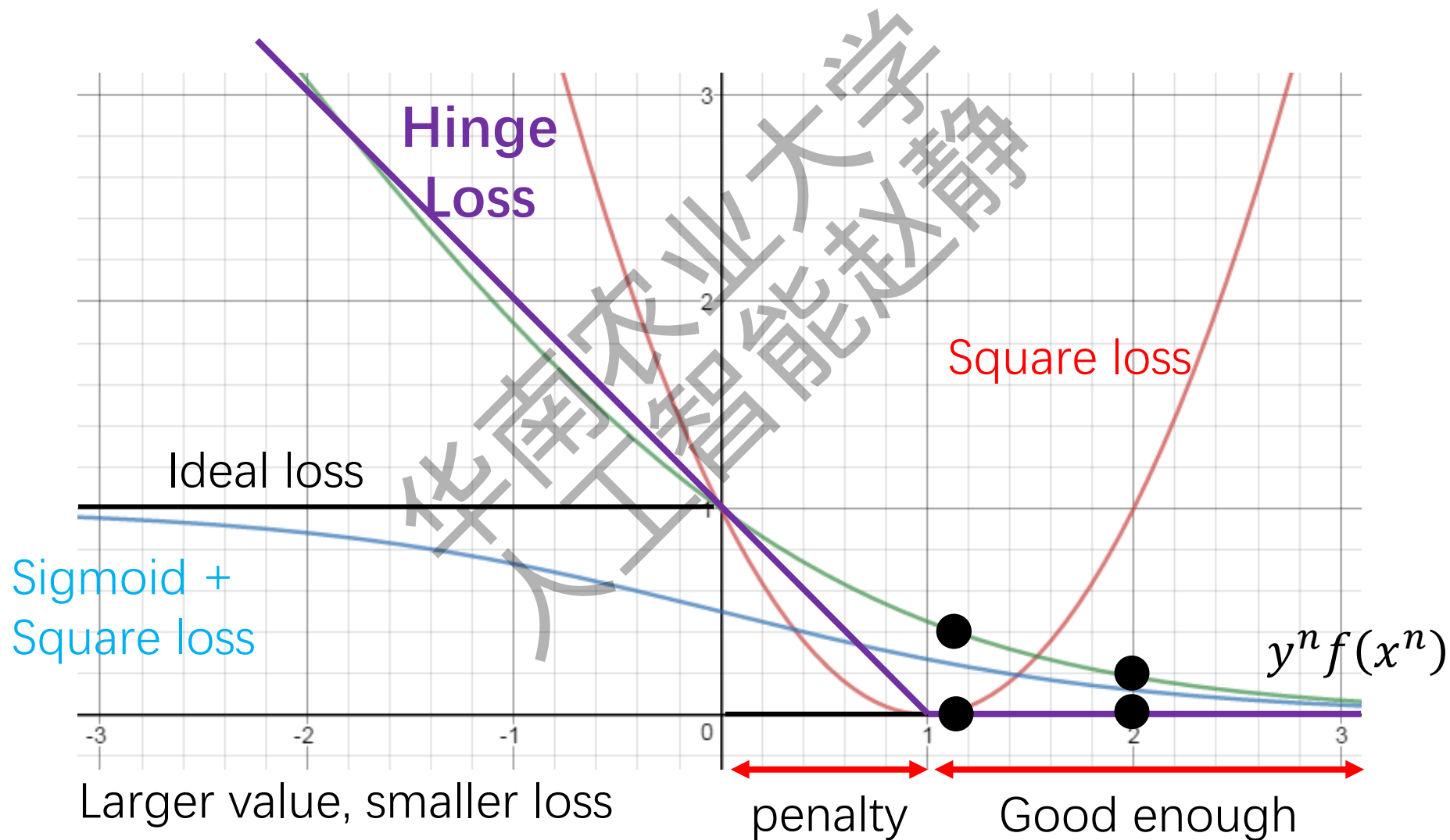
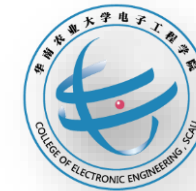
- **-logloss**
cross entropy

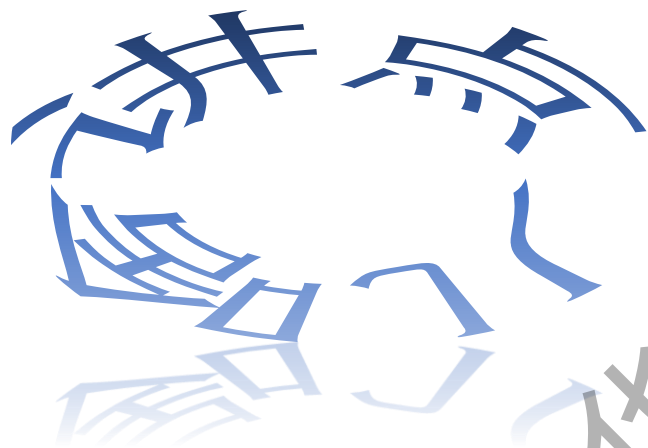
$$\begin{aligned}\ell(\mu) &= -\ln p(\mathcal{D}) \\ &= -\sum_{i=1}^N \ln p(y_i | \mathbf{x}_i) = -\sum_{i=1}^N \ln(\mu(\mathbf{x}_i)^{y_i} (1 - \mu(\mathbf{x}_i))^{(1-y_i)})\end{aligned}$$



- Hinge Loss

$$l(f(x^n), y^n) = \max(0, 1 - y^n f(x^n))$$





- ◆ 线性SVM
- ◆ 带松弛因子的SVM
- ◆ 合页损失函数
- ◆ SVM的对偶问题
- ◆ 核化SVM模型
- ◆ SVM回归 (SVR)



4. 对偶问题

- C-SVM的目标函数为

$$J(\mathbf{w}, b, C) = C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$

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

$$\xi_i \geq 0, \quad i = 1, 2, \dots, N$$

➤ 拉格朗日乘子法

- C-SVM原问题目标函数:

$$J(\mathbf{w}, b, C) = C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$
$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N$$
$$\xi_i \geq 0, \quad i = 1, 2, \dots, N$$

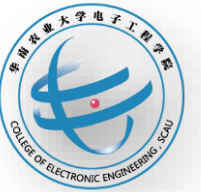
- 写成标准的不等式约束问题:

$$\min C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$
$$\text{s.t. } 1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0, \quad i = 1, 2, \dots, N$$
$$-\xi_i \leq 0, \quad i = 1, 2, \dots, N$$

- 对应的广义拉格朗日函数:

$$L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^N \mu_i \xi_i$$

$$\text{s.t. } \alpha_i \geq 0, \quad \mu_i \geq 0, \quad \xi_i \geq 0, \quad i = 1, 2, \dots, N$$



➤ SVM的对偶问题

- 拉格朗日函数 $L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu})$

$$L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^N \mu_i \xi_i$$

- 原问题(Primal)的最优解: $p^* = \min_{\mathbf{w}, b} \theta_P(\mathbf{w}, b) = \min_{\mathbf{w}, b} \max_{\alpha_i, \mu_i \geq 0} L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\mu})$
- 对偶问题(Dual)的最优解: $d^* = \max_{\alpha_i, \mu_i \geq 0} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\mu})$
- 二者关系 $d^* \leq p^*$
- 满足KKT条件时, $d^* = p^*$

注意: 和原问题相比, 对偶问题交换了max和min的顺序。

- 对 \mathbf{w}, b, ξ_i 求导, 令一阶导数为0, 求 $\min_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha})$

$$L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^N \mu_i \xi_i$$

$$\frac{\partial L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial b} = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$$

$$\frac{\partial L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial \xi_i} = 0 \Rightarrow C = \alpha_i + \mu_i$$

- 将上述结论代入拉格朗日函数 $L(\mathbf{w}, b, \alpha, \xi, \mu)$

$$L(\mathbf{w}, b, \alpha, \xi, \mu) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^N \mu_i \xi_i$$

$$C = \alpha_i + \mu_i$$

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

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$= \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i) + \sum_{i=1}^N \alpha_i \xi_i$$

$$= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_{i=1}^N \alpha_i y_i b + \sum_{i=1}^N \alpha_i$$

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

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

- 从而得到对偶变量 α 的优化问题:

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



w, b 的计算

得到最佳的 α 后, 可计算 $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$

用任意一个支持向量即可求得 b : $b = y_i - \mathbf{w}^T \mathbf{x}_i$

为了得到更稳定的解, 对所有支持向量求平均

$$b = \frac{1}{N_S} \sum_{m \in S} \left(y_m - \sum_{m' \in S} \alpha_{m'} y_{m'} \langle \mathbf{x}_m, \mathbf{x}_{m'} \rangle \right)$$

得到 α, \mathbf{w}, b 后, 可计算SVM模型:

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b \\ &= \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)^T \mathbf{x} + b \\ &= \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \\ &= \sum_{i=1}^N \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b \end{aligned} \quad \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

点 \mathbf{x} 的标签可根据 $f(\mathbf{x})$ 的符号得到: $\hat{y} = \text{sign}(f(\mathbf{x}))$

➤ α 的稀疏性

SVM目标函数对应的KKT条件中，每个训练样本点：

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

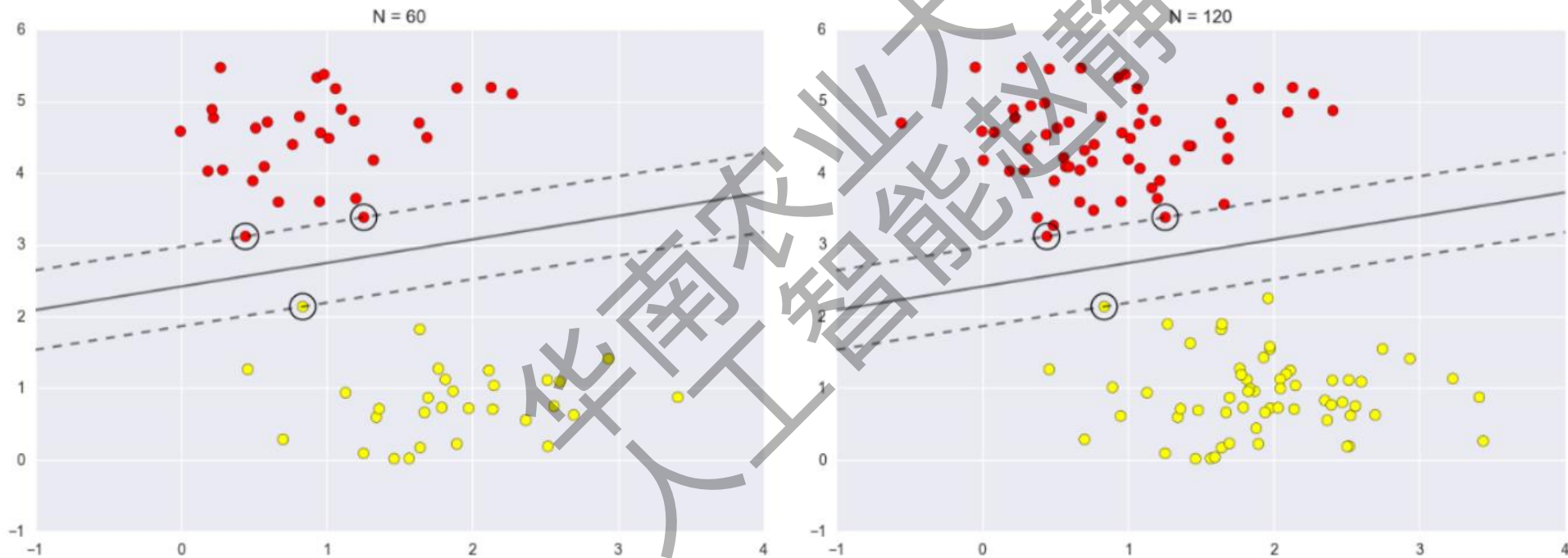
因此， $\alpha_i = 0$ 或 $1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) = 0$

- 当 $\alpha_i = 0$ 时，该点在决策函数中不起作用（非支持向量）

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{i=1}^N \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b$$

- 当 $\alpha_i \neq 0$ 时， $1 - \xi_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) = 0$ ，这些样本点被称为支持向量。

支持向量：真正发挥作用的数据点， α 值不为0的点



例：数据为3个样本，其中正例 $X_1(3,3)$ ， $X_2(4,3)$ ，负例 $X_3(1,1)$

解：
$$\min_{\alpha} \left(\frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^m \alpha_i \right)$$

s. t.
$$\sum_{i=1}^m \alpha_i y_i = 0,$$

$$\alpha_i \geq 0$$

$$\frac{1}{2} (18\alpha_1^2 + 25\alpha_2^2 + 2\alpha_3^2 + 42\alpha_1\alpha_2 - 12\alpha_1\alpha_3 - 14\alpha_2\alpha_3) - \alpha_1 - \alpha_2 - \alpha_3$$

由于： $\alpha_1 + \alpha_2 = \alpha_3$ 化简可得： $4\alpha_1^2 + \frac{13}{2}\alpha_2^2 + 10\alpha_1\alpha_2 - 2\alpha_1 - 2\alpha_2$

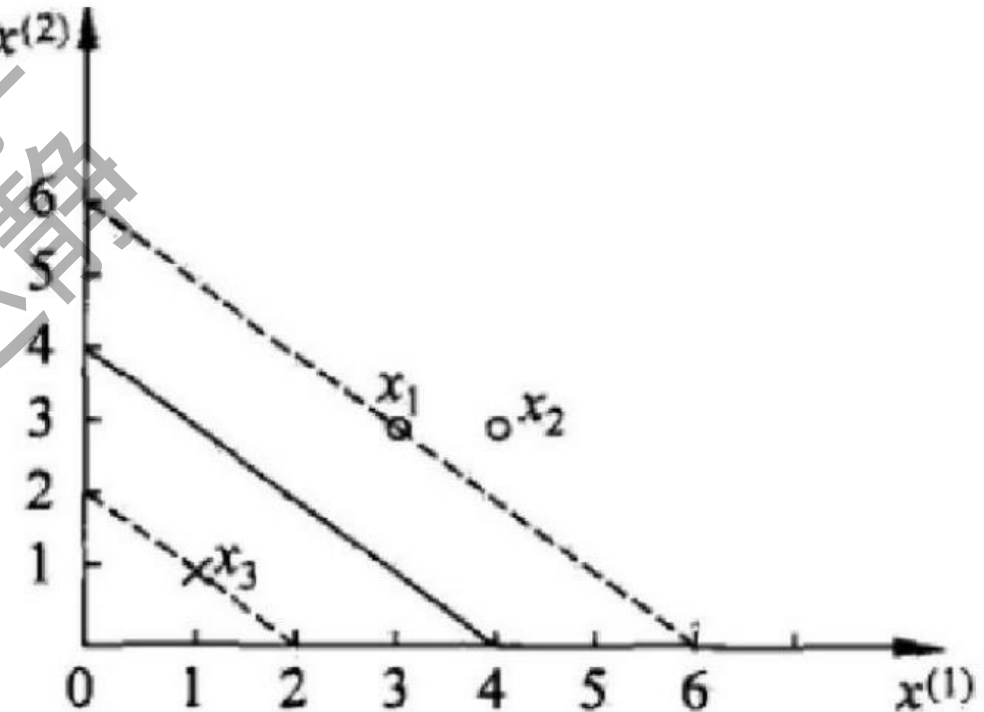
分别对 α_1 和 α_2 求偏导，偏导等于0可得， $\alpha_1 = 1.5 \alpha_2 = -1$

所以解应在边界上 $\alpha_1 = 0 \alpha_2 = -2/13$

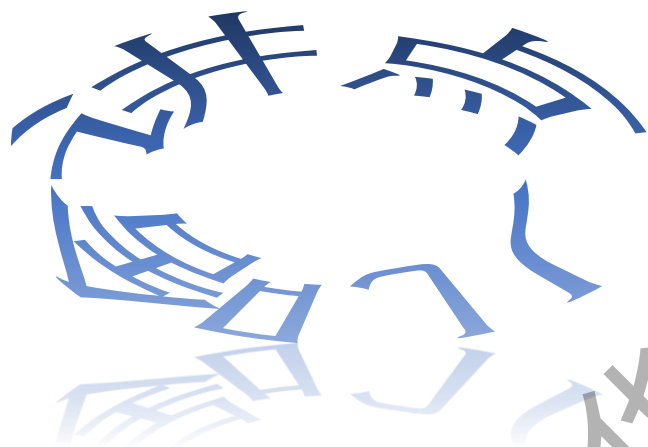
$\alpha_1 = 0.25 \alpha_2 = 0, \alpha_3 = 0.25$ ✓

$w^* = \sum_{i=1}^m \alpha_i^* y_i x_i = 1/4 * 1 * (3,3) + 1/4 * (-1) * (1,1) = (1/2, 1/2)$

$b^* = y - w^{*T} x = 1 - (1/2, 1/2)^T (3,3) = -2$

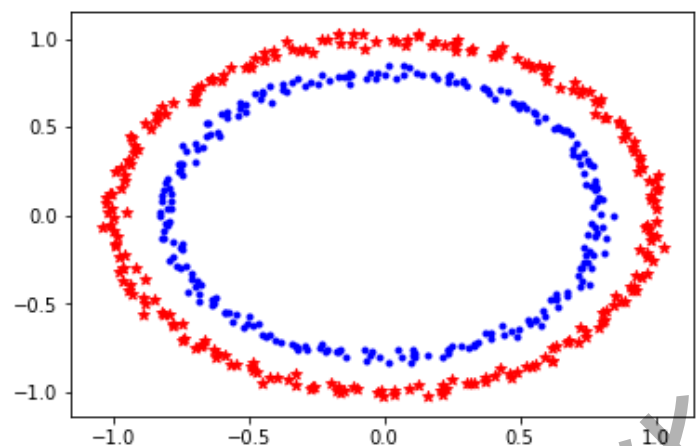
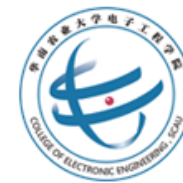


$1/2 x_1 + 1/2 x_2 - 2 = 0$

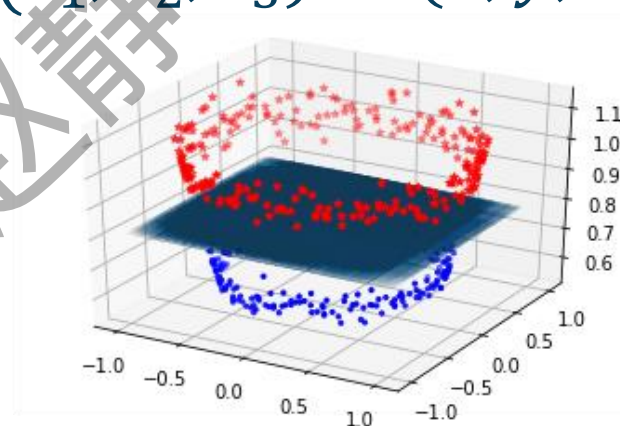


- ◆ 线性SVM
- ◆ 带松弛因子的SVM
- ◆ 合页损失函数
- ◆ SVM的对偶问题
- ◆ **核化SVM模型**
- ◆ SVM回归 (SVR)

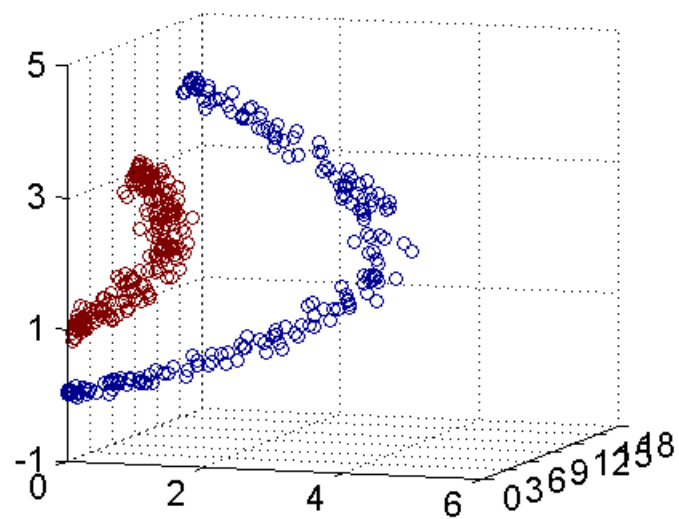
5. 核方法



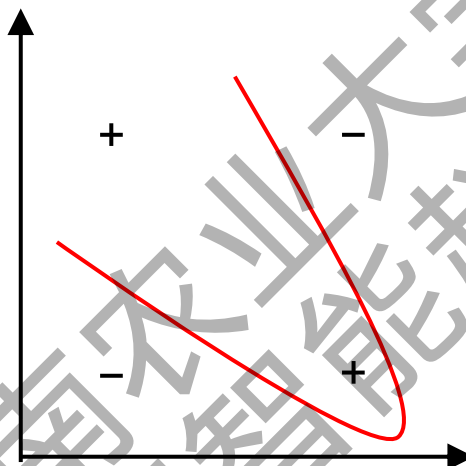
$$(z_1, z_2, z_3) = (x, y, x^2 + y^2)$$



将原始空间映射到一个更高维特征空间，使得在这个特征空间数据线性可分。

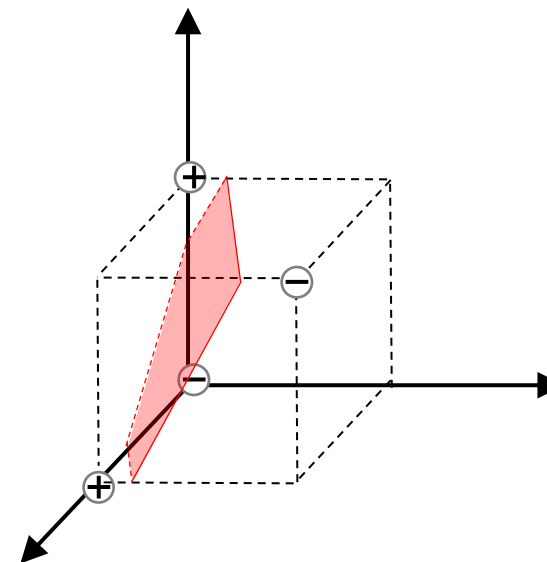


高维下线性可分



XOR数据集

$$\mathbf{x} \mapsto \phi(\mathbf{x})$$





构造一个5维函数

$$\varphi(\mathbf{x}): \mathbf{x} = \begin{bmatrix} a \\ b \end{bmatrix} \longrightarrow \varphi(\mathbf{x}) = \begin{bmatrix} a \\ b \\ ab \\ ab^2 \\ ab^3 \end{bmatrix}$$

$$\varphi(\mathbf{x}_1) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\varphi(\mathbf{x}_2) = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\varphi(\mathbf{x}_3) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\varphi(\mathbf{x}_4) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

设:

$$\omega = \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \\ 6 \end{bmatrix}$$

$$b=1$$

$$\omega^T \varphi(\mathbf{x}_1) + b = 1 \geq 0$$

$$\omega^T \varphi(\mathbf{x}_2) + b = 3 \geq 0$$

$$\omega^T \varphi(\mathbf{x}_3) + b = -1 < 0$$

$$\omega^T \varphi(\mathbf{x}_4) + b = -1 < 0$$

➤ 目标函数

- 令 $\phi(\mathbf{x})$ 表示将 \mathbf{x} 映射后的特征向量，则在特征空间中划分超平面对应的模型可表示为

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

- 根据之前SVM的推导，得到特征映射后的SVM目标函数

$$\min C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } 1 - \xi_i \leq y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b), \quad i = 1, 2, \dots, N$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, N$$

$$x \mapsto \phi(x)$$

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$s. t. \quad \textcircled{1} \quad \xi_i \geq 0$$

$$\textcircled{2} \quad y_i(w^T x_i + b) \geq 1 - \xi_i$$



$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$s. t. \quad \textcircled{1} \quad \xi_i \geq 0$$

$$\textcircled{2} \quad y_i(w^T \phi(x_i) + b) \geq 1 - \xi_i$$

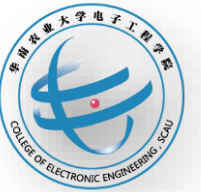
➤ 核方法——对偶

- 相应的对偶问题为：

$$\begin{aligned} \max & \left(\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \right) \\ \text{s. t. } & \sum_{i=1}^N \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, N \end{aligned}$$

- 求得对偶问题的解 α 后，可计算 \mathbf{w}, b ，从而得到分类判别函数：

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^N \alpha_i y_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle + b \end{aligned}$$



➤ 核技巧 (Kernel Trick)

- 判别函数为: $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=1}^N \alpha_i y_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle + b$
- **核函数**: 高维空间中的点积可写成**核**(*kernel*)的形式

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

- SVM核化目标函数为

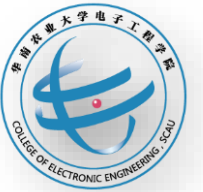
$$W(\boldsymbol{\alpha}) = \sum_{i=0}^N \alpha_i - \frac{1}{2} \sum_{i=0}^N \sum_{j=0}^N \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

- 判别函数为: $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=1}^N \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i) + b$

Kernel Trick

$$\begin{aligned}x &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\K(x, z) &= \phi(x) \cdot \phi(z) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix} \\&= x_1^2z_1^2 + 2x_1x_2z_1z_2 + x_2^2z_2^2 \\ \phi(x) &= \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} = (x_1z_1 + x_2z_2)^2 = \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \right)^2 \\&= (x \cdot z)^2\end{aligned}$$

Kernel Trick



$$K(x, z) = (x \cdot z)^2$$

$$= (x_1 z_1 + x_2 z_2 + \cdots + x_k z_k)^2$$

$$= \underline{x_1^2 z_1^2} + \underline{x_2^2 z_2^2} + \cdots + \underline{x_k^2 z_k^2}$$

$$+ 2\underline{x_1 x_2 z_1 z_2} + 2\underline{x_1 x_3 z_1 z_3} + \cdots$$

$$+ 2\underline{x_2 x_3 z_2 z_3} + 2\underline{x_2 x_4 z_2 z_4} + \cdots$$

$$= \phi(x) \cdot \phi(z)$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix} \quad z = \begin{bmatrix} z_1 \\ \vdots \\ z_k \end{bmatrix}$$

$$\phi(x) = \begin{bmatrix} x_1^2 \\ \vdots \\ x_k^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \vdots \\ \sqrt{2}x_2x_3 \\ \vdots \end{bmatrix}$$

➤ 构造核函数

核函数 K 和映射 ϕ 是一一对应关系

核函数的形式不能随意取

满足一定的条件

两个 ϕ 内积的形式

- 令 \mathcal{X} 为输入空间, $k(\cdot, \cdot)$ 是定义在 $\mathcal{X} \times \mathcal{X}$ 上的对称函数, 则 k 是核函数的充要条件是对任意数据 $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$, 核矩阵 \mathbf{K} 总是半正定的:

$$\mathbf{K} = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \kappa(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \kappa(\mathbf{x}_2, \mathbf{x}_1) & \kappa(\mathbf{x}_2, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) & \kappa(\mathbf{x}_N, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

✓ 多项式核

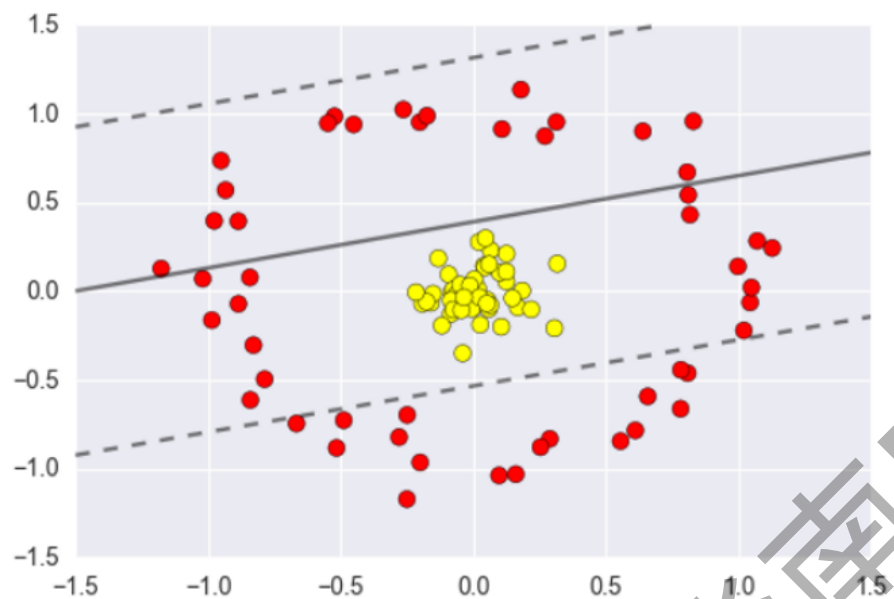
- 多项式核: $k(\mathbf{x}, \mathbf{x}') = (\gamma \mathbf{x}^T \mathbf{x}' + r)^M$
- 当 $M = 1$, $\gamma = 1$, $r = 0$ 时, 为线性核,

$$\phi(\mathbf{x}) = \mathbf{x}$$

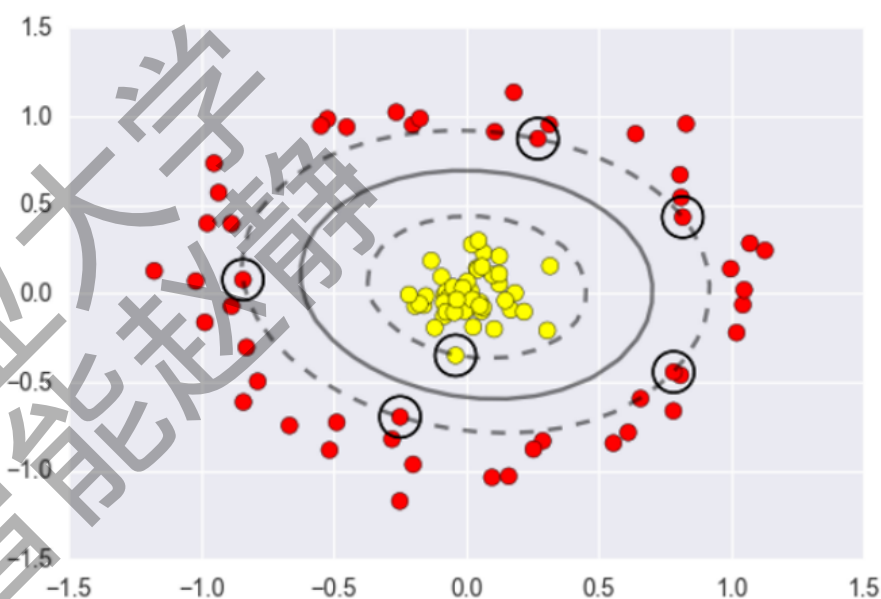
$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

✓ 高斯核函数 (径向基核函数 RBF)

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



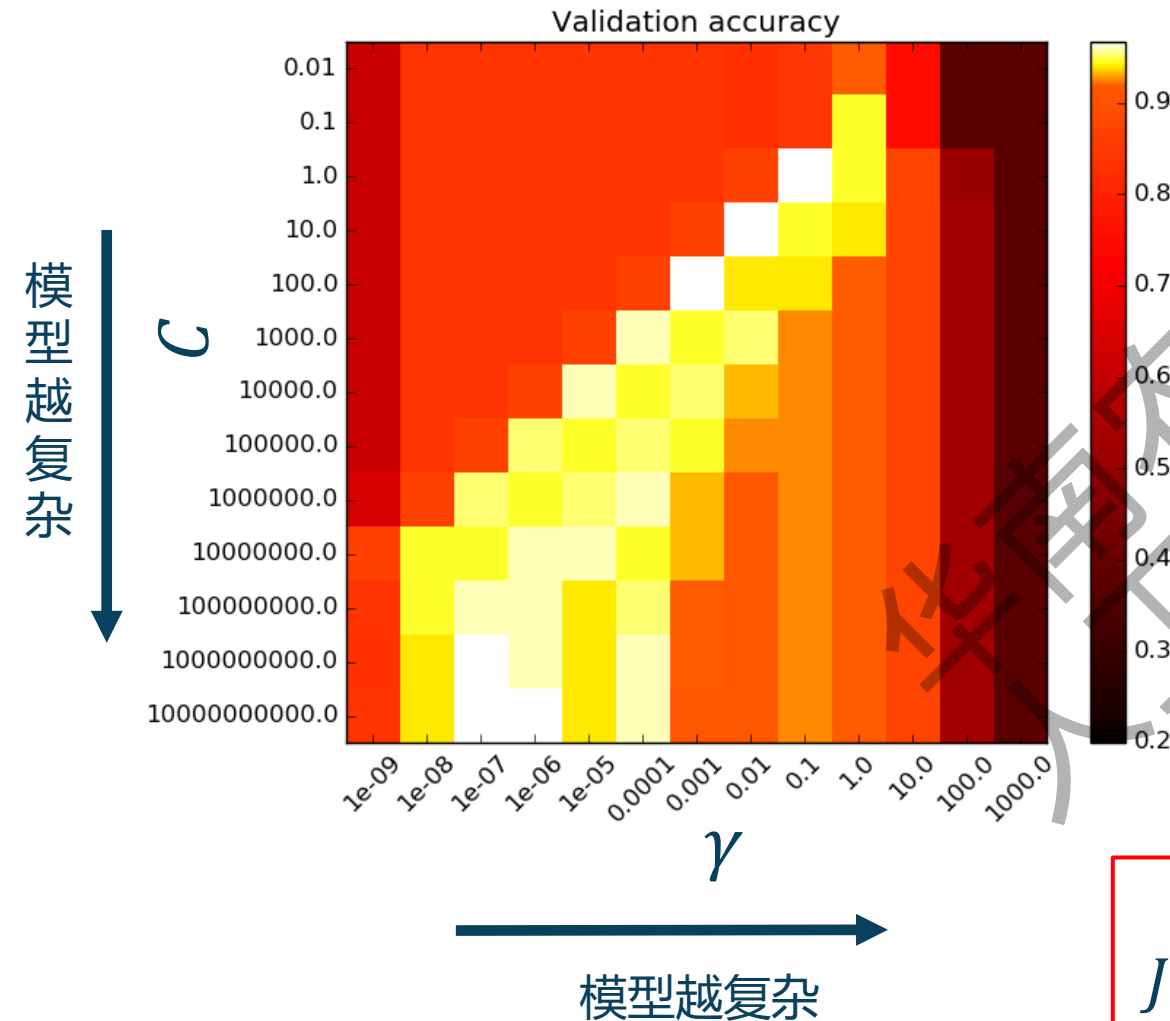
线性核函数



高斯核函数

例：RBF核

在鸢尾花分类任务上（只取了前2维特征），不同参数值的RBF核SVM分类器的交叉验证精度



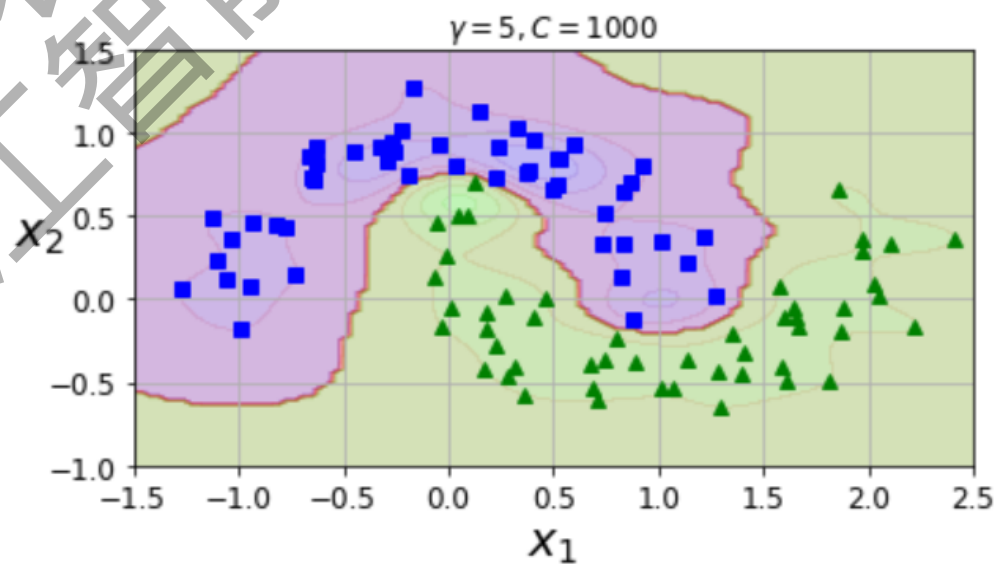
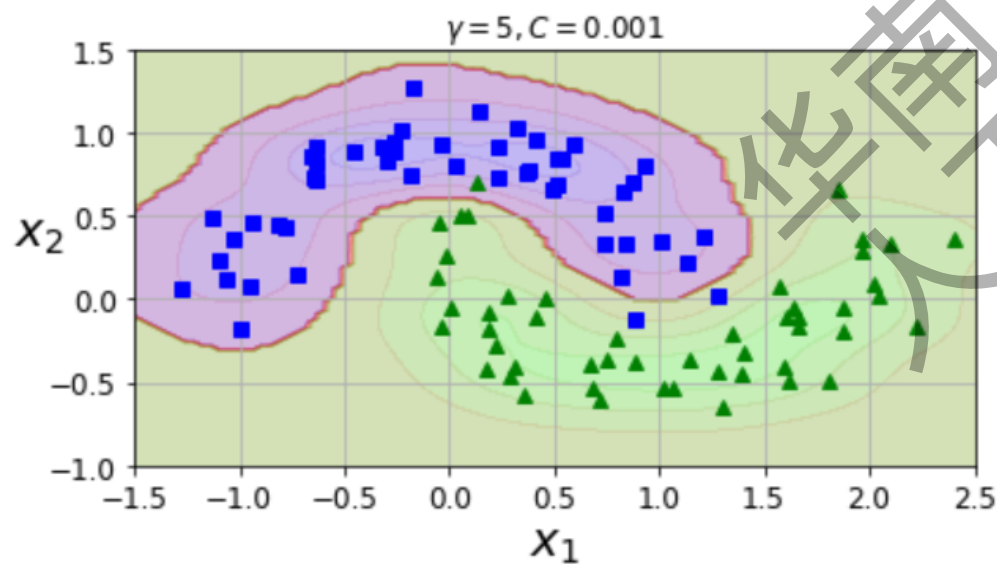
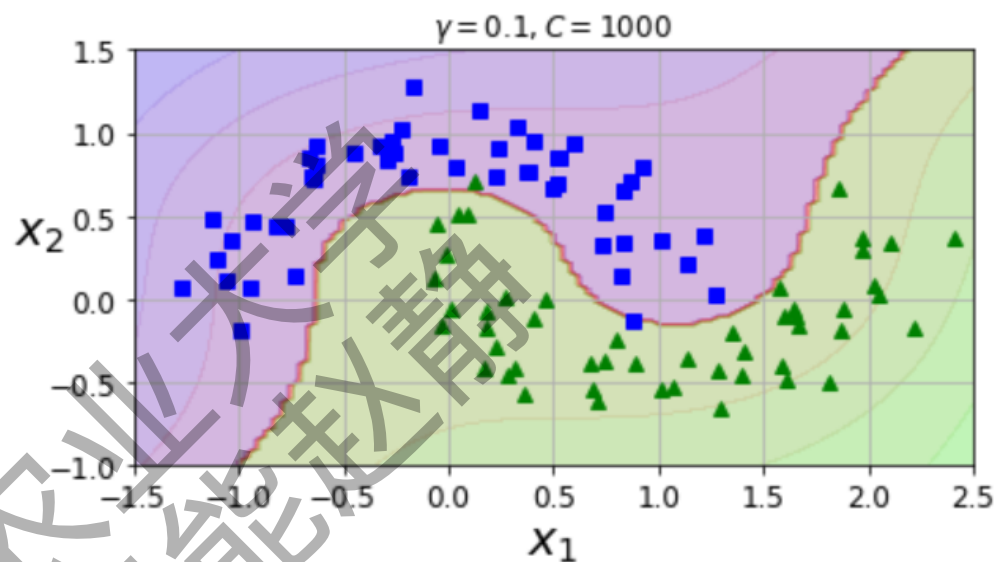
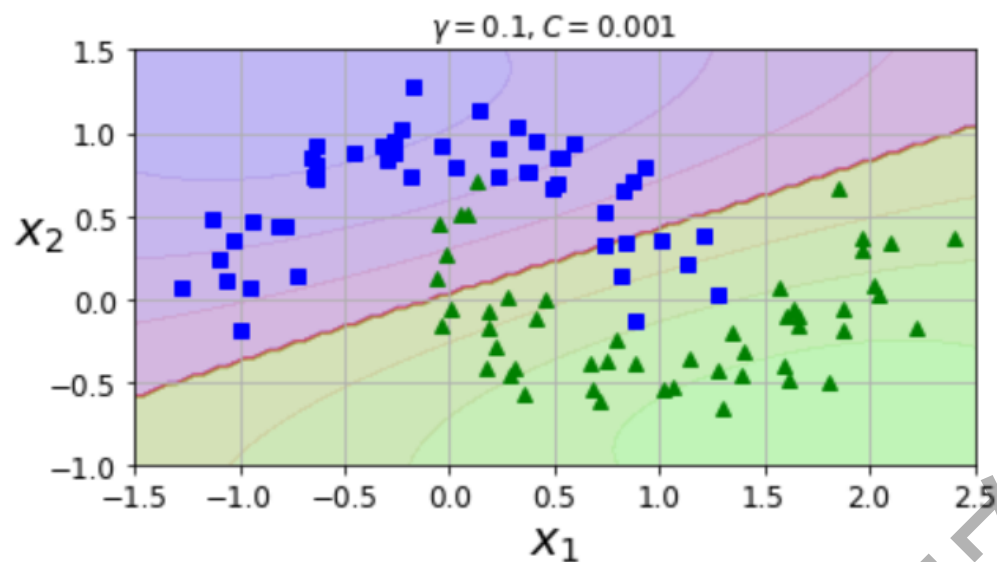
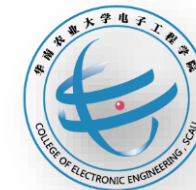
C : 训练样本有多重要
越大, 越看重训练样本, 模型越复杂

γ : RBF核宽度的倒数,
表示每个训练样本的影响范围

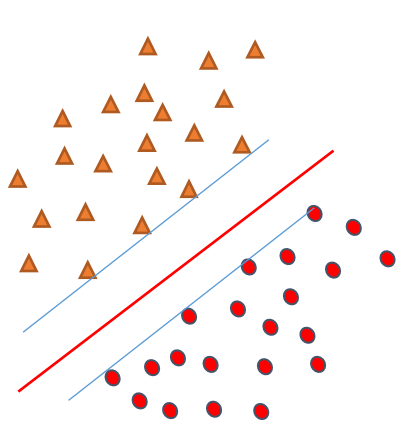
越大, RBF核宽度越小,
每个训练样本/支持向量的范围越小,
模型越复杂

$$J(\mathbf{w}, b; C) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N L_{Hinge}(y_i, f(\mathbf{x}_i; \mathbf{w}, b))$$

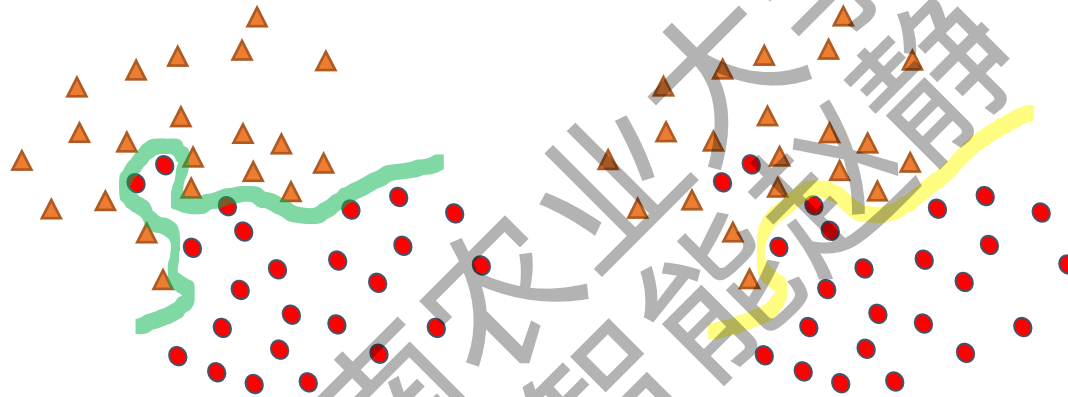
例：RBF核



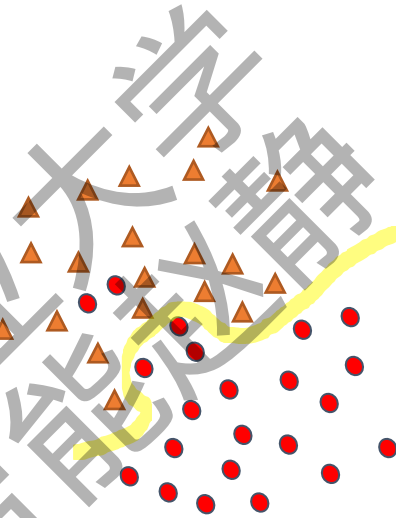
硬间隔、软间隔和非线性 SVM



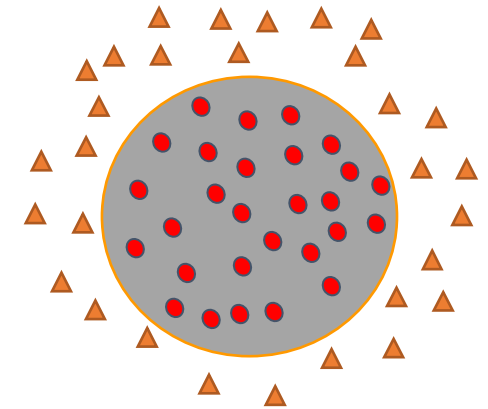
硬间隔



软间隔



软间隔



线性不可分

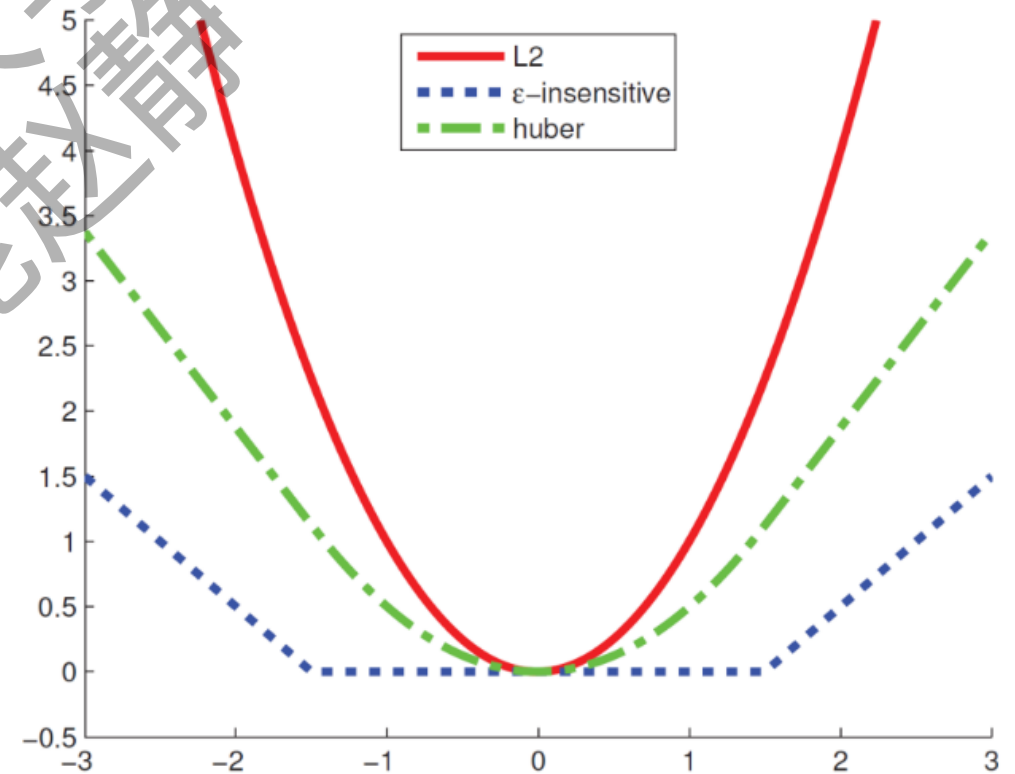


- ◆ 线性SVM
- ◆ 带松弛因子的SVM
- ◆ 合页损失函数
- ◆ SVM的对偶问题
- ◆ 核化SVM模型
- ◆ SVM回归 (SVR)

➤ ϵ 不敏感损失函数 (ϵ insensitive loss)

$$L_{\epsilon}(y, \hat{y}) = \begin{cases} 0 & |y - \hat{y}| \leq \epsilon \\ |y - \hat{y}| - \epsilon & \text{otherwise} \end{cases}$$

不惩罚小的损失





➤ 支持向量回归 (Support Vector Regression, SVR)

- 假设回归函数为线性模型: $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, SVR的目标函数为

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

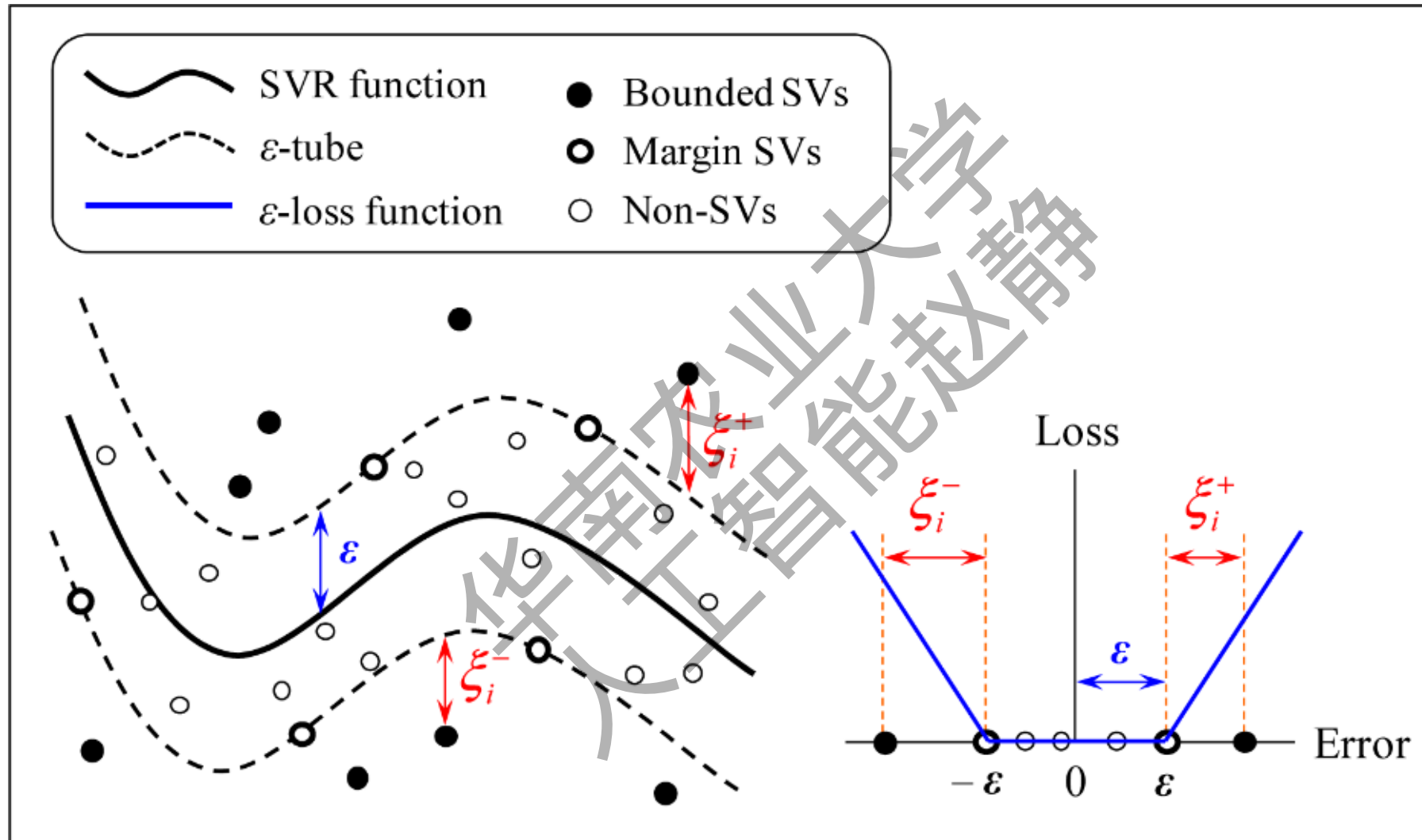
$$\text{s. t. } |y_i - (\mathbf{w}^T \mathbf{x}_i + b)| \leq \epsilon, \quad i = 1, 2, \dots, N$$

- 加入松弛变量 $\xi_i \geq 0$, 用于表示每个点在 ϵ 管道外的程度, SVR模型的目标函数变为

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N (\xi_i^V + \xi_i^A)$$

$$\text{s. t. } -\epsilon - \xi_i^V \leq y_i - (\mathbf{w}^T \mathbf{x}_i + b) \leq \epsilon + \xi_i^A, \quad i = 1, 2, \dots, N$$

$$\xi_i^A \geq 0, \quad \xi_i^V \geq 0, \quad i = 1, 2, \dots, N$$



➤ 拉格朗日函数

- 与SVM类似，SVR的拉格朗日函数为：

$$L(\mathbf{w}, b, \boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge, \boldsymbol{\xi}^V, \boldsymbol{\xi}^\wedge, \boldsymbol{\mu}^V, \boldsymbol{\mu}^\wedge) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N (\xi_i^V + \xi_i^\wedge)$$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } -\epsilon - \xi_i^V \leq y_i - (\mathbf{w}^T \mathbf{x}_i + b) \leq \epsilon + \xi_i^\wedge$$

$$\xi_i^\wedge \geq 0, \quad \xi_i^V \geq 0$$

$$+ \sum_{i=1}^N \alpha_i^V (-\epsilon - \xi_i^V - y_i + \mathbf{w}^T \mathbf{x}_i + b)$$

$$+ \sum_{i=1}^N \alpha_i^\wedge (y_i - \mathbf{w}^T \mathbf{x}_i - b - \epsilon - \xi_i^\wedge)$$

$$- \sum_{i=1}^N \mu_i^V \xi_i^V - \sum_{i=1}^m \mu_i^\wedge \xi_i^\wedge$$

➤ SVR的对偶问题

- 拉格朗日函数 $L(\mathbf{w}, b, \boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge, \boldsymbol{\xi}^V, \boldsymbol{\xi}^\wedge, \boldsymbol{\mu}^V, \boldsymbol{\mu}^\wedge)$ 对 $\mathbf{w}, b, \xi_i^\wedge, \xi_i^V$ 的一阶导数:

$$\frac{\partial L(\mathbf{w}, b, \boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge, \boldsymbol{\xi}^V, \boldsymbol{\xi}^\wedge, \boldsymbol{\mu}^V, \boldsymbol{\mu}^\wedge)}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^N (\alpha_i^\wedge - \alpha_i^V) \mathbf{x}_i$$

$$\frac{\partial L(\mathbf{w}, b, \boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge, \boldsymbol{\xi}^V, \boldsymbol{\xi}^\wedge, \boldsymbol{\mu}^V, \boldsymbol{\mu}^\wedge)}{\partial b} = 0 \Rightarrow \sum_{i=1}^N (\alpha_i^\wedge - \alpha_i^V) = 0$$

$$\frac{\partial L(\mathbf{w}, b, \boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge, \boldsymbol{\xi}^V, \boldsymbol{\xi}^\wedge, \boldsymbol{\mu}^V, \boldsymbol{\mu}^\wedge)}{\partial \xi_i^\wedge} = 0 \Rightarrow C - \alpha_i^\wedge + \mu_i^\wedge$$

$$\frac{\partial L(\mathbf{w}, b, \boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge, \boldsymbol{\xi}^V, \boldsymbol{\xi}^\wedge, \boldsymbol{\mu}^V, \boldsymbol{\mu}^\wedge)}{\partial \xi_i^V} = 0 \Rightarrow C - \alpha_i^V + \mu_i^V$$

- 将上述结论代入拉格朗日函数 $L(\mathbf{w}, b, \boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge, \boldsymbol{\xi}^V, \boldsymbol{\xi}^\wedge, \boldsymbol{\mu}^V, \boldsymbol{\mu}^\wedge)$ ，得到对偶问题：

$$\begin{aligned} & \max_{\boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge} \left(- \sum_{i=1}^N ((\epsilon - y_i)\alpha_i^\wedge + (\epsilon + y_i)\alpha_i^V) - \sum_{i=1}^N \frac{1}{2} (\alpha_i^\wedge - \alpha_i^V)(\alpha_j^\wedge - \alpha_j^V) \mathbf{x}_i^T \mathbf{x}_j \right) \\ \text{s. t. } & \sum_{i=1}^N (\alpha_i^\wedge - \alpha_i^V) = 0 \\ & 0 \leq \alpha_i^\wedge \leq C, \quad i = 1, 2, \dots, N \\ & 0 \leq \alpha_i^V \leq C, \quad i = 1, 2, \dots, N \end{aligned}$$

- 去掉负号，将上述目标函数中的max换成min，得到等价问题：

$$\min_{\boldsymbol{\alpha}^V, \boldsymbol{\alpha}^\wedge} \left(\sum_{i=1}^N ((\epsilon - y_i)\alpha_i^\wedge + (\epsilon + y_i)\alpha_i^V) + \sum_{i=1}^N \frac{1}{2} (\alpha_i^\wedge - \alpha_i^V)(\alpha_j^\wedge - \alpha_j^V) \mathbf{x}_i^T \mathbf{x}_j \right)$$

➤ SVR模型

- 求得 α 的最优值后，可计算 \mathbf{w}, b ，从而得到SVR模型

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b = \sum_{i=1}^N (\alpha_i^{\wedge} - \alpha_i^{\vee}) \mathbf{x}_i^T \mathbf{x} + b \\ &= \sum_{i=1}^N (\alpha_i^{\wedge} - \alpha_i^{\vee}) \langle \mathbf{x}_i, \mathbf{x} \rangle + b \end{aligned}$$

- \mathbf{x} 与训练数据的点积 $\langle \mathbf{x}_i, \mathbf{x} \rangle$ 换成核函数 $k(\mathbf{x}_i, \mathbf{x})$ ，得到核化SVR模型

$$f(\mathbf{x}) = \sum_{i=1}^N (\alpha_i^{\wedge} - \alpha_i^{\vee}) k(\mathbf{x}_i, \mathbf{x}) + b$$

➤ SVR的支持向量

- 根据KKT条件:

$$\alpha_i^V (\epsilon + \xi_i^V + y_i - (\mathbf{w}^T \mathbf{x}_i + b)) = 0$$

$$\alpha_i^A (\epsilon + \xi_i^A - y_i + (\mathbf{w}^T \mathbf{x}_i + b)) = 0$$

支持向量:

仅当样本 (\mathbf{x}_i, y_i) 落在 ϵ 间隔中, α_i^V 、 α_i^A 才取非0值。

