

1.

1) 最小错误率决策: 应当寻找最大的后验概率, 来使错误率最小.

$$P(\text{error}|\mathbf{x}) = 1 - P(\mathbf{w}_i|\mathbf{x}) \quad \text{if you decide } \mathbf{w}_i$$

• Minimum error decision: Maximum a posteriori (MAP)

$$\text{Decide } \omega_i \text{ if } P(\omega_i|\mathbf{x}) > P(\omega_j|\mathbf{x}) \quad \text{for all } j \neq i$$

最小风险决策: 通过计算不同类别的最小风险决策, 选择最小的  $R(\alpha_i|\mathbf{x})$  所在的类别.

贝叶斯最小风险决策:

Condition risk

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j) P(\omega_j|\mathbf{x})$$

Minimum risk decision (Bayes decision)

$$\arg \min_i R(\alpha_i|\mathbf{x})$$

2)

• Formulation (Problem 13, Chapter 2)

– C+1 classes

$$\lambda(\alpha_i|\omega_j) = \begin{cases} 0, & i = j \\ \lambda_s, & i \neq j \\ \lambda_r, & \text{reject} \end{cases}$$

替典错误

$$\lambda_r < \lambda_s$$

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j) P(\omega_j|\mathbf{x})$$

$$\Rightarrow R_i(\mathbf{x}) = \begin{cases} \lambda_s [1 - P(\omega_i|\mathbf{x})], & i = 1, \dots, c \\ \lambda_r, & \text{reject} \end{cases}$$

这里不是公式, 只是一个事实陈述

$$\arg \min_i R_i(\mathbf{x}) = \begin{cases} \arg \max_i P(\omega_i|\mathbf{x}), & \text{if } \max_i P(\omega_i|\mathbf{x}) > 1 - \lambda_r / \lambda_s \\ \text{reject}, & \text{otherwise} \end{cases}$$

2.

1)

Case 2:  $\Sigma_i = \Sigma$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mu_i)^t \Sigma^{-1}(\mathbf{x} - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma| + \ln P(\omega_i)$$

$$\Rightarrow g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mu_i)^t \Sigma^{-1}(\mathbf{x} - \mu_i) + \ln P(\omega_i)$$

– 展开二次式  $(\mathbf{x} - \mu_i)^t \Sigma^{-1}(\mathbf{x} - \mu_i)$

线性判别函数!  $g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0}$

$$\mathbf{w}_i = \Sigma^{-1} \mu_i \quad w_{i0} = -\frac{1}{2} \mu_i^t \Sigma^{-1} \mu_i + \ln P(\omega_i)$$

– 二类决策面  $g_i(\mathbf{x}) = g_j(\mathbf{x})$

$$\Rightarrow \mathbf{w}^t (\mathbf{x} - \mathbf{x}_0) = 0 \quad \mathbf{w} = \Sigma^{-1}(\mu_i - \mu_j)$$

替为任意值, 不一定准确

$$\mathbf{x}_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\ln [P(\omega_i)/P(\omega_j)]}{(\mu_i - \mu_j)^t \Sigma^{-1}(\mu_i - \mu_j)} (\mu_i - \mu_j)$$

二类决策面带入具体的数值进行计算得

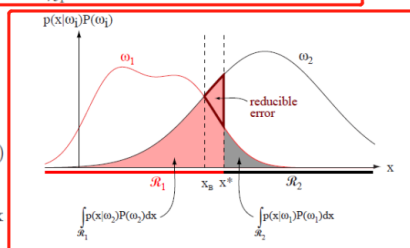
$$X_2 = 2 \times X_1$$

## 2类情况

$$\begin{aligned} P(\text{error}) &= P(x \in R_2, \omega_1) + P(x \in R_1, \omega_2) \\ &= P(x \in R_2 | \omega_1) P(\omega_1) + P(x \in R_1 | \omega_2) P(\omega_2) \\ &= \int_{R_2} p(x | \omega_1) P(\omega_1) dx + \int_{R_1} p(x | \omega_2) P(\omega_2) dx. \end{aligned}$$

## 一般情况

$$\begin{aligned} P(\text{correct}) &= \sum_{i=1}^c P(x \in R_i, \omega_i) \\ &= \sum_{i=1}^c P(x \in R_i | \omega_i) P(\omega_i) \\ &= \sum_{i=1}^c \int_{R_i} p(x | \omega_i) P(\omega_i) dx \end{aligned}$$



决策面为 $x_0$ 时为最小错误率分类

$R_2$  与  $R_1$  的范围由判别面方程可求，概率密度和先验可以带入具体数值。

2)

利用风险公式 后验概率\*损失 列出来，然后列个等式就好了  
可以取对数进行约简计算

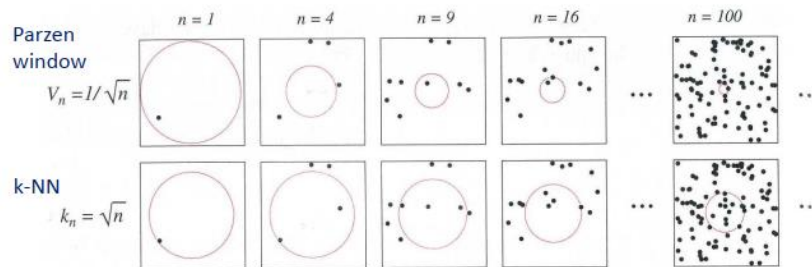
Handwritten calculations for a linear discriminant analysis problem:

$$\begin{aligned} &P(\omega_1 | x) = P(\omega_2 | x) \\ &\ln 2 + \ln P(\omega_1 | x) = \ln P(\omega_2 | x) \\ &\ln 2 + w_1^T x + w_0 = w_2^T x + w_0 \\ &\text{其中 } w_2^T = \begin{bmatrix} -\frac{4}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{4}{3} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -\frac{4}{3} \\ -\frac{1}{3} \end{bmatrix} \\ &w_0 = -\frac{1}{2} w_1^T x + \ln 2 + \ln P(\omega_1) \\ &= -\frac{1}{2} \begin{bmatrix} -\frac{4}{3} & -\frac{1}{3} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \ln 2 + \ln \frac{1}{2} \\ &= \frac{2}{3} + \ln \frac{1}{2} \\ &w_1^T = \begin{bmatrix} -\frac{4}{3} & -\frac{1}{3} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -\frac{4}{3} \\ -\frac{1}{3} \end{bmatrix} \\ &w_0 = \frac{1}{2} \begin{bmatrix} -\frac{4}{3} & -\frac{1}{3} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \ln 2 + \ln \frac{1}{2} \\ &= -\frac{2}{3} + \ln \frac{1}{2} \\ &\text{所以 } \ln 2 + \frac{4}{3} x_1 + \frac{1}{3} x_2 + (-\frac{2}{3}) + \ln \frac{1}{2} = \frac{4}{3} x_1 - \frac{2}{3} x_2 + (-\frac{2}{3}) + \ln \frac{1}{2} \\ &\frac{4}{3} x_1 - \frac{1}{3} x_2 = \ln 2 \end{aligned}$$

3.  
1)

### • 非参数概率密度估计

- Parzen window: 固定局部区域体积  $V$ ,  $k$  变化
- k-nearest neighbor: 固定局部样本数  $k$ ,  $V$  变化



2)

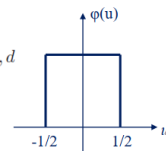
### Parzen Window

#### • 窗函数: hypercube

$$\varphi(u) = \begin{cases} 1 & |u_j| \leq 1/2 \\ 0 & \text{otherwise.} \end{cases} \quad j = 1, \dots, d$$

- 满足条件

$$\varphi(u) \geq 0 \quad \int \varphi(u) du = 1$$



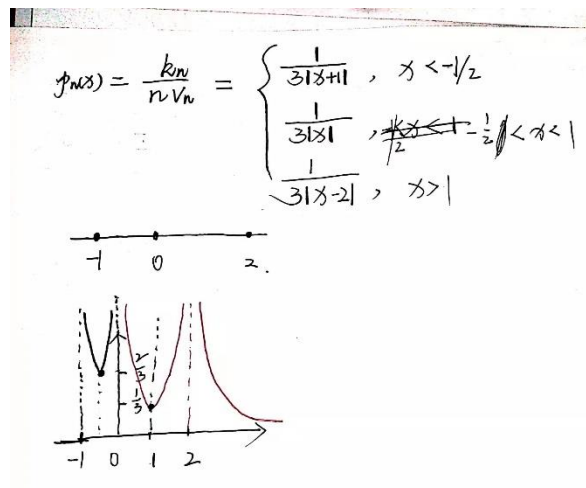
- 以  $x$  为中心、体积为  $V_n = h_n^d$  的局部区域内样本数

$$k_n = \sum_{i=1}^n \varphi\left(\frac{x - x_i}{h_n}\right)$$

$x_i$  训练样本  $x$  中心  
 $h_n$  窗口宽度

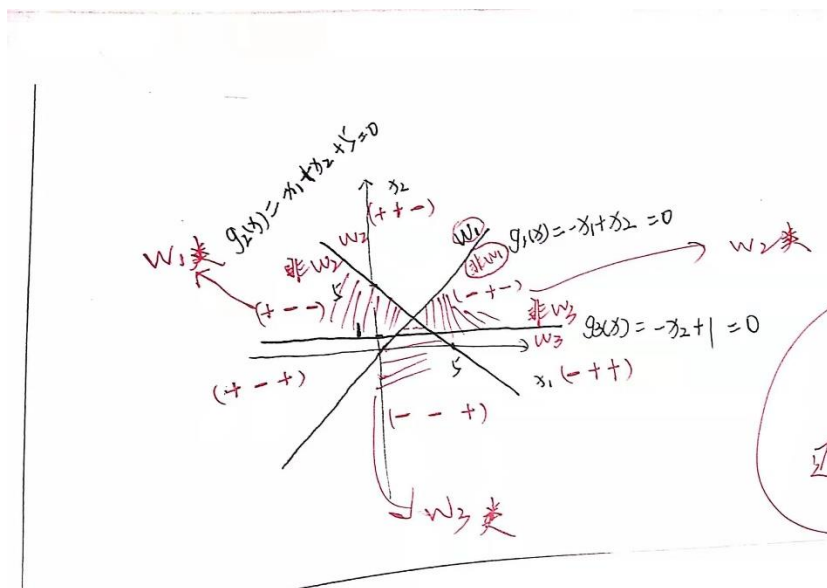
- 概率密度估计  $k_n/nV_n$  除以  $n$  再除以体积

$$p_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \varphi\left(\frac{x - x_i}{h_n}\right)$$



- 无穷的时候趋近于 0, 忘记画

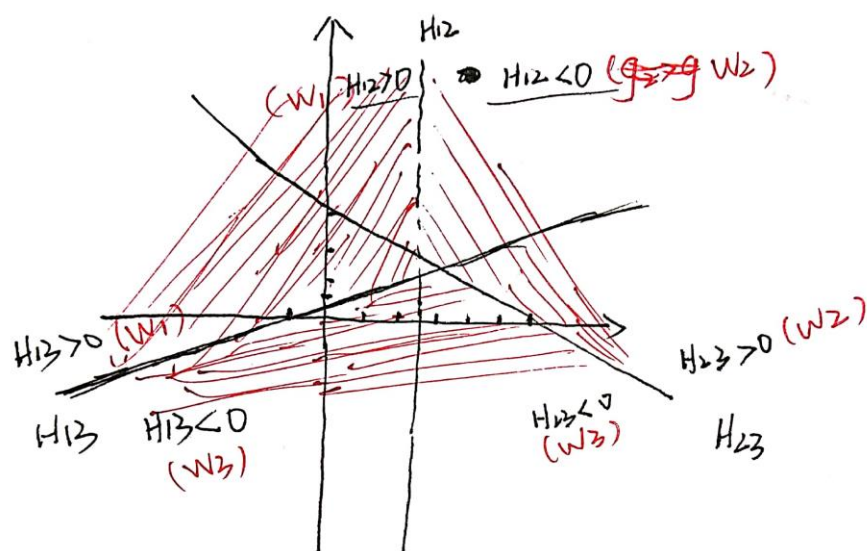
4.  
1)



$$H_{12} = g_1(x) - g_2(x) = -2x_1 + 5$$

$$H_{23} = g_2(x) - g_3(x) = x_1 + 2x_2 - 6$$

$$H_{13} = g_1(x) - g_3(x) = -x_1 + 2x_2 - 1$$



2)

- 感知准则函数—基本思想

- 考虑如下准则函数：

$$J_p(\mathbf{a}) = \sum_{\mathbf{y} \in Y} (-\mathbf{a}^T \mathbf{y}), \text{ 其中, } Y \text{ 为错分样本集合}$$

- 当  $\mathbf{y}$  被错分时,  $\mathbf{a}^T \mathbf{y} \leq 0$ , 则  $-\mathbf{a}^T \mathbf{y} \geq 0$ 。因此  $J_p(\mathbf{a})$  总是大于等于0。在可分情形下, 当且仅当  $Y$  为空集时  $J_p(\mathbf{a})$  将等于零, 这时将不存在错分样本。

- 因此, 目标是最小化  $J_p(\mathbf{a})$ :  $\min_{\mathbf{a}} J_p(\mathbf{a})$

### 可变增量批处理修正方法（另一种表述形式）

---

#### Batch Variable-Increment Perceptron

---

```
1  begin initialize:  $\mathbf{a}, \eta_0, k=0$ 
2  do  $k \leftarrow k+1 \pmod n$ 
3       $Y_k = \{\}$ ,  $j = 0$ 
4      do  $j \leftarrow j + 1$ 
5          if  $\mathbf{y}_j$  is misclassified, then append  $\mathbf{y}_j$  to  $Y_k$ 
6      until  $j = n$ 
7       $\mathbf{a} = \mathbf{a} + \eta_k \sum_{\mathbf{y} \in Y(k)} \mathbf{y}$  //发现所有错分, 然后再修正
8  until  $Y_k = \{\}$  //直到所有样本均正确分类
9  return  $\mathbf{a}$ 
10 end
```

---

5.

1)

1. 原理: 假设各类出现的先验概率均相等, 每个样本点以概率为1属于1个类, 计算数据点到类中心的欧氏距离的平方, 寻找与样本点最近的类中心点, 将其分给最近的类。

2)

5. (2)  $x_1 = (-6, 1)^T$ ,  $x_2 = (-6, -1)^T$ ,  $x_3 = (-4, 0)^T$ ,  $x_4 = (4, 0)^T$ ,  $x_5 = (5, 1)^T$   
 $x_6 = (6, -1)^T$

① 把每个样本看作一类, 计算点对之间的距离

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
$x_1$	0					
$x_2$	$\sqrt{5}$	0				
$x_3$	$\sqrt{5}$	$\sqrt{5}$	0			
$x_4$	$\sqrt{10}$	$\sqrt{10}$	$\sqrt{64}$	0		
$x_5$	$\sqrt{12}$	$\sqrt{12}$	$\sqrt{8}$	$\sqrt{2}$	0	
$x_6$	$\sqrt{13}$	$\sqrt{13}$	$\sqrt{10}$	$\sqrt{5}$	$\sqrt{5}$	0

②  $G_1 = \{x_4, x_5\}$



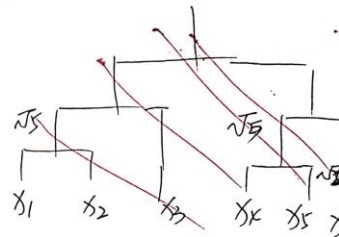
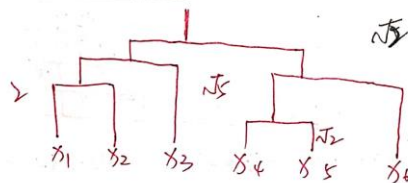
③ 根据最近距离原则,  $G_1 = \{x_4, x_5\}$ ,  $G_2 = \{x_3, x_6\}$

	$G_1$	$x_3$	$G_2$	$x_6$
$G_1$	0			
$x_3$	$\sqrt{5}$	0		
$G_2$	$\sqrt{10}$	$\sqrt{64}$	0	
$x_6$	$\sqrt{13}$	$\sqrt{10}$	$\sqrt{5}$	0

④  $G_3 = G_1 \cup \{x_3\} = \{x_4, x_5, x_3\}$ ,  $G_4 = G_2 \cup \{x_6\} = \{x_3, x_5, x_6\}$

	$G_3$	$G_4$
$G_3$	0	
$G_4$	$\sqrt{10}$	0

决策树.



6.

同 2018-2019 只是有关公式的地方换一下.

$$\frac{\Delta w_{ij}}{\Delta w_{ij}} = -\eta \frac{\partial E}{\partial w_{ij}} = -\eta \cdot \frac{\partial E}{\partial z_j^k} \cdot \dots \dots$$

$$\dots = -\eta \cdot \frac{\partial E}{\partial z_j^k} \cdot \dots \dots$$

$$= -\eta \cdot \frac{t_j^k}{z_j^k}$$

即:  $\Delta w_{ij} = \eta \sum_k \delta_j^k \cdot y_i^k$   $\delta_j^k = f_2'(net_j^k) \cdot \delta_j^k$   $\delta_j^k = \frac{t_j^k}{z_j^k}$

$$\Delta w_{ih} = \eta \sum_k \delta_h^k \cdot x_i^k$$

$$\delta_h^k = f_1'(net_h^k) \cdot \delta_h^k$$

$$\delta_h^k = \sum_j w_{hj} \delta_j^k$$

$$y_h^k = f_1(net_h^k)$$

$$z_h^k = f_2(net_h^k)$$

$$f_1'(net_h^k) = y_h^k(1 - y_h^k)$$

$$f_2'(net_h^k) = z_h^k(1 - z_h^k)$$

7.

- 1) 简述 PCA 的主要思想和求解过程
- 2) 比较 PCA、CCA、LDA、ICA 的区别与适用场景。
- 3) 解释 LDA 所基于的数据分布假设，并阐述其不足之处

解:

- 1) PCA 的主要思想: PCA 就是将高维的数据通过线性变换投影到低维空间上去, 得到一组无关的基。

求解过程:



3. 有  $N$  个样本  $x_1, \dots, x_N$ ，每个样本维数  $D$ ，希望将样本维数降低到  $K$ ，请给出

PCA 算法的计算过程

Step1: 计算数据点的均值

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

Step2: 计算样本的协方差矩阵

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$

Step3: 对  $D \times D$  的  $S$  矩阵进行特征值分解

Step4: 对特征值进行从大到小排序，并选取最大的  $K$  个，记为  $u_1, u_2, \dots, u_K$

Step5:  $U = [u_1, u_2, \dots, u_K]$  就是  $D \times K$  的投影矩阵

Step6: 利用  $Z = U^T X$  实现 PCA 降维

2)

**PCA:** 主成分分析。PCA 属于无监督（训练样本无标签）的降维方法，是一种正交投影，侧重选择样本点投影方差最大的方向，减少特征相关性。

适用场景：

- a) 原始数据特征多而且特征冗余。
- b) 需要对样本进行可视化的时候，三维以上的特征无法进行可视化

**CCA:** 典型相关分析。它选择的投影标准是降维到 1 维后，两组数据的相关系数最大。

适用场景：侧重于两组数据有相关关系的时候。

**LDA:** 线性判别分析，可以降维可以分类，是从更利于分类的角度的，有监督（训练样本有标签）。希望数据投影后类内方差最小，类间方差最大。

适用场景：侧重于分类。

**ICA:** 独立成分分析。ICA 信号需要是非高斯的，寻找的是最能使数据的相互独立的方向。

应用场景：盲信号分离。

3)

假设：

每一个类是单模态高斯分布  $\rightarrow$  多模态 LDA

每一个类的协方差矩阵都相同  $\rightarrow$  异方差 LDA

不足：

类分离问题。

降维维数不能超过  $C-1$

补充知识：

LDA 用于降维，和 PCA 有很多相同，也有很多不同的地方，因此值得好好的比较一下两者的降维异同点。

首先我们看看相同点：

- 1) 两者均可以对数据进行降维。
- 2) 两者在降维时均使用了矩阵特征分解的思想。
- 3) 两者都假设数据符合高斯分布。

我们接着看看不同点：



- 1) LDA 是有监督的降维方法，而 PCA 是无监督的降维方法
- 2) LDA 降维最多降到类别数  $k-1$  的维数，而 PCA 没有这个限制。
- 3) LDA 除了可以用于降维，还可以用于分类。
- 4) LDA 选择分类性能最好的投影方向，而 PCA 选择样本点投影具有最大方差的方向。

向。

LDA 算法的主要优点有：

1) 在降维过程中可以使用类别的先验知识经验，而像 PCA 这样的无监督学习则无法使用类别先验知识。

2) LDA 在样本分类信息依赖均值而不是方差的时候，比 PCA 之类的算法较优。

LDA 算法的主要缺点有：

1) LDA 不适合对非高斯分布样本进行降维，PCA 也有这个问题。

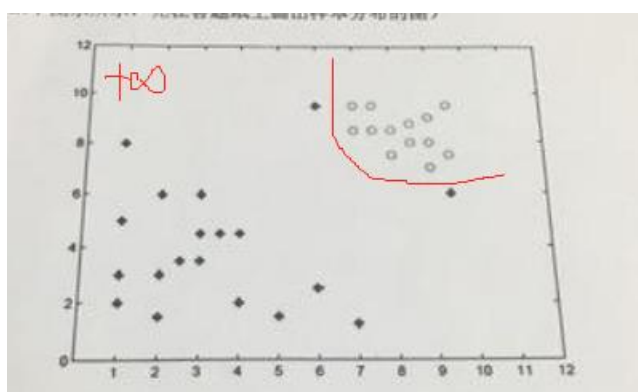
2) LDA 降维最多降到类别数  $k-1$  的维数，如果我们降维的维度大于  $k-1$ ，则不能使用 LDA。当然目前有一些 LDA 的进化版算法可以绕过这个问题。

3) LDA 在样本分类信息依赖方差而不是均值的时候，降维效果不好。

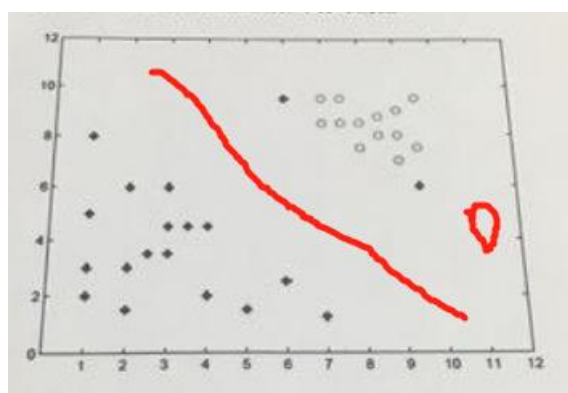
4) LDA 可能过度拟合数据。

## 9.SVM

### 1) C 取无穷大分类边界



### 2) C 取无穷小时的边界



3) 惩罚因子  $C$  越大，则 SVM 会更倾向把所有数据分对，往往出现较小的 margin，最终导致过拟合现象，泛化性能不好。

C 过于小，则惩罚力度不够，SVM 会更倾向实现最大化的 margin，而对样本分对分错不关心，不利于分类。

综上所述，在测试集中，C 取无穷小的时候效果会相对更好一些，泛化性能好。

#### 4) Soft-margin SVM 的原问题和对偶问题，并阐述核方法的基本思想是如何将线性模型转换非线性模型的

原问题：

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y_i(w \cdot x_i + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, N \end{aligned}$$

对偶问题：

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^N \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, N \end{aligned}$$

将在原始低维空间线性不可分的分类问题通过非线性变换成高维空间线性可分的分类问题，在高维空间学习线性支持向量机。在线性支持向量机学习的对偶问题中，把低维到高维的非线性变换的内积形式用核函数表示。