



软件开发环境国家重点实验室
State Key Laboratory of Software Development Environment

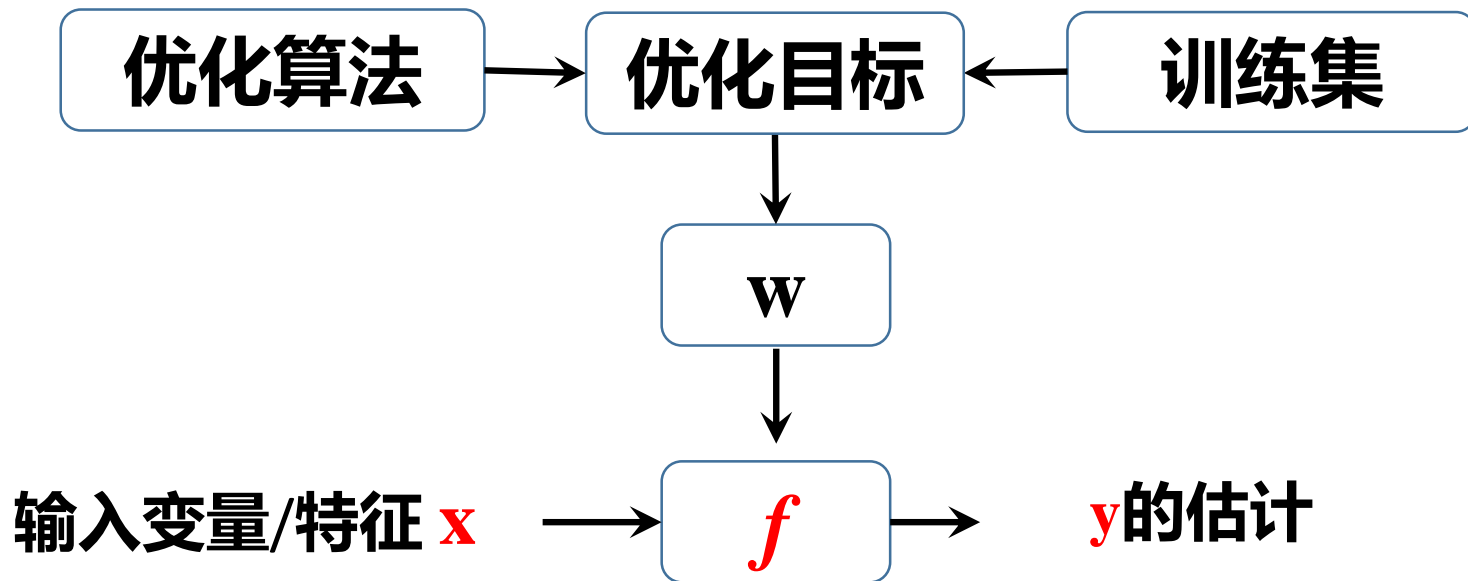
机器学习

刘祥龙

北京航空航天大学计算机学院
软件开发环境国家重点实验室

第四讲

线性回归



$$\begin{aligned} f(\mathbf{x}) &= w_0 + w_1x_1 + \dots + w_mx_m \\ &= w_0 + \sum_{i=1}^m w_ix_i \\ &= \sum_{i=0}^m w_ix_i \quad \boxed{x_0 = 1} \end{aligned}$$

一般地：

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

基函数

$\phi(x)$ 称为基函数 (Basis function)

一般情况下 $\phi_0(x) = 1$, 此时 w_0 为截矩或偏置 (bias)

最简单的情况下: $\phi_j(x) = x_j$

多元线性回归 - 最小二乘法

- 线性模型+最小均方误差+优化：最小二乘法

$$\hat{\mathbf{w}}^* = \arg \min_{\hat{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}^T) (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$$

令 $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$, 对 $\hat{\mathbf{w}}$ 求导得到

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})$$

令上式为零可得 $\hat{\mathbf{w}}$ 最优解的闭式解

多元线性回归 - 最小二乘法

- $\mathbf{X}^T \mathbf{X}$ 是满秩矩阵或正定矩阵, 则

$$\hat{\mathbf{w}}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

其中 $(\mathbf{X}^T \mathbf{X})^{-1}$ 是 $\mathbf{X}^T \mathbf{X}$ 的逆矩阵, 线性回归模型为

$$f(\hat{\mathbf{x}}_i) = \hat{\mathbf{x}}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbf{X}^T \mathbf{X}$ 不是满秩矩阵: 引入正则化

● 假定观测数据由确定的函数加高斯噪声组成：

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{where} \quad p(\epsilon|\beta) = \mathcal{N}(\epsilon|0, \beta^{-1})$$

$$\text{即:} \quad p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}).$$

似然函数

$$\longrightarrow p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}).$$

对数似然

$$\begin{aligned} \longrightarrow \ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \boxed{E_D(\mathbf{w})} \end{aligned}$$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$



● 求梯度并求零点的值：

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T = \mathbf{0}.$$

解得 \mathbf{w} ：

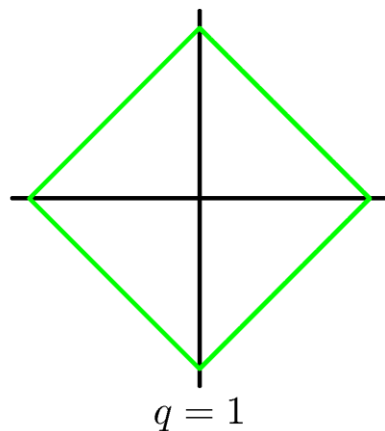
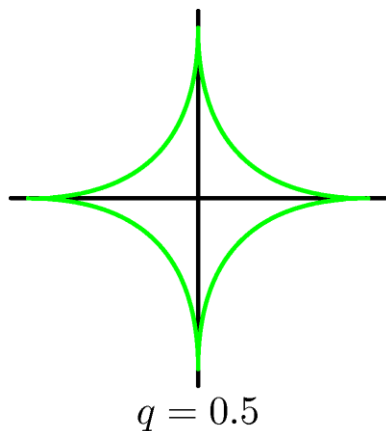
$$\mathbf{w}_{\text{ML}} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}.$$

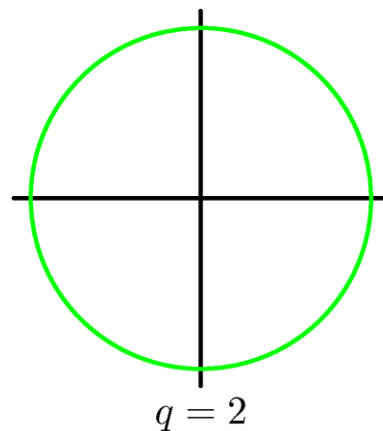
如何防止过拟合

- 加入正则项

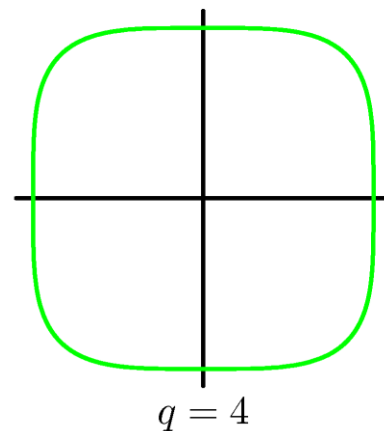
$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



Lasso



Quadratic



如何防止过拟合

- 加入正则项

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

二次正则项

求得 \mathbf{w} : $\mathbf{w} = \left(\lambda \mathbf{I} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}.$



软件开发环境国家重点实验室
State Key Laboratory of Software Development Environment

分类问题

离散值预测：二分类

- 线性回归： $z = \boldsymbol{w}^T \boldsymbol{x} + b$
- 二分类： $y \in \{0, 1\}$
- 如何建立分类与线性回归的联系？

- 最理想的函数——单位阶跃函数

$$y = \begin{cases} 0, & z < 0; \\ 0.5, & z = 0; \\ 1, & z > 0, \end{cases}$$

- 预测值大于零就判为正例，小于零就判为反例，预测值为临界值零则可任意判别
 - 缺点：不连续

二分类

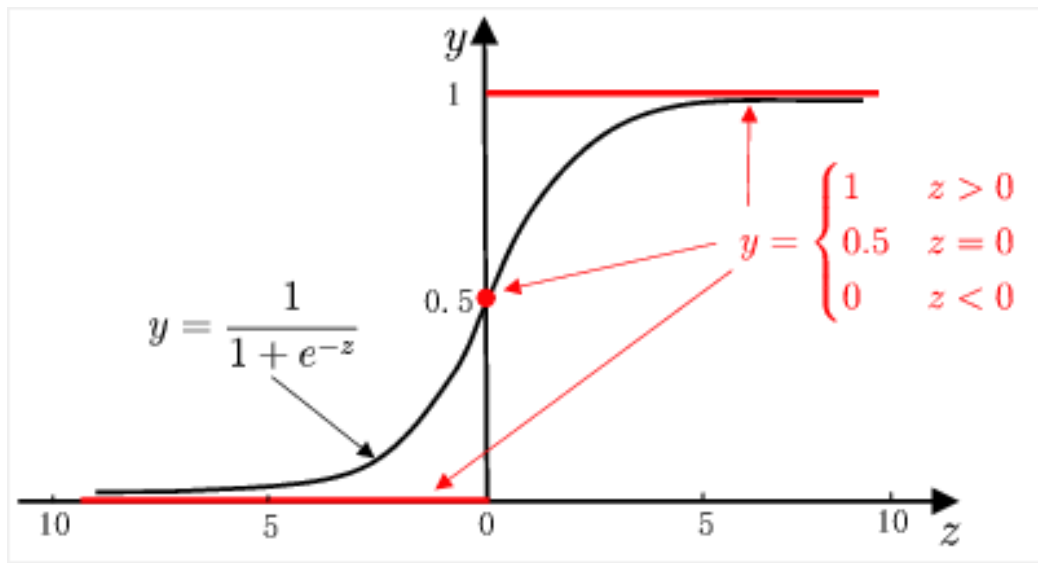
• 替代函数——逻辑函数 (logistic function)

□ 单调可微、任意阶可导

S函数

$$y = \frac{1}{1 + e^{-z}}$$

皮埃尔·弗朗索瓦·韦吕勒在
1844或1845年在研究它与人口
增长的关系时命名的



逻辑回归 (logistic regression)

- 运用逻辑函数

$$y = \frac{1}{1 + e^{-z}} \quad \text{变为} \quad y = \frac{1}{1 + e^{-(w^T x + b)}}$$

- 逻辑回归的优点

- 无需事先假设数据分布
- 可得到“类别”的近似概率预测
- 可直接应用现有数值优化算法求取最优解

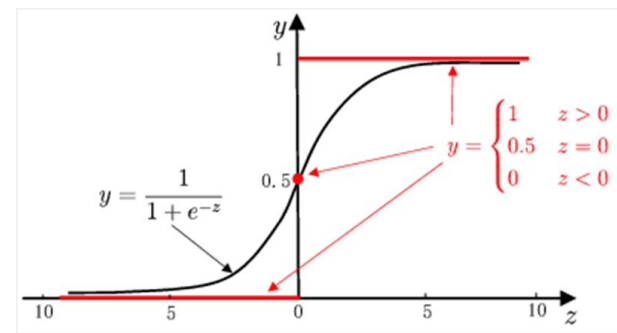
逻辑回归 - 极大似然法

- 分类概率

$$\ln \frac{p(y = 1 \mid \mathbf{x})}{p(y = 0 \mid \mathbf{x})} = \mathbf{w}^T \mathbf{x} + b$$

$$p(y = 1 \mid \mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{x} + b}}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$$

$$p(y = 0 \mid \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$$



- 极大似然法 (maximum likelihood)

- 给定数据集

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^m$$

- 最大化样本属于其真实标记的概率

- 最大化似然函数

$$\ell(\mathbf{w}, b) = \prod_{i=1}^m p(y_i | \mathbf{x}_i; \mathbf{w}_i, b)$$

- 最大化对数似然函数

$$\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i | \mathbf{x}_i; \mathbf{w}_i, b)$$

逻辑回归 - 极大似然法

- 转化为最小化负对数似然函数求解

- 令 $\beta = (w; b)$, $\hat{x} = (x; 1)$, 则 $w^T x + b$ 可简写为 $\beta^T \hat{x}$

- 再令 $p_1(\hat{x}_i; \beta) = p(y = 1 \mid \hat{x}_i; \beta)$

$$p_0(\hat{x}_i; \beta) = p(y = 0 \mid \hat{x}_i; \beta) = 1 - p_1(\hat{x}_i; \beta)$$

则似然项可重写为

$$p(y_i \mid x_i; w_i, b) = y_i p_1(\hat{x}_i; \beta) + (1 - y_i) p_0(\hat{x}_i; \beta)$$

- 目标函数等价于

$$\ell(\beta) = \sum_{i=1}^m \left(-y_i \beta^T \hat{x}_i + \ln \left(1 + e^{\beta^T \hat{x}_i} \right) \right)$$

逻辑回归 – 极大似然法

□ 优化求解 $\beta^* = \arg \min_{\beta} \ell(\beta)$

□ 牛顿法第t+1轮迭代解的更新公式

$$\beta^{t+1} = \beta^t - \left(\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\beta)}{\partial \beta}$$

其中关于 β 的一阶、二阶导数分别为

$$\frac{\partial \ell(\beta)}{\partial \beta} = - \sum_{i=1}^m \hat{\mathbf{x}}_i (y_i - p_1(\hat{\mathbf{x}}_i; \beta))$$

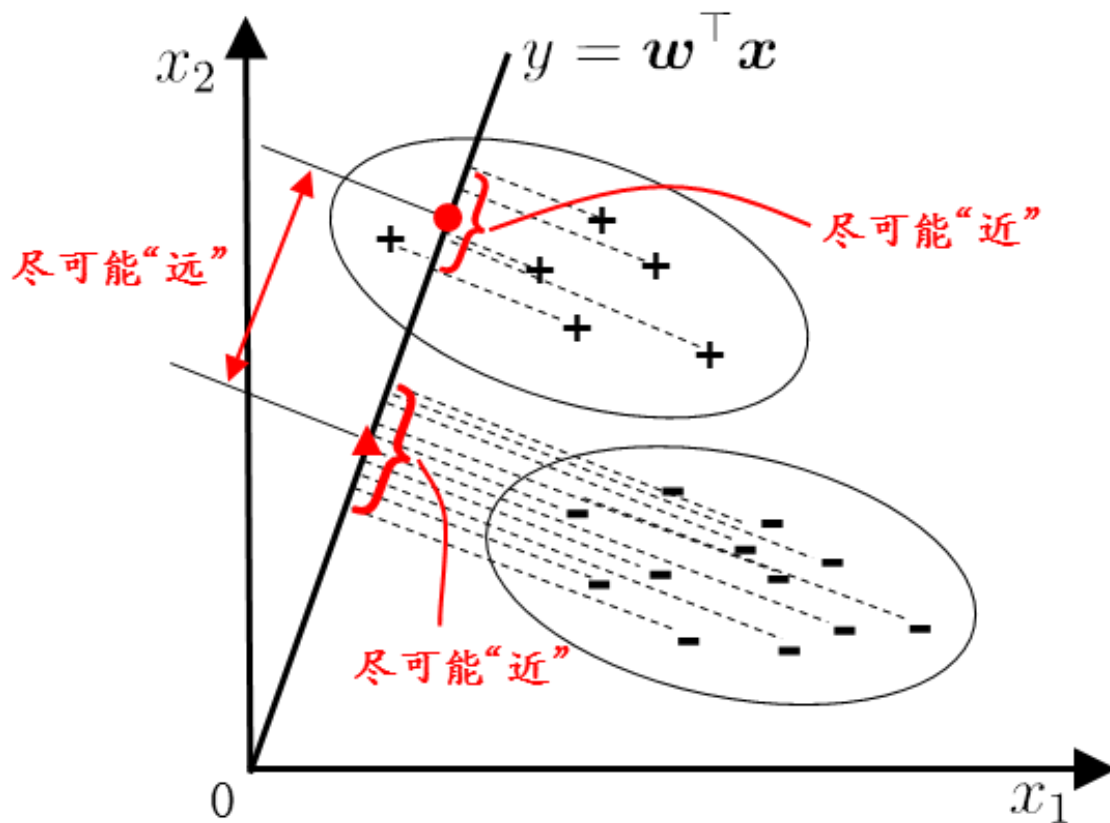
$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^m \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^T p_1(\hat{\mathbf{x}}_i; \beta) (1 - p_1(\hat{\mathbf{x}}_i; \beta))$$

高阶可导连续凸函数，梯度下降法/牛顿法 [Boyd and Vandenberghe, 2004]

其他二分类任务 - 线性判别分析

- 线性判别分析 (Linear Discriminant Analysis)

[Fisher, 1936]



LDA也可被视为一种
监督降维技术

• LDA的思想

- 欲使同类样例的投影点尽可能接近：可以让同类样例投影点的协方差尽可能小
- 欲使异类样例的投影点尽可能远离：可以让类中心之间的距离尽可能大

• 变量说明

- 第 i 类示例的集合 X_i
- 第 i 类示例的均值向量 μ_i
- 第 i 类示例的协方差矩阵 Σ_i
- 两类样本的中心在直线上的投影： $w^T \mu_0$ 和 $w^T \mu_1$
- 两类样本的协方差： $w^T \Sigma_0 w$ 和 $w^T \Sigma_1 w$

- 最大化目标

$$\begin{aligned} J &= \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} \\ &= \frac{w^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} \end{aligned}$$

- 类内散度矩阵

$$\begin{aligned} S_w &= \Sigma_0 + \Sigma_1 \\ &= \sum_{x \in X_0} (x - \mu_0) (x - \mu_0)^T + \sum_{x \in X_1} (x - \mu_1) (x - \mu_1)^T \end{aligned}$$

- 类间散度矩阵 $S_b = (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T$

- 广义瑞利商 (generalized Rayleigh quotient)

$$J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

- 令 $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$, 最大化广义瑞利商等价形式为

$$\begin{aligned} \min_{\mathbf{w}} \quad & -\mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1 \end{aligned}$$

- 运用拉格朗日乘子法 $\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$ $\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w} = \lambda \mathbf{w}$

线性判别分析

- 同向向量

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \quad S_B W = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T W = (\mu_1 - \mu_2) * \lambda_w$$

$$\underline{S_b w} = \lambda (\underline{\mu_0 - \mu_1})$$

- 结果

$$S_w^{-1} S_b w = \lambda w \longrightarrow w = S_w^{-1} (\mu_0 - \mu_1)$$

- 求解

- 奇异值分解

$$S_w = U \Sigma V^T$$

LDA推广 - 多分类任务

- 优化目标

$$\max_{\mathbf{W}} \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}$$

其中 $\mathbf{W} \in \mathbb{R}^{d \times (N-1)}$



$$\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$$

$$\mathbf{S}_b = \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$$

$$\mathbf{S}_w = \sum_{i=1}^N \mathbf{S}_{w_i}$$

$$\mathbf{S}_{w_i} = \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$$

\mathbf{W} 的闭式解则是 $\mathbf{S}_w^{-1} \mathbf{S}_b$ 的 $N-1$ 个最大广义特征值所对应的特征向量组成的矩阵

- 多分类LDA将样本投影到 $N-1$ 维空间, $N-1$ 通常远小于数据原有的属性数, 因此LDA也被视为一种监督降维技术

• 多分类学习方法

- 二分类学习方法推广到多类
- 利用二分类学习器解决多分类问题 (常用)
 - 对问题进行拆分, 为拆出的每个二分类任务训练一个分类器
 - 对于每个分类器的预测结果进行集成以获得最终的多分类结果

• 拆分策略

- 一对一 (One vs. One, OvO)
- 一对其余 (One vs. Rest, OvR)
- 多对多 (Many vs. Many, MvM)



• 拆分阶段

- N个类别两两配对
 - $N(N-1)/2$ 个二类任务
- 各个二类任务学习分类器
 - $N(N-1)/2$ 个二类分类器

• 测试阶段

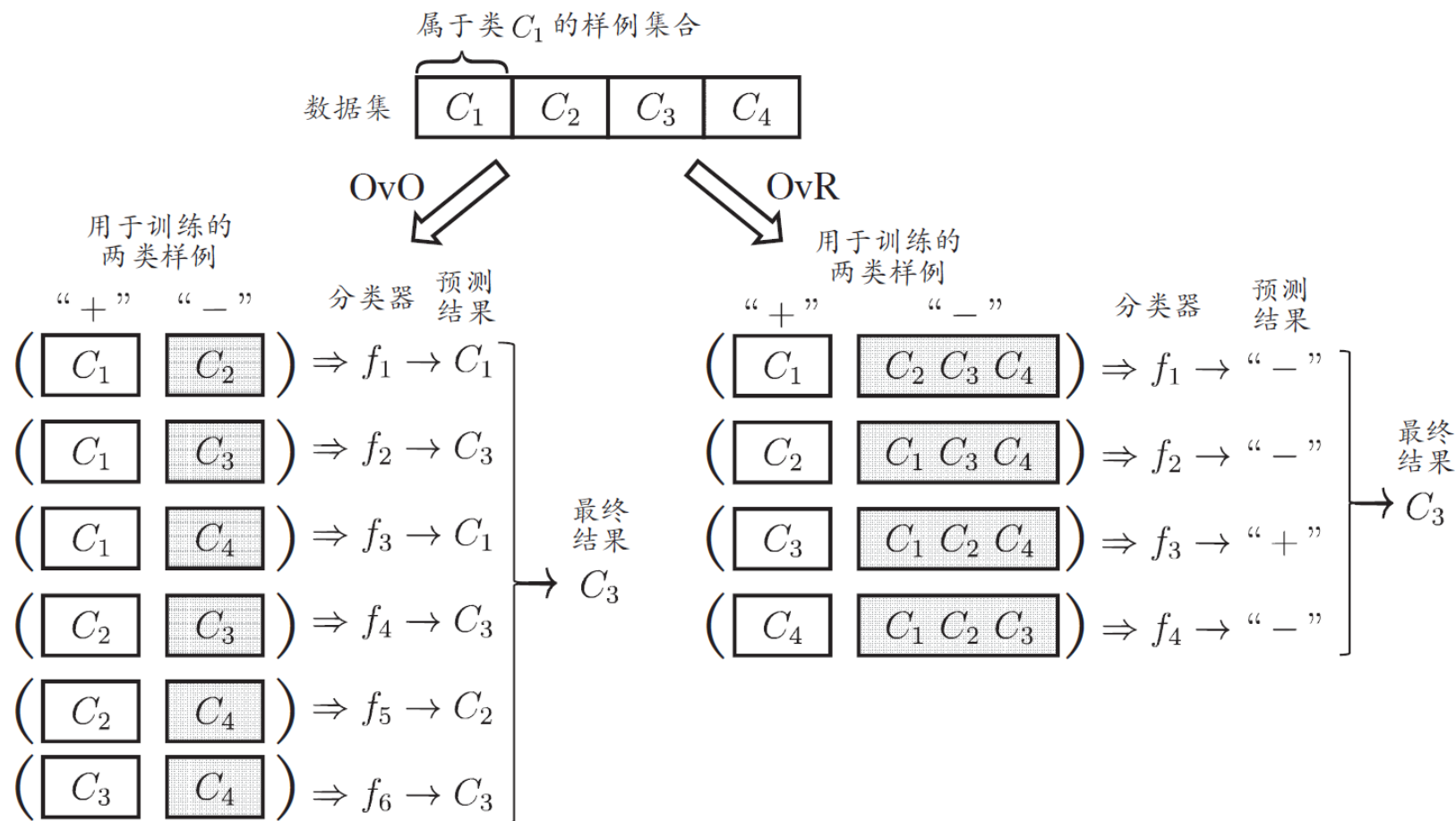
- 新样本提交给所有分类器预测
 - $N(N-1)/2$ 个分类结果
- 投票产生最终分类结果
 - 被预测最多的类别为最终类别

• 任务拆分

- ❑ 某一类作为正例，其他反例
 - N 个二类任务
- ❑ 各个二类任务学习分类器
 - N 个二类分类器

• 测试阶段

- ❑ 新样本提交给所有分类器预测
 - N 个分类结果
- ❑ 比较各分类器预测置信度
 - 置信度最大类别作为最终类别



多分类学习- 两种策略比较

- 一对一
- 训练 $N(N-1)/2$ 个分类器，存储开销和测试时间大
- 训练只用两个类的样例，训练时间短
- 一对其余
- 训练 N 个分类器，存储开销和测试时间小
- 训练用到全部训练样例，训练时间长

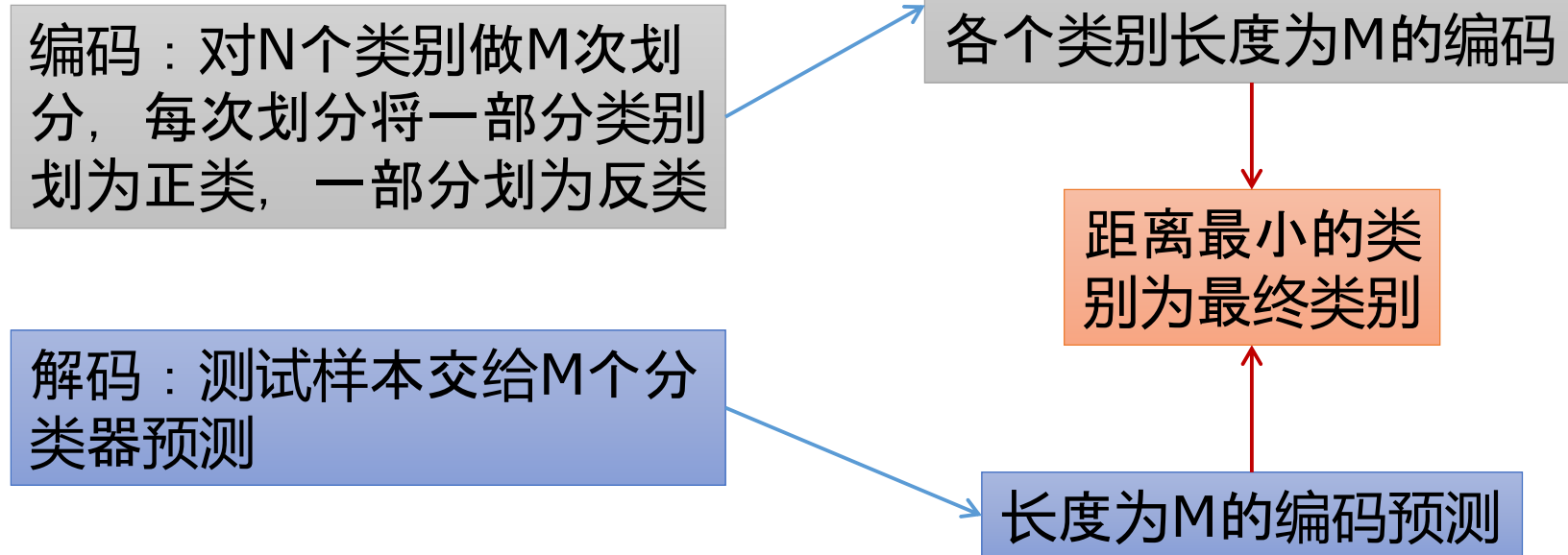
预测性能取决于具体数据分布，多数情况下两者差不多

多分类学习 - 多对多

- 多对多 (Many vs Many, MvM)

- 若干类作为正类, 若干类作为反类

- 纠错输出码 (Error Correcting Output Code, ECOC)



• 纠错输出码(Error Correcting Output Code, ECOC)

	f_1	f_2	f_3	f_4	f_5	海明距离	欧氏距离
	↓	↓	↓	↓	↓	↓	↓
$C_1 \rightarrow$	-1	+1	-1	+1	+1	3	$2\sqrt{3}$
$C_2 \rightarrow$	+1	-1	-1	+1	-1	4	4
$C_3 \rightarrow$	-1	+1	+1	-1	+1	1	2
$C_4 \rightarrow$	-1	-1	+1	+1	-1	2	$2\sqrt{2}$
测试示例 \rightarrow	-1	-1	+1	-1	+1	↑	↑

(a) 二元 ECOC 码

[Dietterich and Bakiri, 1995]

	f_1	f_2	f_3	f_4	f_5	f_6	f_7	海明距离	欧氏距离
	↓	↓	↓	↓	↓	↓	↓	↓	↓
$C_1 \rightarrow$	-1	-1	+1	+1	-1	+1	+1	4	4
$C_2 \rightarrow$	-1	0	0	0	+1	-1	0	2	2
$C_3 \rightarrow$	+1	+1	-1	-1	-1	+1	-1	5	$2\sqrt{5}$
$C_4 \rightarrow$	-1	+1	0	+1	-1	0	+1	3	$\sqrt{10}$
测试示例 \rightarrow	-1	+1	+1	-1	+1	-1	+1	↑	↑

(b) 三元 ECOC 码

[Allwein et al. 2000]

- ECOC编码对分类器错误有一定容忍和修正能力，编码越长、纠错能力越强
- 对同等长度的编码，理论上来说，任意两个类别之间的编码距离越远，则纠错能力越强



软件开发环境国家重点实验室
State Key Laboratory of Software Development Environment

支持向量机

支持向量机

- C. Cortes和V. Vapnik (1995年提出)
 - ❑ 支持向量机是基于统计学习理论(Statistical Learning Theory, SLT)发展起来的一种新的机器学习的方法。
 - ❑ 统计学习理论主要创立者是Vladimir N. Vapnik。



Google



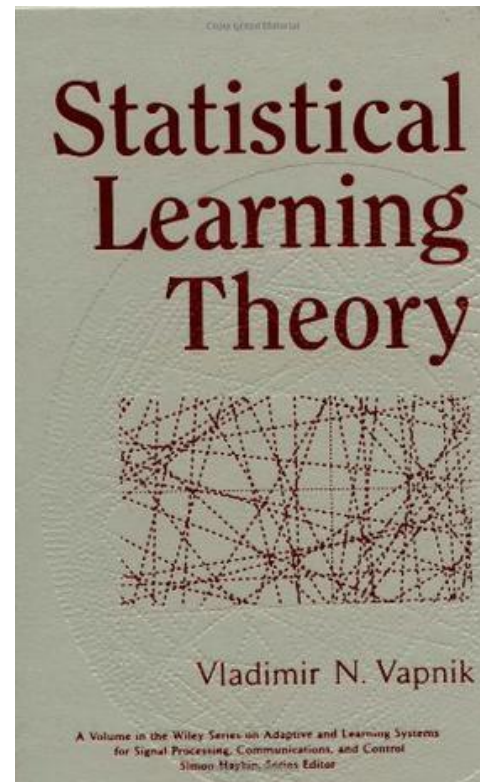
Vladimir N. Vapnik

- ❑ 1936年 出生于苏联
- ❑ 1958年 乌兹别克国立大学 硕士
- ❑ 1964年 莫斯科控制科学学院 博士
- ❑ 1964-1990年 莫斯科控制科学学院
- ❑ 1991-2001年 美国AT&T贝尔实验室
发明支持向量机理论
- ❑ 2002-2014年 NEC实验室(美国)
从事机器学习研究
- ❑ 1995年和2003年 伦敦大学皇家霍洛威学院和美国哥伦比亚大学计算机专业的教授
- ❑ 2006年 美国国家工程院院士
- ❑ 2014年至今 美国Facebook公司 从事人工智能研究
- ❑ 2017年 IEEE John von Neumann Medal



V. Vapnik对于统计机器学习的贡献

- 1968年，Vapnik和Chervonenkis提出了VC熵和VC维的概念，这些是统计学习理论的核心概念。同时，他们发现了泛函空间的大数定理，得到了关于收敛速度的非渐进界的主要结论。
- 1974年，Vapnik和Chervonenkis提出了结构风险最小化归纳原则。
- 1989年，Vapnik和Chervonenkis发现了经验风险最小化归纳原则和最大似然方法一致性的充分必要条件，完成了对经验风险最小化归纳推理的分析。
- 90年代中期，有限样本情况下的机器学习理论研究逐渐成熟起来，形成了较完善的理论体系——统计学习理论。



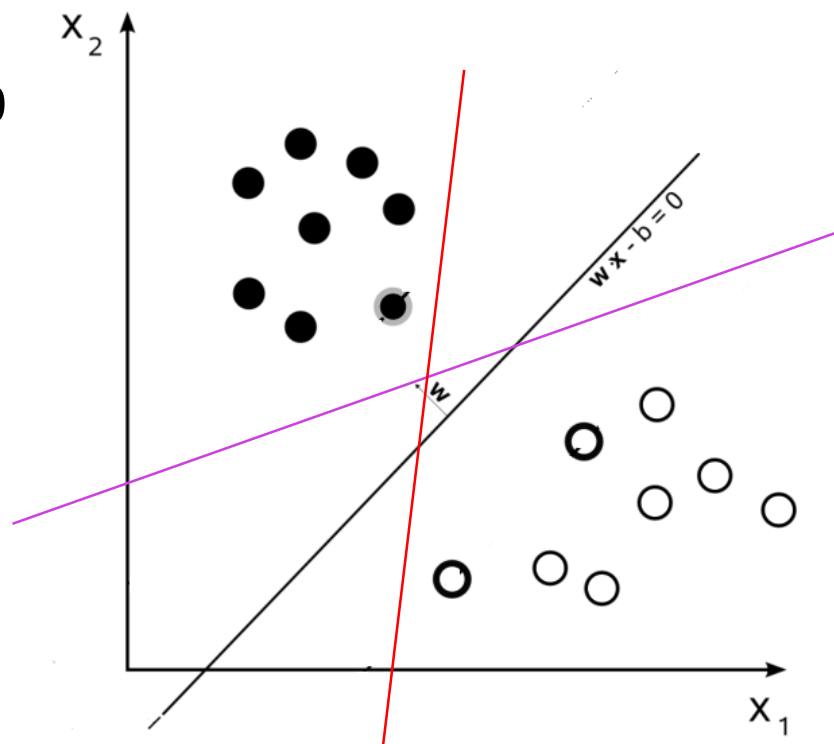
线性可分支持向量机：函数

• 训练数据集

- $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}, x_i \in R^n, y_i \in \{+1, -1\}$

• 二分类目标

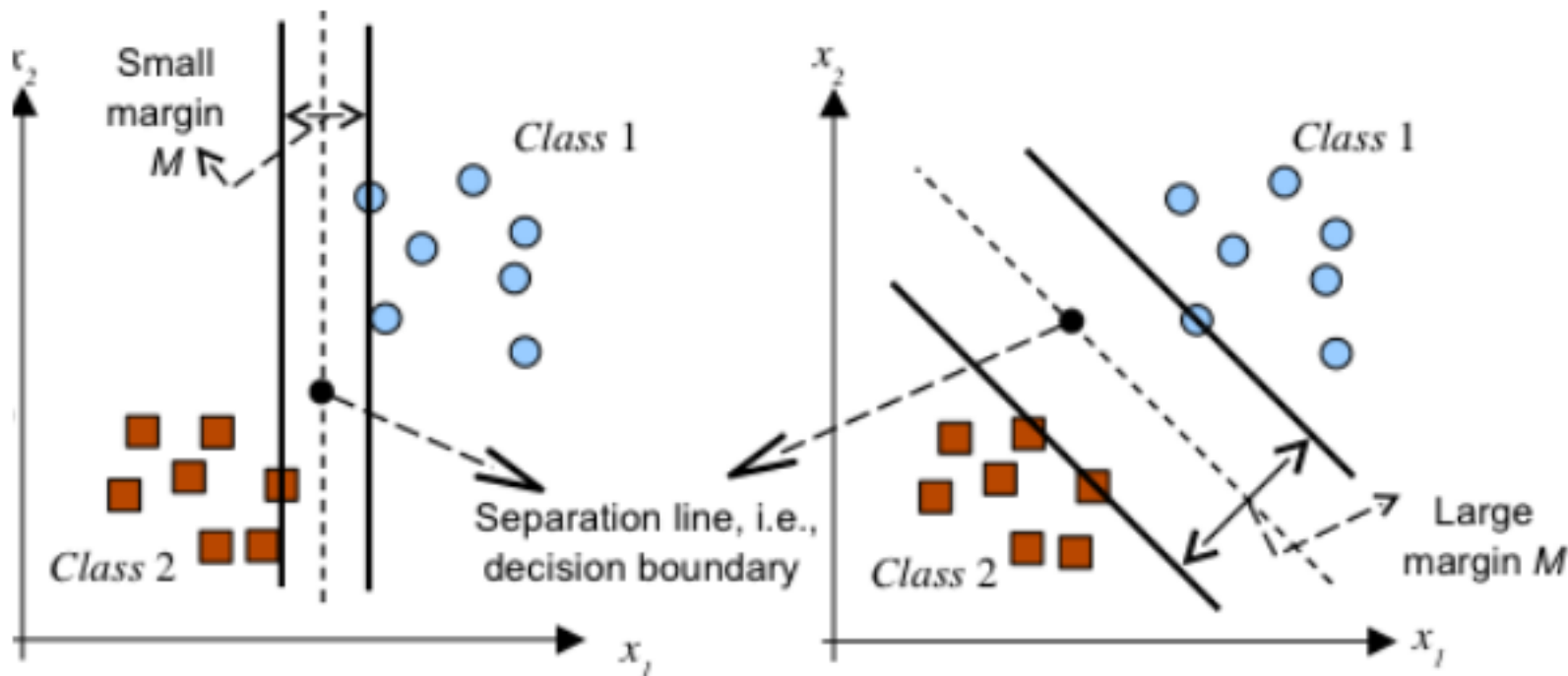
- 函数：分离超平面 $w^T x_i + b = 0$
- 二分类函数 $y_i = \text{sign}(w^T x_i + b)$
- 若线性可分
 - 感知机：存在无数分离超平面
 - 线性可分支持向量机：唯一解



支持向量机：优化目标

• 线性可分的最优分类面

- 最优分类面就是要求分类线不但能将两类正确分开(训练错误率为0)，且使分类间隔最大。SVM考虑寻找一个满足分类要求的超平面，并且使训练集中的点距离分类面尽可能的远，也就是寻找一个分类面使它两侧的空白区域(Margin)最大。



• 函数间隔

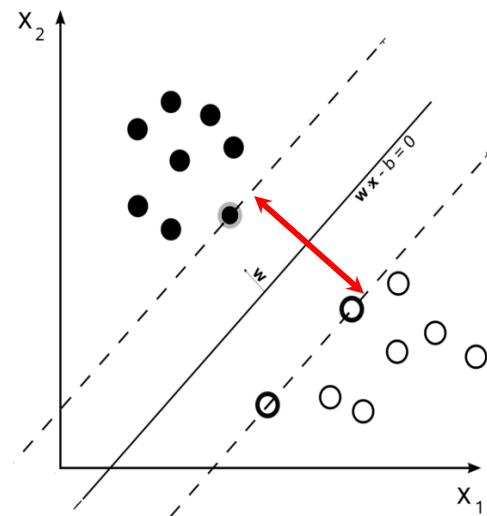
- $|w^T x + b|$ 表示点相对分离超平面的距离，反映分类确信度
- $y(w^T x + b)$ 同时反映预测与实际分类是否一致及确信度
- 数据集 T 上函数间隔

$$\min_{i=1,2,\dots,N} y_i(w^T x_i + b)$$

• 几何间隔

- $(w, b) \rightarrow (2w, 2b)$ 时超平面不变，函数间隔发生变化
- 函数间隔 \rightarrow 几何间隔，可约束 $\|w\| = 1$

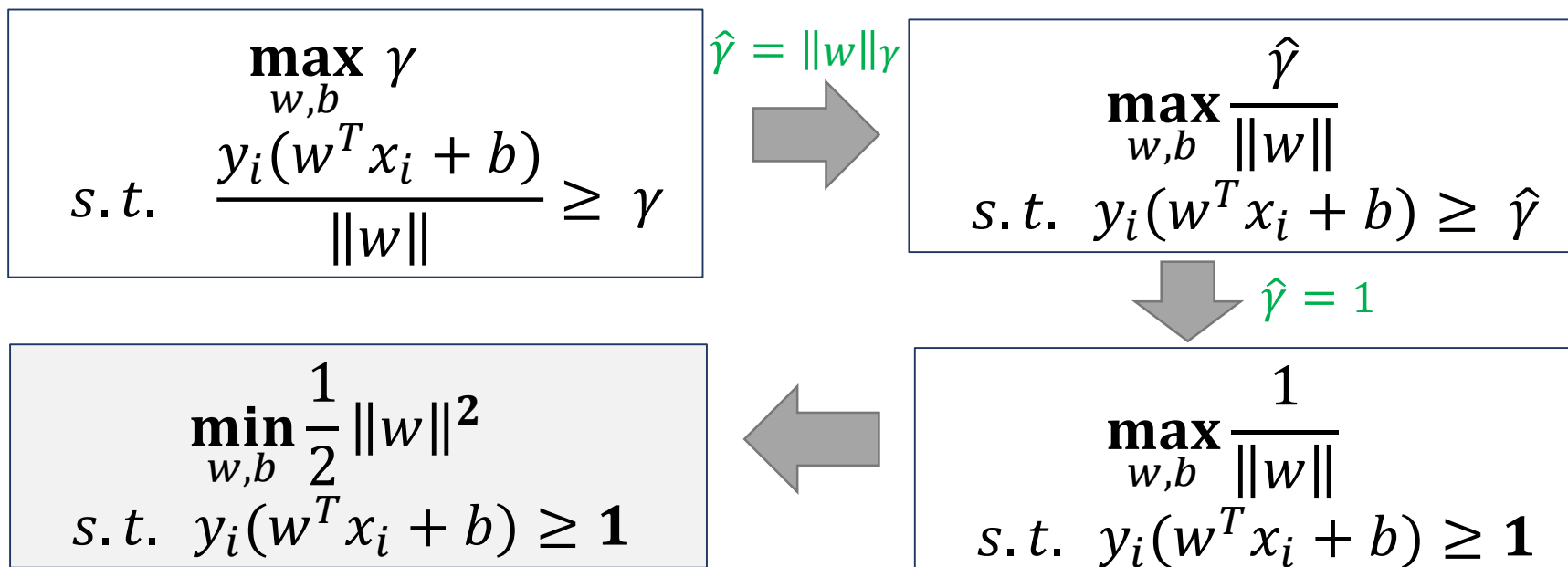
$$\gamma = \min_{i=1,2,\dots,N} \frac{y_i(w^T x_i + b)}{\|w\|}$$



• 间隔最大化

- 正确划分, 几何间隔最大
- 直观解释: 充分大的确信度对数据进行分类
 - 最难分的点以足够大的确信度分开
 - 对未知新数据有很好的泛化能力

• 最大间隔分类超平面



线性可分支持向量机：优化求解

• 可解

- 凸优化: f, g 凸函数, h 仿射函数

$$\begin{aligned} & \min_w f(w) \\ \text{s.t. } & g_i(w) \leq 0 \\ & h_i(w) = 0 \end{aligned}$$

$$\begin{aligned} & \min_{w,b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } & y_i(w^T x_i + b) \geq 1 \end{aligned}$$

• 解存在唯一

- 存在性：目标函数有解 (w^*, b^*) , 且非 $(w^*, b^*) \neq (0, b)$
- 唯一性：假设存在两个解 $\rightarrow w^* = 0$ 或者两个解相同 $\rightarrow w^*$ 唯一

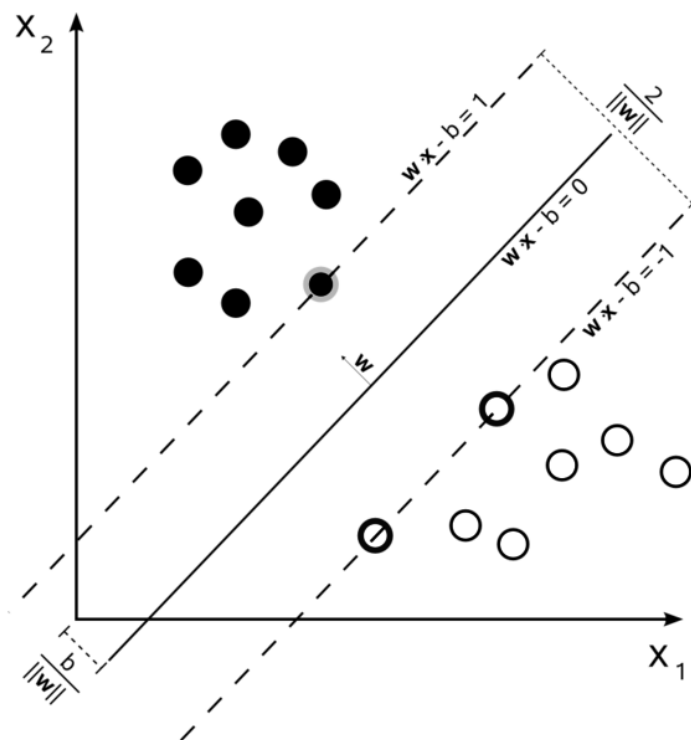
线性可分支持向量机

• 支持向量

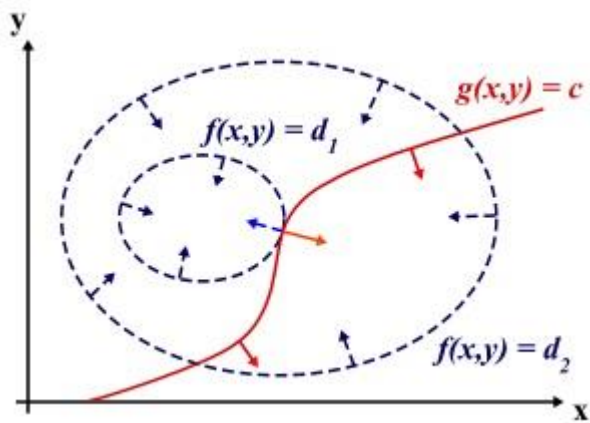
- 使 $y_i(w^T x_i + b) = 1$ 成立的 x_i
- 支持向量所在平面为 $w^T x + b = \pm 1$
- 两平面距离称为间隔

• 特点

- 只有支持向量决定分离平面
- 支持向量较少，具有稀疏性



线性可分支持向量机

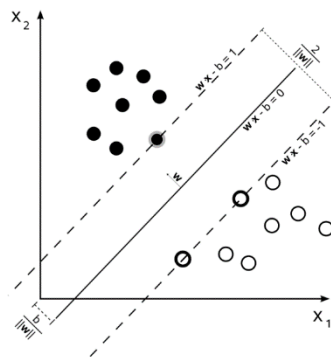


等式约束的优化问题

$$\begin{aligned} \min \quad & y = f(x_1, x_2) \\ \text{s.t.} \quad & g(x_1, x_2) = 0 \end{aligned}$$



$$y_\lambda = f(x_1, x_2) + \lambda(g(x_1, x_2))$$



不等式约束的优化问题?

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1 \end{aligned}$$

$$\begin{aligned} \min \quad & f_0(x) \\ \text{s.t.} \quad & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p \end{aligned}$$



$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

对偶问题

$$\min f_0(x)$$

$$\text{s. t. } f_i(x) \leq 0, \quad i = 1, \dots, m$$

$$h_i(x) = 0, \quad i = 1, \dots, p$$

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

定义 Lagrange 对偶函数为 Lagrange 函数关于 x 取最小值，即对于 $\lambda \in \mathbf{R}^m$ ， $\nu \in \mathbf{R}^p$ ，有

$$g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu) = \inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x))。$$

Lagrange 对偶函数是一族关于 (λ, ν) 的仿射函数的逐点下确界，即使原问题不是凸的，对偶函数也是凹函数。

对偶问题

可以证明, Lagrange 对偶函数构成了原问题(1)最优值 p^* 的下界,

即对任意 $\lambda \succeq 0$ 和 ν 有

$$\inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)) \quad g(\lambda, \nu) \leq p^*$$

$$\begin{aligned} \min f_0(x) \\ \text{s.t. } f_i(x) \leq 0, \quad i = 1, \dots, m \\ h_i(x) = 0, \quad i = 1, \dots, p \end{aligned}$$

求解从 Lagrange 对偶函数求得原函数最优值的最好下界, 构成

Lagrange 对偶问题

$$\begin{aligned} \max g(\lambda, \nu) \\ \text{s.t. } \lambda \succeq 0 \end{aligned}$$

Lagrange 对偶函数是凹函数, 所以 Lagrange 对偶问题是凸优化问题。

我们用 d^* 表示 Lagrange 对偶问题的最优解，则有

弱对偶性

$$d^* \leq p^*$$

即使原问题不是凸优化问题，这个不等式也成立。这个性质称作弱对偶性。

强对偶性

如果 $d^* = p^*$ ，则称强对偶性成立。如果原问题是凸优化问题，则通常（但不总是）强对偶性成立。

KKT条件

$$\inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p v_i h_i(x)) = p^*$$

$$\begin{aligned} \min f_0(x) \\ \text{s. t. } f_i(x) \leq 0, \quad i = 1, \dots, m \\ h_i(x) = 0, \quad i = 1, \dots, p \end{aligned}$$

当强对偶性成立时, x^* 和 (λ^*, v^*) 分别是原问题和对偶问题的最优解, 则其必须满足 KKT 条件:

$$f_i(x^*) \leq 0, \quad i = 1, \dots, m$$

$$h_i(x^*) = 0, \quad i = 1, \dots, p$$

$$\lambda_i^* \geq 0, \quad i = 1, \dots, m$$

$$\lambda_i^* f_i(x^*) = 0, \quad i = 1, \dots, m$$

$$\nabla f_0(x^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(x^*) + \sum_{i=1}^p v_i^* \nabla h_i(x^*) = 0$$

$$\begin{aligned} \min_{w,b} & \frac{1}{2} \|w\|^2 \\ \text{s. t. } & y_i(w^T x_i + b) \geq 1 \end{aligned}$$

• 对偶算法

- 拉格朗日函数：拉格朗日乘子 $\alpha = (\alpha_1, \dots, \alpha_N)^T$

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i y_i (w^T x_i + b) + \sum_{i=1}^N \alpha_i$$

- 对偶问题： $\min_{w,b} \max_{\alpha} L(w, b, \alpha) \rightarrow \max_{\alpha} \min_{w,b} L(w, b, \alpha)$

• 对偶算法

□ 对偶问题

$$\max_{\alpha} \min_{w, b} L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i y_i (w^T x_i + b) + \sum_{i=1}^N \alpha_i$$

(1) 求 $\min_{w, b} L(w, b, \alpha)$: 对 w, b 分别求导

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

$$L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i^T x_j) + \sum_{i=1}^N \alpha_i$$

• 对偶算法

□ 对偶问题: $\max_{\alpha} \min_{w,b} L(w, b, \alpha)$

(2) 求 $\max_{\alpha} \min_{w,b} L(w, b, \alpha)$

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i^T x_j) + \sum_{i=1}^N \alpha_i$$

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

□ 最优解: 由KKT条件 $w^* = \sum_{i=1}^N \alpha_i^* y_i x_i$ 和 $\alpha_i^* (y_i (w^{*T} x_i + b) - 1) = 0$, 存在 $\alpha_i^* > 0$

$$w^* = \sum_{i=1}^N \alpha_i^* y_i x_i, \quad b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i (x_i^T x_j)$$

超平面法向量 w^* 是支持向量的线性组合

□ 决策函数: $f(x) = \text{sign}(\sum_{i=1}^N \alpha_i^* y_i (x^T x_i) + b^*), \alpha_i^* > 0$

- 支持向量

- KKT互补条件 $\alpha_i^*(y_i(w^T x_i + b) - 1) = 0$
- 决策函数 $f(x) = \text{sign}(\sum_{i=1}^N \alpha_i^* y_i (x^T x_i) + b^*)$, $\alpha_i^* > 0$
- 对应 $\alpha_i^* > 0$ 的样例: $y_i(w^T x_i + b) = 1$, 落在间隔边界

- 示例

<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

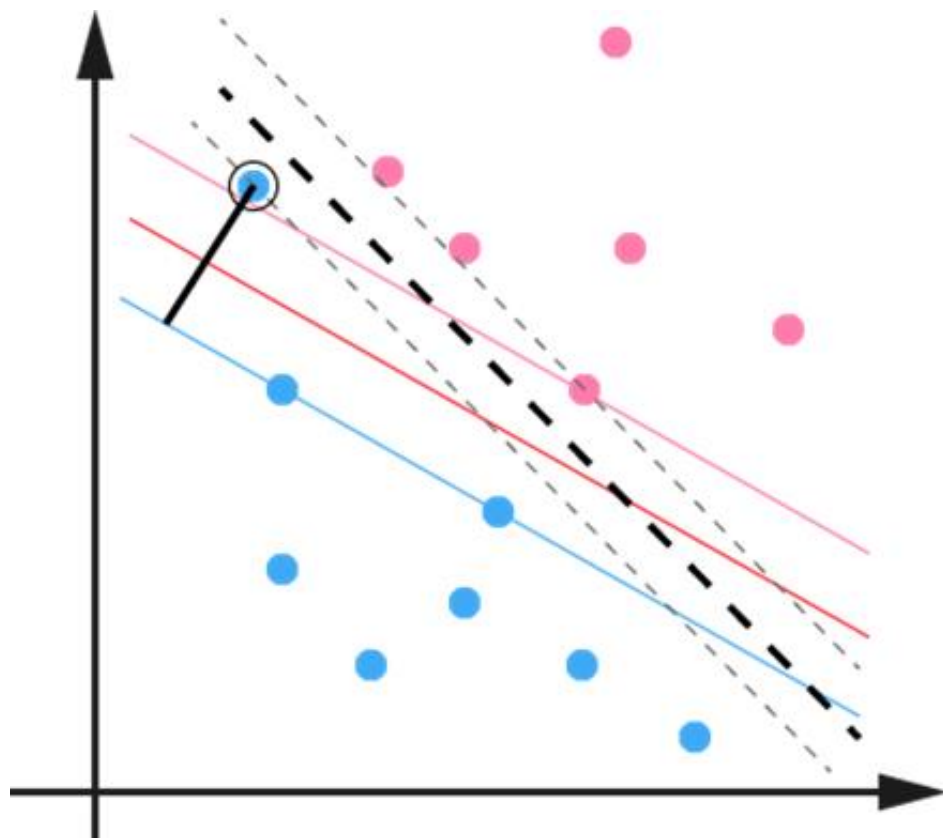
• 处理噪声和离群点

- 求解最优分类面的时间代价大还可能导致泛化性能差。因此，对于分布有交集的数据需要有一定范围内的“错分”，又有较大分界区域的广义最优分类面。

准确性



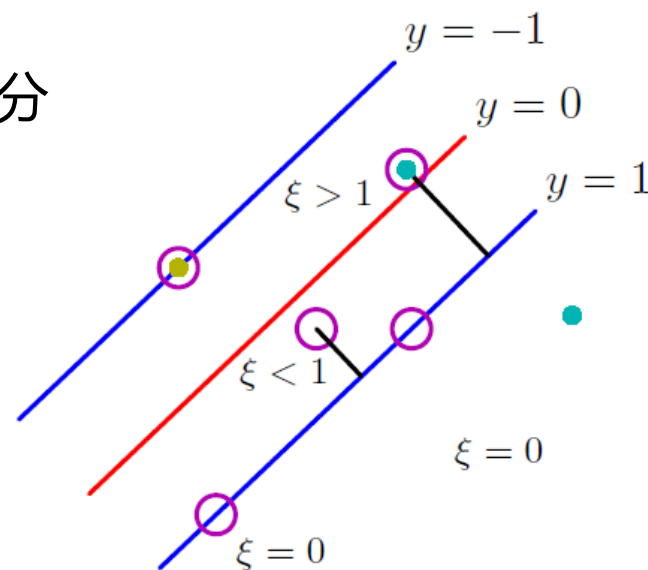
泛化性



线性支持向量机

- 线性不可分训练数据集

- 通常存在异常点，去除后可线性可分
- 硬间隔无法满足



- 目标函数

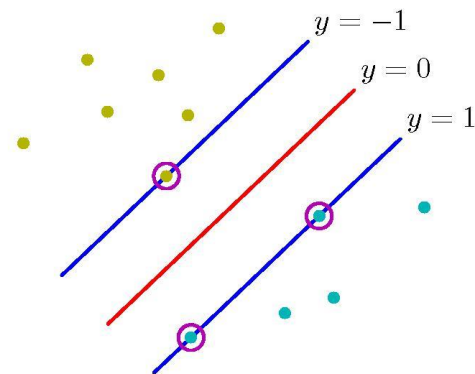
$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 + c \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1 - \xi_i \end{aligned}$$

• 处理噪声和离群点

- 这种处理方式也被视为是从硬间隔(Hard Margin)向软间隔(Soft Margin)的转变。

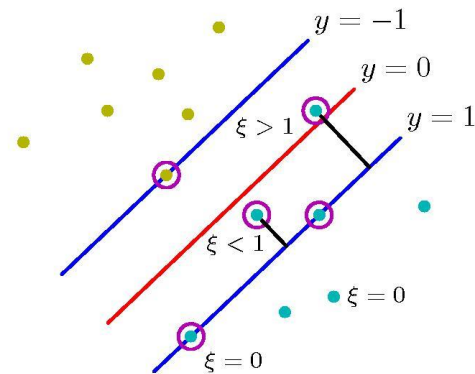
硬间隔

$$\min_{w,b} \frac{1}{\|w\|^2}$$
$$s.t. \ t_n(w^T x_n + b) \geq 1, \ n = 1, \dots, N$$



软间隔

$$\min_{w,b,\xi} \frac{1}{\|w\|^2} + C \sum_{n=1}^N \xi_n$$
$$s.t. \ t_n(w^T x_n + b) \geq 1 - \xi_n, \ n = 1, \dots, N$$
$$\xi_n \geq 0$$



• 对偶算法

- 拉格朗日函数求 $L(w, b, \xi, \alpha, \mu)$
- 求解对偶
 - 求 $\min_{w, b, \xi} L(w, b, \xi, \alpha, \mu)$
 - 求 $\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i^T x_j) + \sum_{i=1}^N \alpha_i$, 满足 $C - \alpha_i - \mu_i = 0$
- 最优解: 由KKT条件 $w^* = \sum_{i=1}^N \alpha_i^* y_i x_i$, $\mu_i^* \xi_i^* = 0$, 及 $\alpha_i^* (y_i (w^T x_i +$

超平面法向量 w^* 是支持向量的线性组合

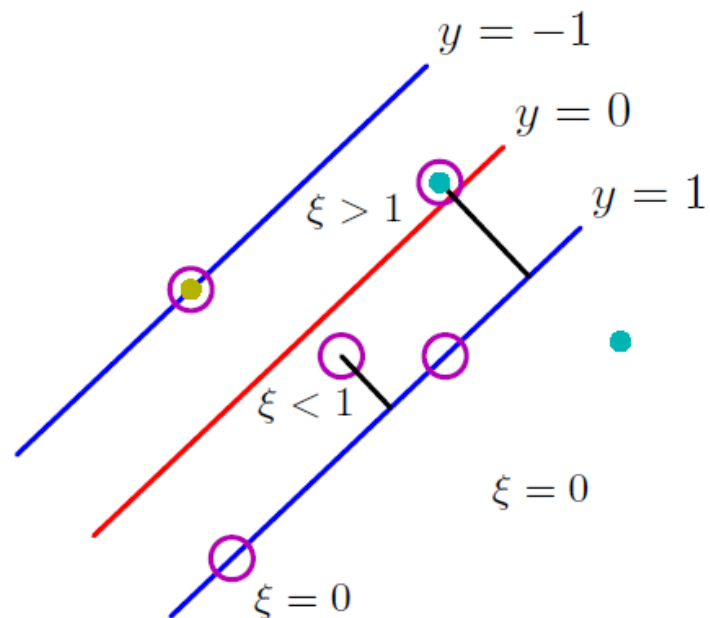
线性支持向量机

• 支持向量 $\mu_i^* \xi_i^* = 0, \quad \alpha_i^* (y_i (w^T x_i + b) - 1 + \xi_i) = 0, \quad C - \alpha_i - \mu_i = 0$

- $0 < \alpha_i^* < C: y_i (w^T x_i + b) - 1 + \xi_i = 0$
- $\alpha_i^* < C \rightarrow \xi_i = 0, x_i$ 落在间隔边界
- $\alpha_i^* = C \rightarrow \xi_i < 1, \quad \xi_i = 1, \quad \xi_i > 1$

• 示例

$$C - \alpha_i - \mu_i = 0$$
$$\mu_i^* \xi_i^* = 0$$



• 凸二次规划对偶问题

- 样本容量大时，直接求解效率较低
- 序列最小最优化 (SMO) 算法求解

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

• 启发式算法

- 所有变量满足KKT条件，选择最不满足条件的两个变量，固定其他变量，迭代求解。
- 两变量构成二次规划问题

- J. C. Platt(1999年提出)

- 支持向量机的学习问题可以形式化为求解具有全局最优解的凸二次规划问题。许多方法可以用于求解这一问题，但当训练样本容量很大时，这些算法往往效率较低，以致无法使用。
- 序列最小优化算法(Sequential Minimal Optimization, SMO)是一种启发式算法。基本思想是：如果所有变量都满足此优化问题的KKT条件，那么这个问题的解就得到了。
- SMO算法的特点是不不断地将原二次规划问题分解为只有两个变量的二次规划问题，并对子问题进行解析求解，直到所有变量都满足KKT条件为止。因为子问题解析解存在，所以每次计算子问题都很快，虽然子问题次数很多，但是总体上还是高效的。

• 两变量二次规划

$$\min_{\alpha_1, \alpha_2} W(\alpha_1, \alpha_2)$$

$$s.t. \alpha_1 y_1 + \alpha_2 y_2 = -\sum_{i=3} y_i \alpha_i, \quad C \geq \alpha_i \geq 0, i = 1, 2$$

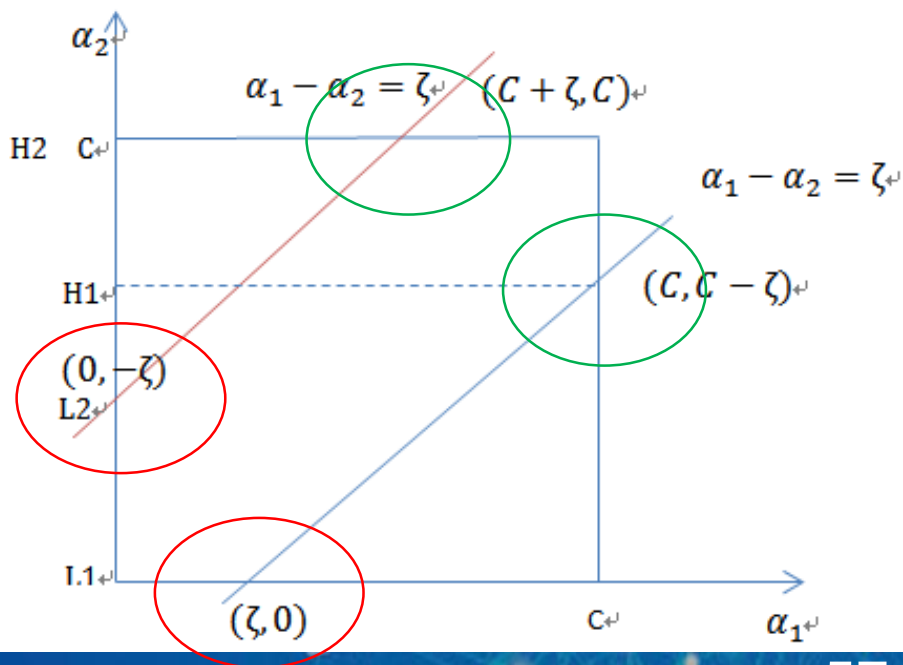
□ 迭代求解: $(\alpha_1^{old}, \alpha_2^{old}) \rightarrow (\alpha_1^{new}, \alpha_2^{new})$

□ $\frac{\partial W}{\partial \alpha_2} = 0 \Rightarrow \alpha_2^{new} = \alpha_2^{old} + \frac{y_2(E_1 - E_2)}{\eta}$

• 剪辑 α_2^{new} , 变量约束范围

□ y_1, y_2 异号: 最小值, 最大值

□ y_1, y_2 同号: 最大值, 最小值



序列最小优化算法

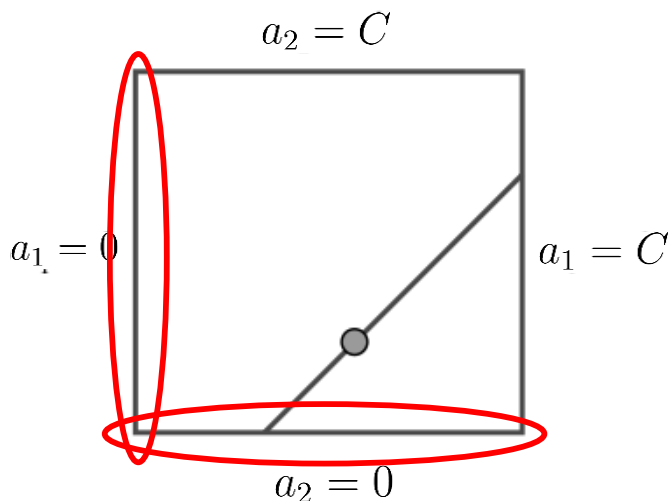
• 两个变量二次规划的解析方法

- 两个变量(a_1, a_2)的约束可用二维空间中的图形表示

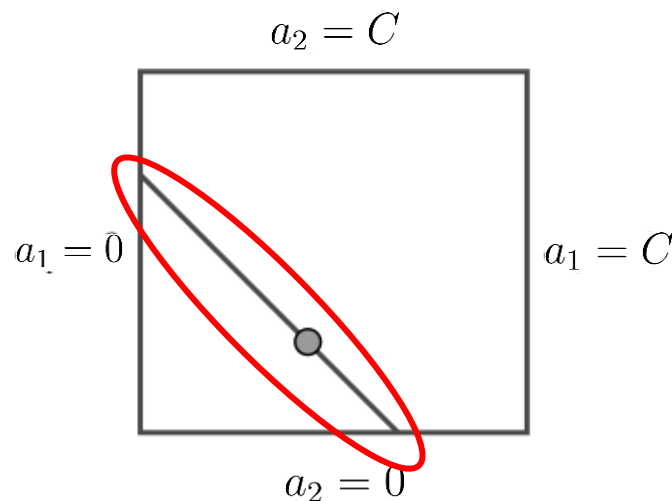
$$a_1 t_1 + a_2 t_2 = - \sum_{n=3}^N t_n a_n = \zeta$$

$$0 \leq a_n \leq C, i = 1, 2$$

目标函数在一条平行于对角线的线段上的最优值
两个变量的最优化问题转化为单变量最优化问题



$$t_1 \neq t_2 \Rightarrow a_1 - a_2 = k$$



$$t_1 = t_2 \Rightarrow a_1 + a_2 = k$$

• 变量选取

- 第一个变量选择（外层循环）：违反KKT条件最严重的样本点
 - $0 < \alpha_i^* < C$: $y_i(w^T x_i + b) = 1$
 - $\alpha_i^* = 0$: $y_i(w^T x_i + b) \geq 1$
 - $\alpha_i^* = C$: $y_i(w^T x_i + b) \leq 1$
- 第二个变量选择（内层循环）：使得 α_2 变化足够大
 - $\alpha_2^{new} = \alpha_2^{old} + \frac{y_2(E_1 - E_2)}{\eta}$, $E_i = [\sum_{i=1}^N \alpha_i^* y_i K(x, x_i) + b^*] - y_i$
 - $E_1 - E_2$ 足够大
- 更新 b 和 E_i
- 重复上述步骤，直到KKT条件满足

• SMO算法解凸二次规划问题

$$\min_a \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(x_n, x_m) - \sum_{n=1}^N a_n$$

$$s.t. \ 0 \leq a_n \leq C, n = 1, 2, \dots, N$$

$$\sum_{n=1}^N a_n t_n = 0$$

- 子问题有两个变量，一个是违反KKT条件最严重的，另一个有约束条件自动确定。两个变量中只有一个是自由变量。假设 a_1, a_2 为两个变量， a_3, a_4, \dots, a_N 固定，那么：

$$a_1 = -t_1 \sum_{n=2}^N a_n t_n = 0$$

- 即 a_2 确定， a_1 也随之确定。
- SMO算法包括：求解两个变量二次规划的解析方法和选择变量的启发式方法。

• 两个变量二次规划的解析方法

- 不失一般性，假设选择的两个变量是 a_1 和 a_2 ，其他变量 a_i ($i=3, 4, \dots, N$)是固定的，原问题的子问题可以写成：

$$\begin{aligned} \min_{a_1, a_2} W(a_1, a_2) &= \frac{1}{2}K_{11}a_1^2 + \frac{1}{2}K_{22}a_2^2 + t_1t_2K_{12}a_1a_2 \\ &\quad - (a_1 + a_2) + t_1a_1 \sum_{n=3}^N t_na_nK_{n1} + t_2a_2 \sum_{n=3}^N t_na_nK_{n2} + const \\ s.t. \quad a_1t_1 + a_2t_2 &= - \sum_{n=3}^N t_na_n = \zeta \\ 0 \leq a_n &\leq C, i = 1, 2 \end{aligned}$$

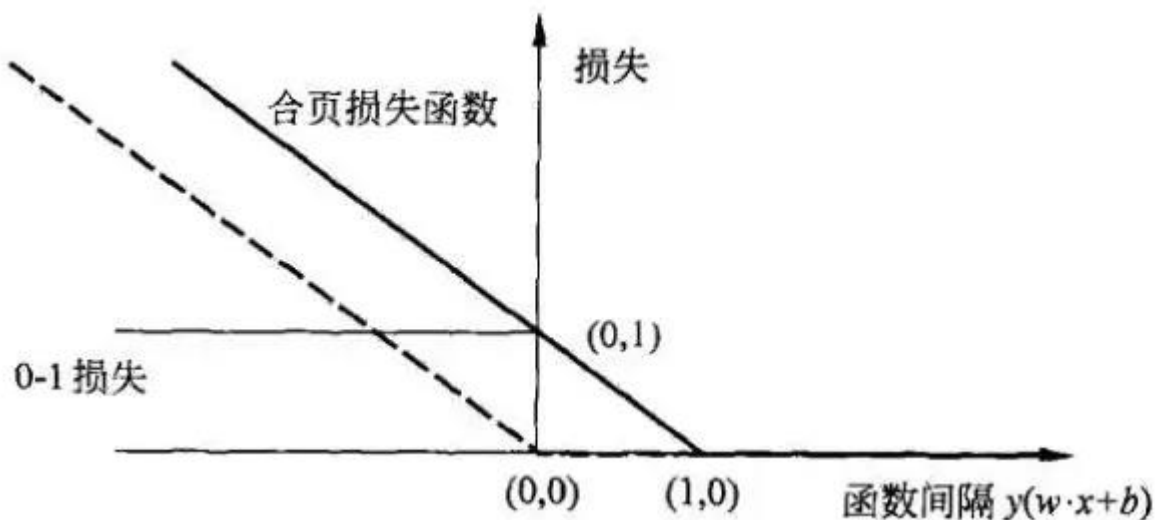
- 其中， $K_{mn} = K(x_m, x_n), m, n = 1, 2, \dots, N$ ， ζ 是常数。

合页损失函数

- 另外一种解释

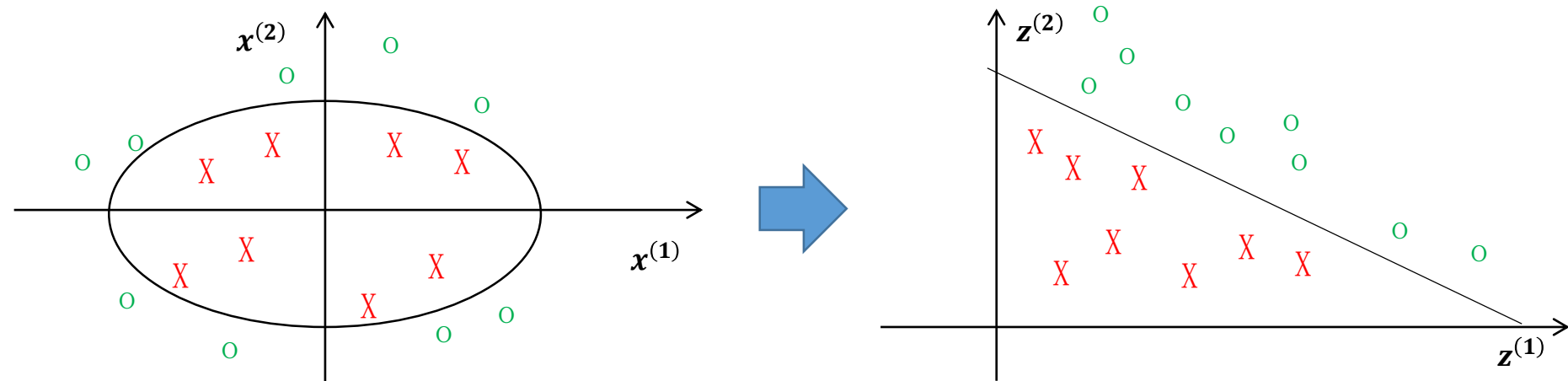
$$\begin{aligned} & \min_{w,b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } & y_i(w^T x_i + b) \geq 1 \end{aligned} \quad \Rightarrow \quad \min_{w,b} \sum_{i=1}^N [1 - y_i(w^T x_i + b)]_+ + \lambda \|w\|^2$$

$$[z]_+ = \begin{cases} z, & z > 0 \\ 0, & z \leq 0 \end{cases}$$



• 核技巧

- ❑ 线性不可分，可采取非线性映射
- ❑ 例子：线性不可分 $(x^{(1)}, x^{(2)})$ 分类平面 $w_1(x^{(1)})^2 + w_2(x^{(2)})^2 + b = 0$
- ❑ 非线性映射 $z^{(1)} = (x^{(1)})^2$, $z^{(2)} = (x^{(2)})^2$, 线性可分
 $w_1z^{(1)} + w_2z^{(2)} + b = 0$



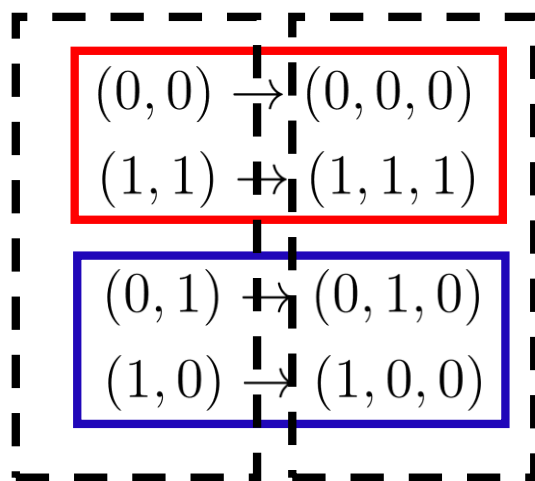
• XOR问题

- 二维样本集 $x = (x_1, x_2)$
- 第一类(0, 0)和 (1, 1), 第二类(1, 0)和 (0, 1)

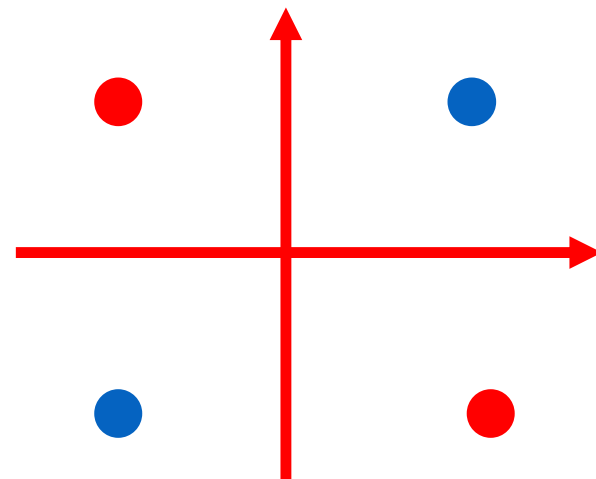
将二维数据映射到三维

- 映射函数

$$\phi(x) = (x_1, x_2, x_1x_2)$$



线性不可分 线性可分



• 定义:

- 映射 $\phi(x)$: 输入空间 \rightarrow 特征空间 (希尔伯特空间)
- 核函数 $K(x, z) = \phi(x)^T \phi(z)$

• 思想

- 显式定义核函数, 而不显式定义映射函数
- 直接计算 $K(x, z)$ 容易, $\phi(x)$ 通常无穷维, 难以定义

• 支持向量机

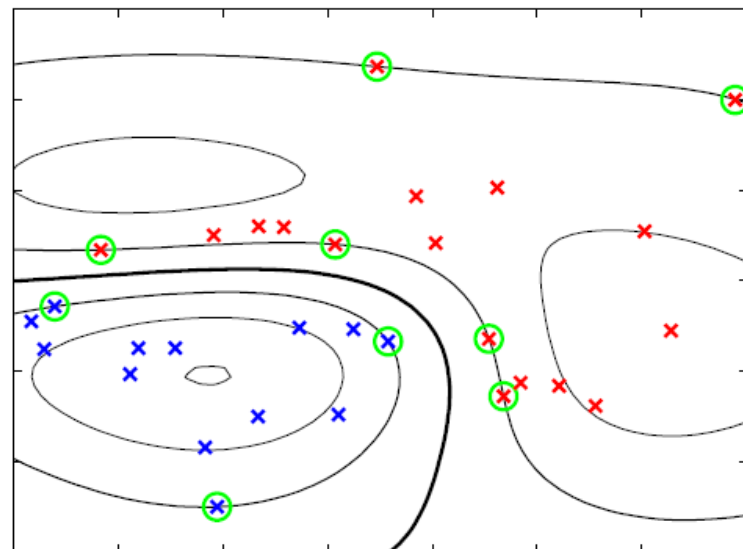
- 决策函数 $f(x) = \text{sign}(\sum_{i=1}^N \alpha_i^* y_i (x^T x_i) + b^*)$
- 内积 $x^T x_i$ 由核函数代替 $K(x, x_i)$, 实际上进行了特征映射
- 决策函数 $f(x) = \text{sign}(\sum_{i=1}^N \alpha_i^* y_i K(x, x_i) + b^*)$

• 正定核

- 核函数 $K(x, z)$ 应该满足什么性质？→ 构造或判断
- Gram矩阵：一组向量的内积构成矩阵的元素
- 正定：对任意的向量 v , $v^T K v \geq 0$
- 核函数： $K(x, z)$ 正定核函数, $\Leftrightarrow K(x, z)$ 对应的Gram矩阵半正定
$$K = [K(x_i, x_j)]_{m \times m}$$
 - 必要性： $K = [K(x_i, x_j)]_{m \times m} = [\phi(x_i)^T \phi(x_j)]_{m \times m} \Rightarrow v^T K v = [\sum_i v_i \phi(x_i)]^2 \geq 0$
 - 充分性： Gram矩阵可定义 $\phi: x \rightarrow K(., x) \Rightarrow K(x, z) = \phi(x)^T \phi(z)$
- 正定核函数判定：
 - 对称函数
 - 对任意的 $K = [K(x_i, x_j)]_{m \times m}$ 半正定

• 常用核函数

- 多项式核函数 $K(x, z) = (x^T z + 1)^P$
- 高斯核函数 $K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$
- 字符串核函数



• 非线性支持向量机

- $f(x) = \text{sign}(\sum_{i=1}^N \alpha_i^* y_i K(x, x_i) + b^*)$
- 算法

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

支持向量机

- 支持向量机工具

- LibSVM: <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>



- SLT中衡量函数集性能的指标
- 为了研究经验风险最小化函数集的学习一致收敛速度和推广性，统计学力理论定义了一些指标来衡量函数集的性能，其中最重要的就是VC维(Vapnik-Chervonenkis Dimension)

VC维定义：

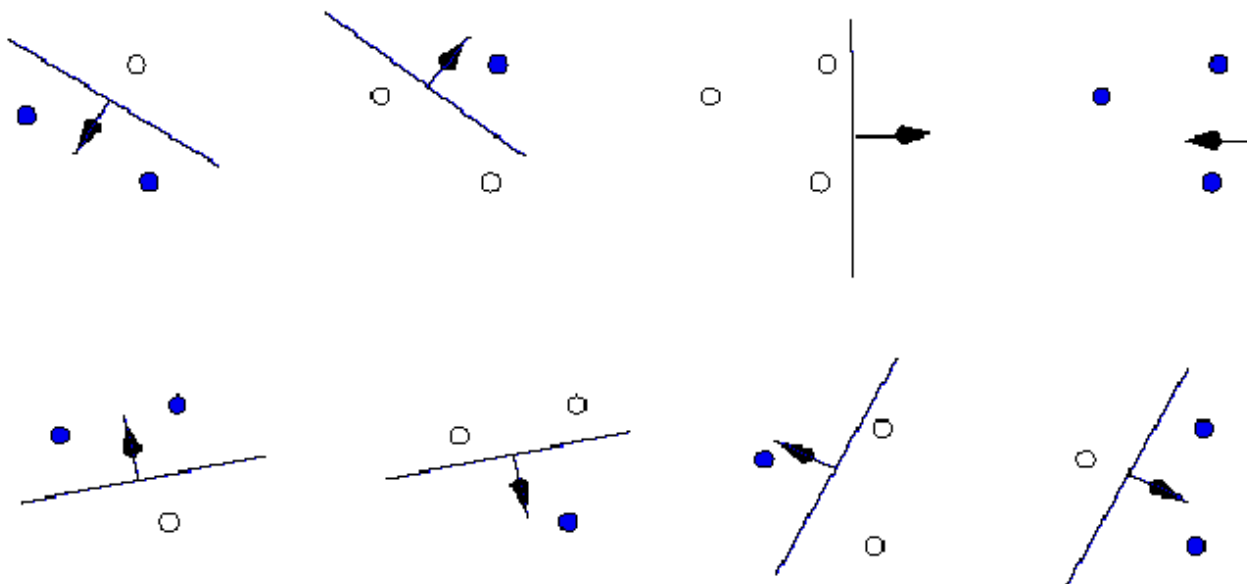
对于一个指示函数(即只有0和1两种取值的函数)集，如果存在 h 个样本能够被函数集里的函数按照所有可能的 2^h 种形式分开，则称函数集能够把 h 个样本打散，函数集的VC维就是能够打散的最大样本数目。

• SLT中衡量函数集性能的指标

□ 对于2维空间的线性分类面函数

□ 其能够以任意方式打散(分类)的最大样本数为3。

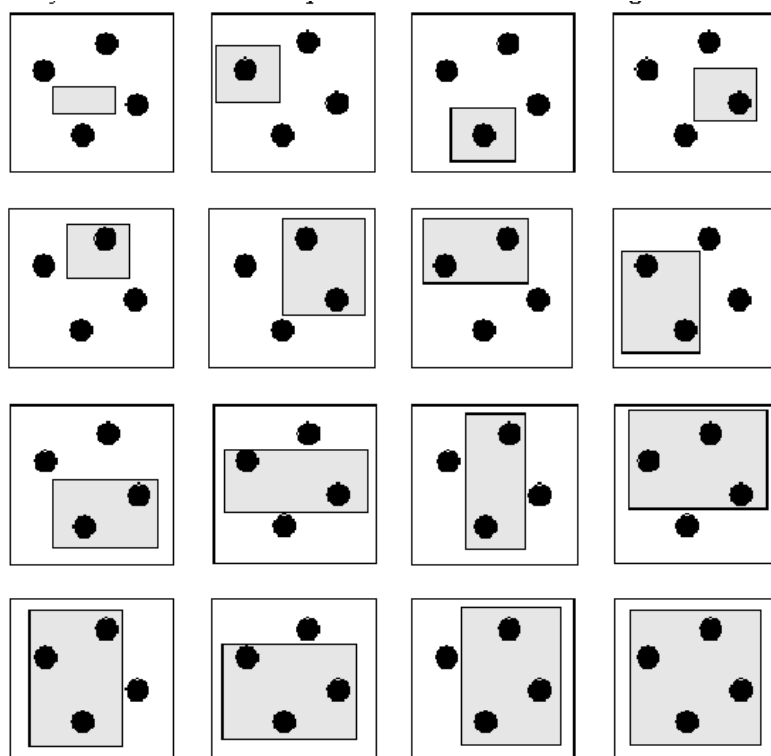
$$y = w_0 + w_1x_1 + w_2x_2$$

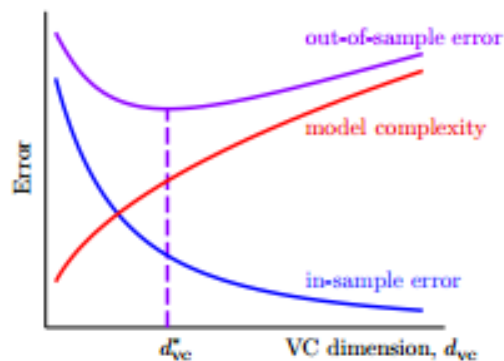


• SLT中衡量函数集性能的指标

- 一般来说,对于n维空间的线性分类面函数
- 的VC维为n+1。

$$y = w_0 + w_1x_1 + \dots + w_nx_n$$





- $d_{VC} \uparrow$: $E_{in} \downarrow$ but $\Omega \uparrow$
- $d_{VC} \downarrow$: $\Omega \downarrow$ but $E_{in} \uparrow$
- best d_{VC}^* in the middle

powerful \mathcal{H} not always good!

$$\mathbb{P}_{\mathcal{D}} \left[|E_{in}(g) - E_{out}(g)| > \epsilon \right] \leq 4(2N)^{d_{VC}} \exp \left(-\frac{1}{8} \epsilon^2 N \right)$$

theory: $N \approx 10,000 d_{VC}$; practice: $N \approx 10 d_{VC}$

given **specs** $\epsilon = 0.1$, $\delta = 0.1$, $d_{VC} = 3$, want $4(2N)^{d_{VC}} \exp \left(-\frac{1}{8} \epsilon^2 N \right) \leq \delta$

N	bound
100	2.82×10^7
1,000	9.17×10^9
10,000	1.19×10^8
100,000	1.65×10^{-38}
29,300	9.99×10^{-2}

sample complexity:
need $N \approx 10,000 d_{VC}$ in theory

• SLT中衡量函数集性能的指标

- 一般而言，VC维反映了函数集的学习能力，VC维越大则学习机器越复杂，学习内容量就越大。
- 目前没有通用的关于任意函数集VC维计算的理论，只对一些特殊的函数集知道其函数维。
- 在 n 维实数空间中线性分类器和线性实函数的VC维是 $n+1$
- 而 $f(x, a) = \sin(ax)$ 的VC维则为无穷大
- 如何用理论和试验的方法计算其VC维是当前SLT中一个待研究的问题。



软件开发环境国家重点实验室
State Key Laboratory of Software Development Environment

本节课结束