

MACHINE LEARNING

机器学习

Model Evaluation

模型评估与选择

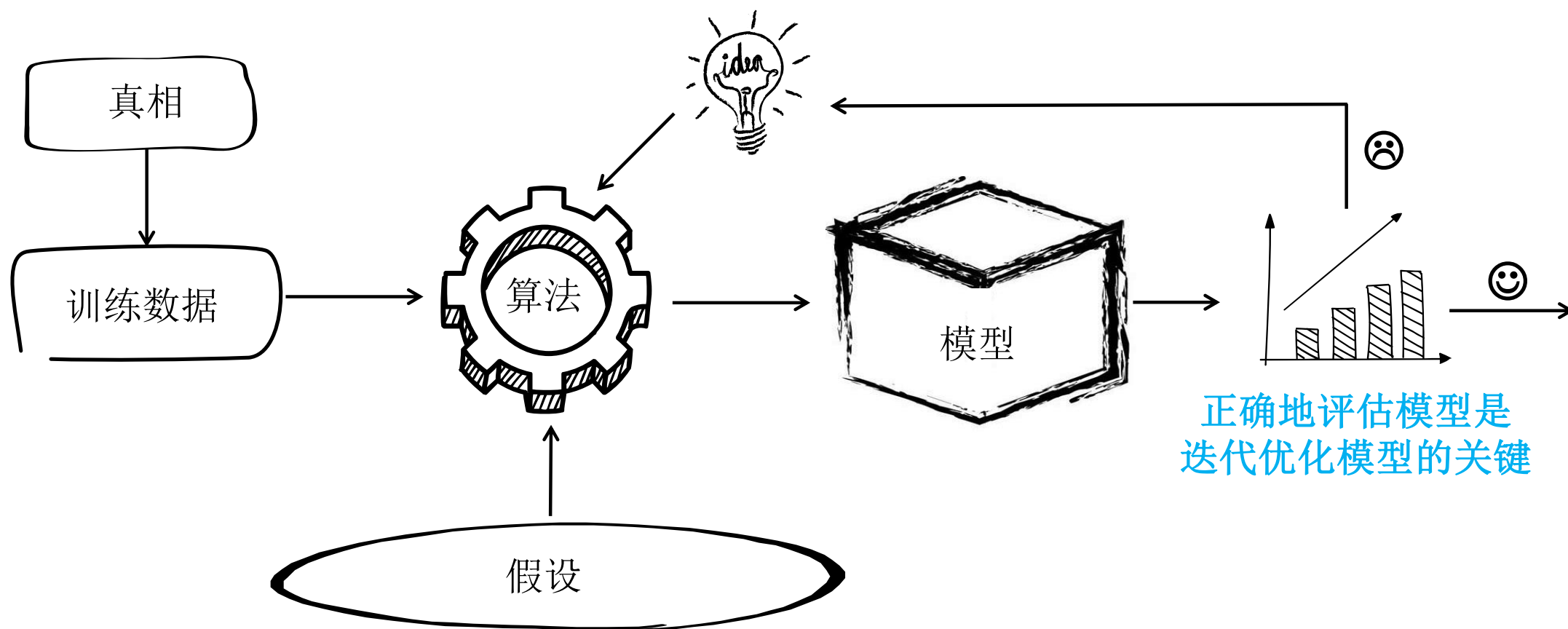


参考：
《机器学习》

Machine Learning Course
Copyright belongs to Wenting Tu.

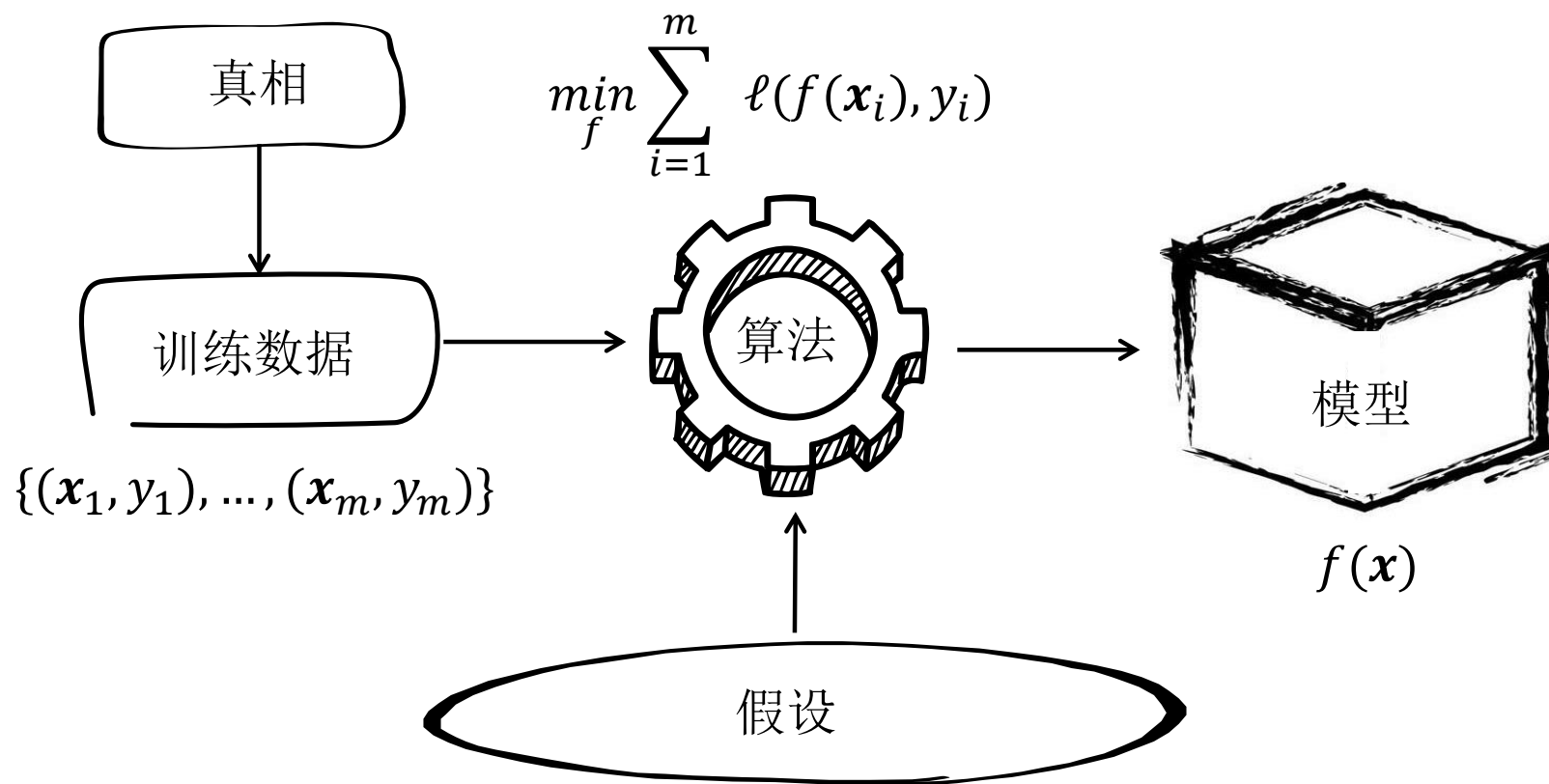
模型评估

- 评估与调优



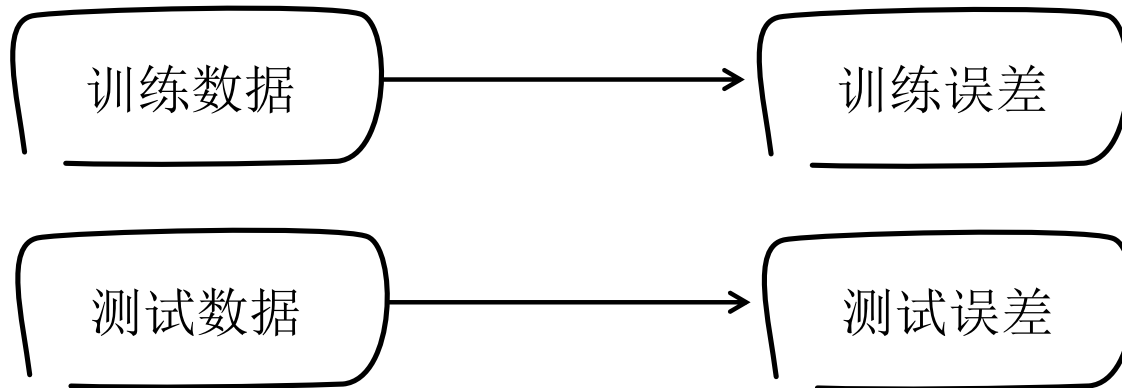
评估的对象

- 经验风险最小化



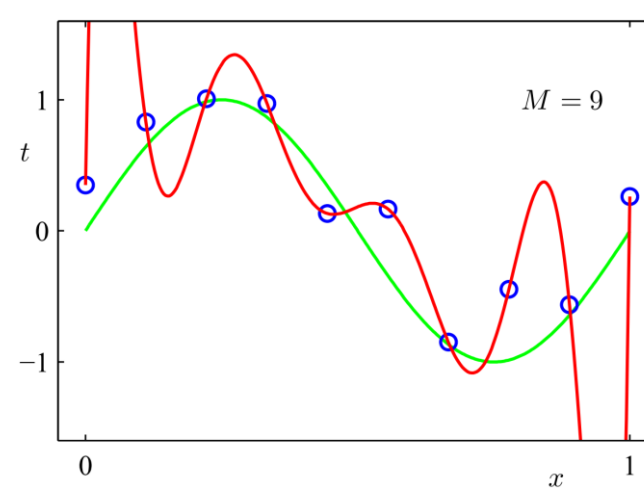
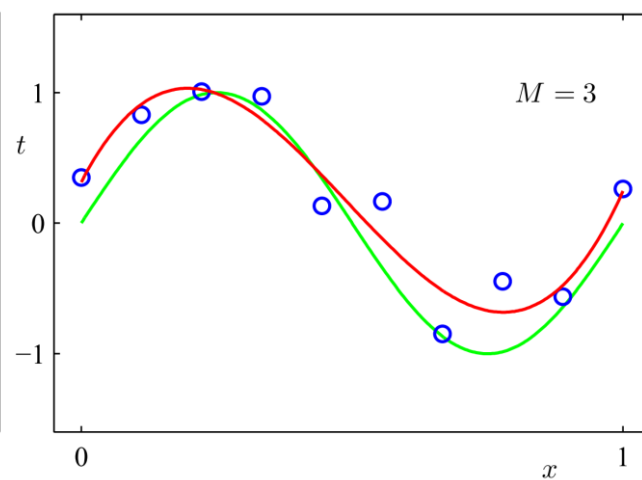
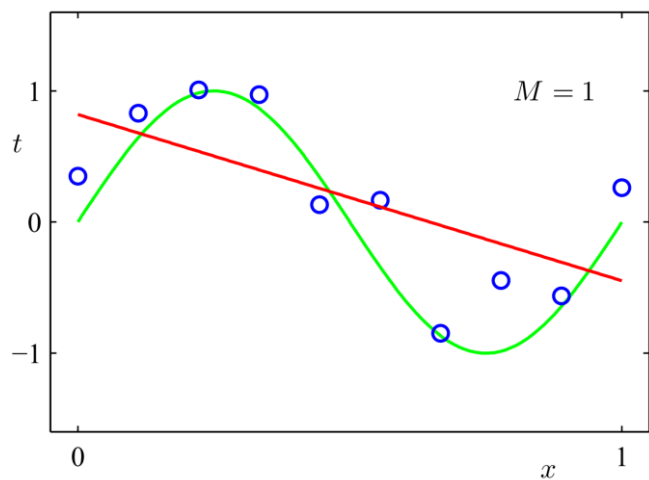
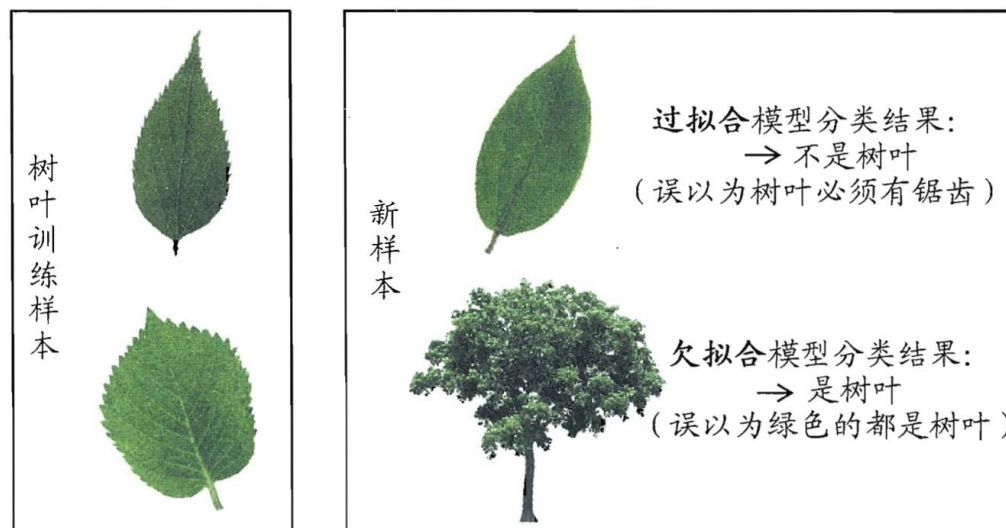
评估的对象

- 经验误差 vs 泛化误差



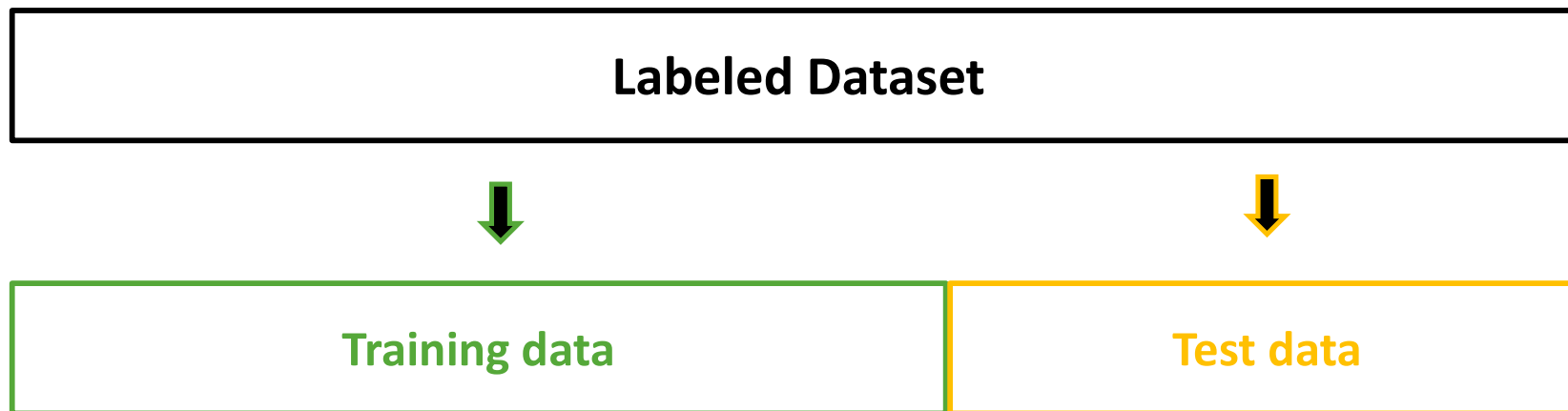
欠拟合 vs 过拟合

• 示例



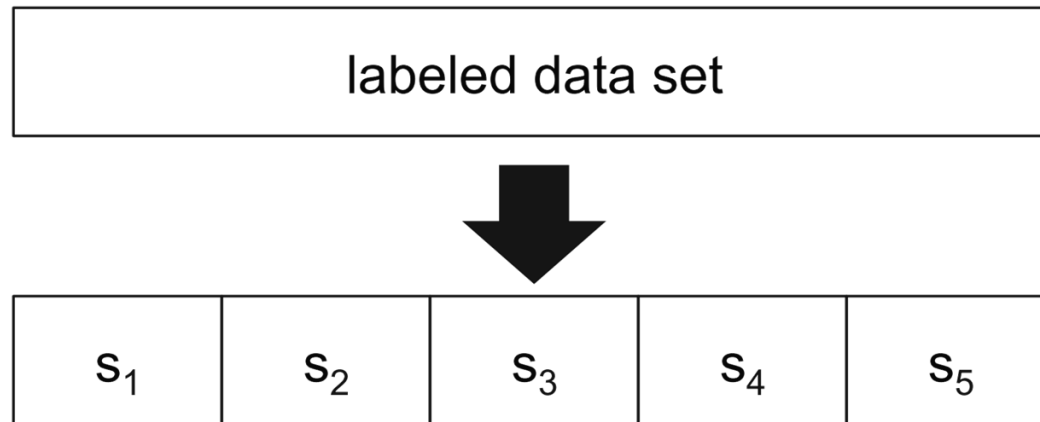
数据切分方法

- 留出法



数据切分方法

- 交叉验证法

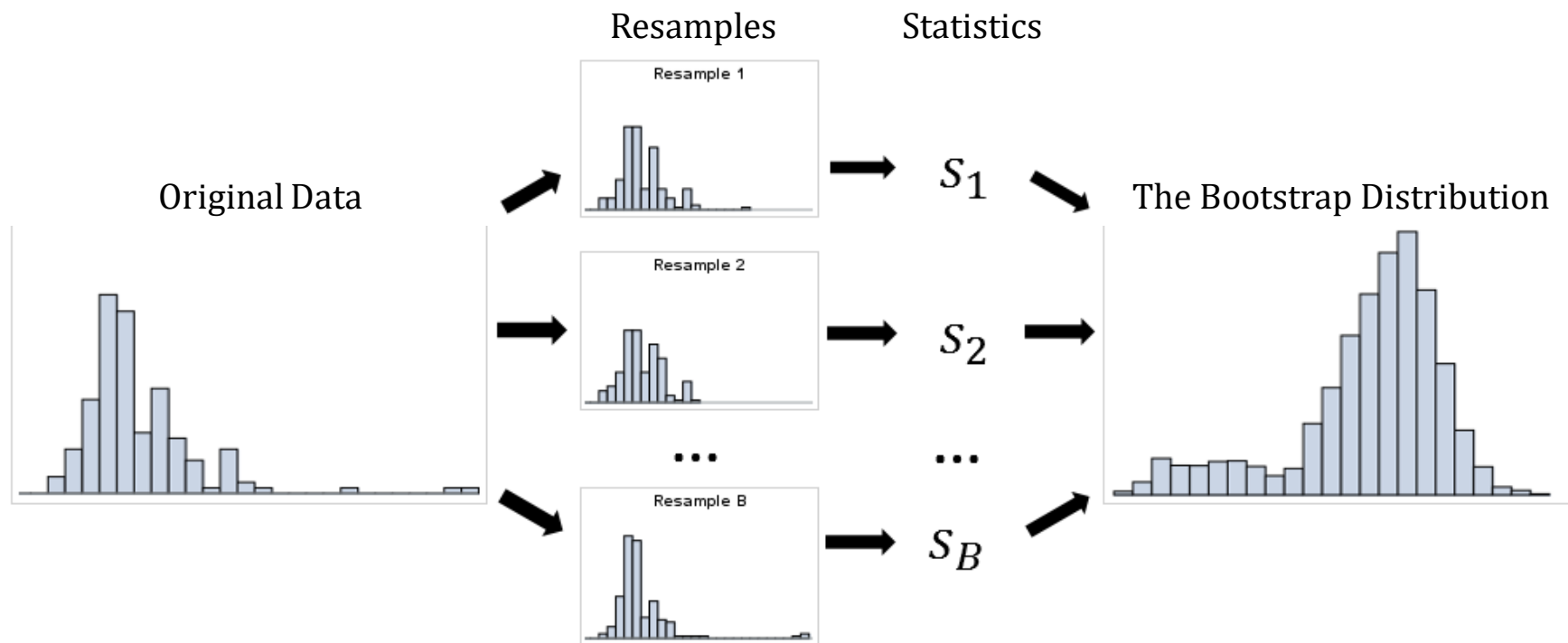


iteration	train on	test on
1	s_2 s_3 s_4 s_5	s_1
2	s_1 s_3 s_4 s_5	s_2
3	s_1 s_2 s_4 s_5	s_3
4	s_1 s_2 s_3 s_5	s_4
5	s_1 s_2 s_3 s_4	s_5

数据切分方法

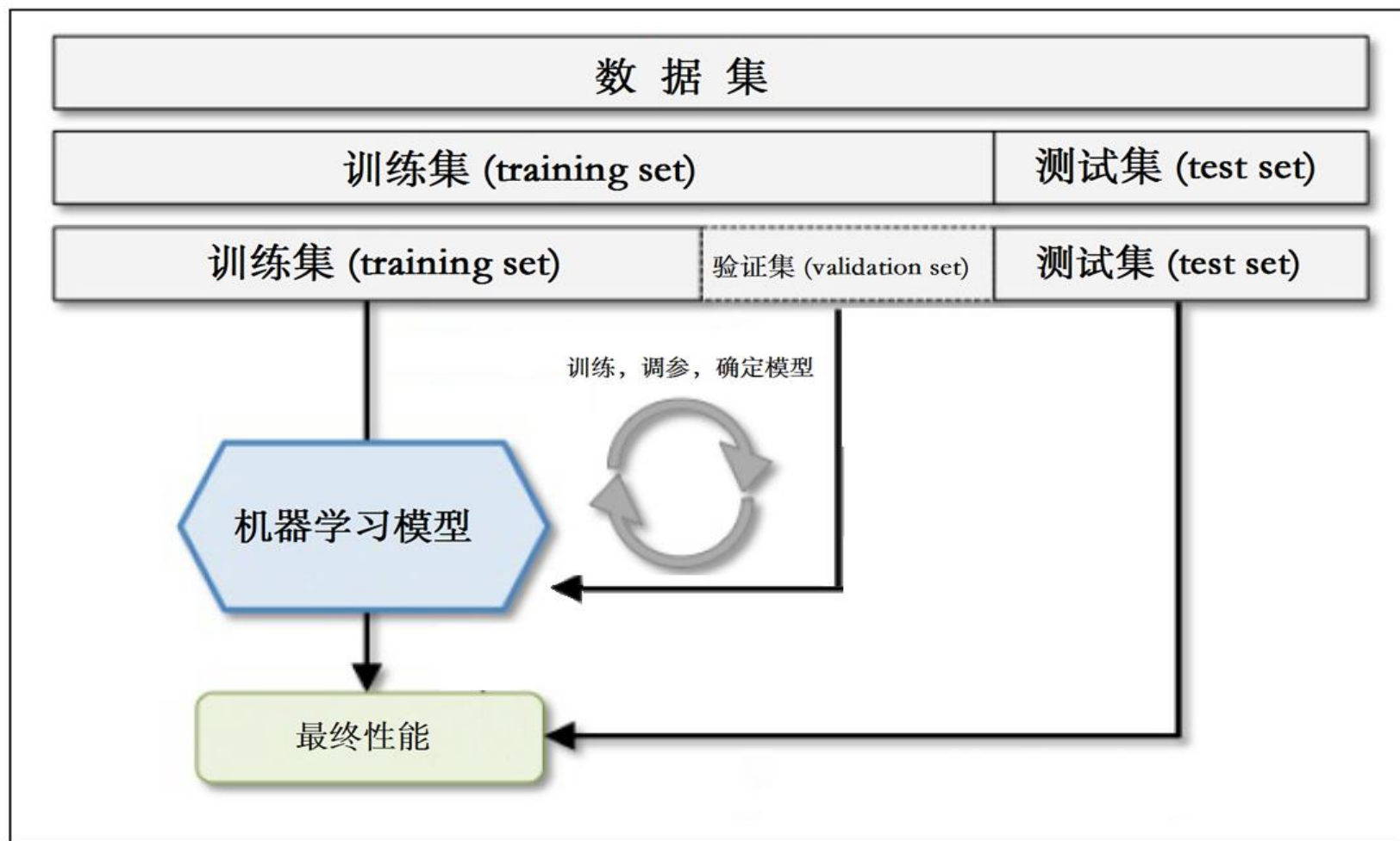
- 自助法

$$\lim_{m \rightarrow \infty} (1 - \frac{1}{m})^m \mapsto \frac{1}{e} \approx 0.368$$



调参

- 训练集 + 验证集 + 测试集



性能度量

- 回归模型的性能度量

$$E(f; \mathcal{D}) = \int_{\mathbf{x} \sim \mathcal{D}} (f(\mathbf{x}) - y)^2 p(\mathbf{x}) d\mathbf{x}$$

$$E(f; D) = \frac{1}{m} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2$$

性能度量

- 分类模型的性能度量

$$\begin{array}{ll} E(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(\mathbf{x}_i) \neq y_i) & \xleftrightarrow[p(\cdot)]{D} E(f; \mathcal{D}) = \int_{\mathbf{x} \sim \mathcal{D}} \mathbb{I}(f(\mathbf{x}) \neq y) p(\mathbf{x}) d\mathbf{x} \\ \text{acc}(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(\mathbf{x}_i) = y_i) & \text{acc}(f; \mathcal{D}) = \int_{\mathbf{x} \sim \mathcal{D}} \mathbb{I}(f(\mathbf{x}) = y) p(\mathbf{x}) d\mathbf{x} \\ & = 1 - E(f; \mathcal{D}) \\ & = 1 - E(f; D) \end{array}$$

性能度量

- 分类模型的性能度量
 - 混淆矩阵

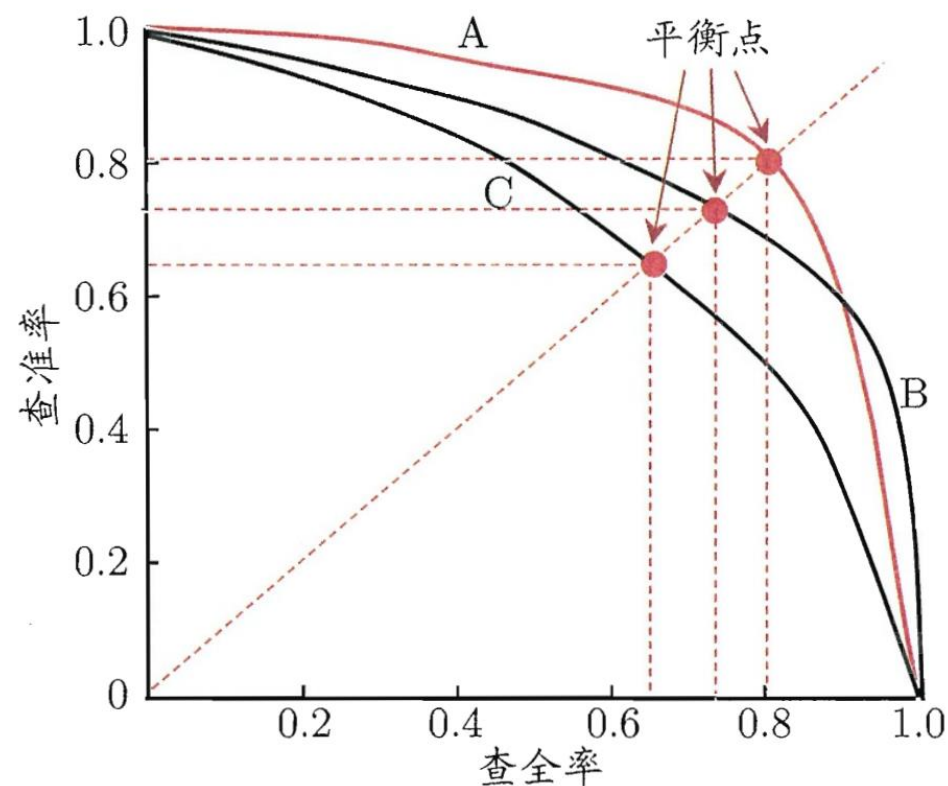
真实情况	预测结果	
	正例	反例
正例	TP (真正例)	FN (假反例)
反例	FP (假正例)	TN (真反例)

- 查准率、查全率与F1、PR曲线

$$P = \frac{TP}{TP + FP} \quad R = \frac{TP}{TP + FN}$$

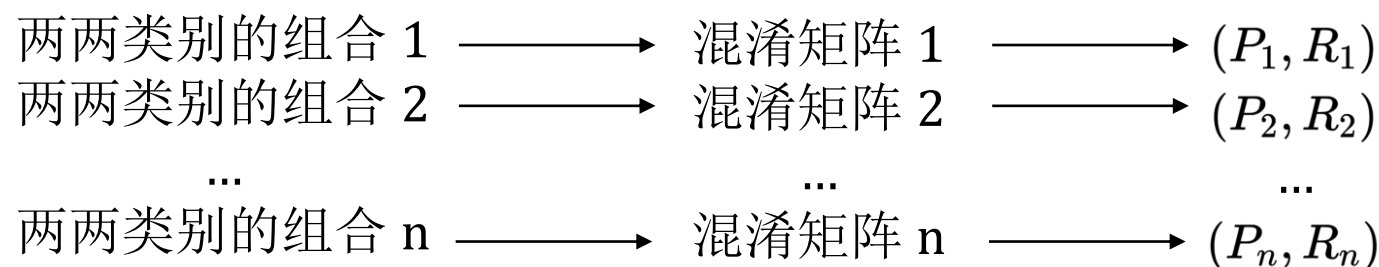
$$F1 = \frac{2 \times P \times R}{P + R} = \frac{2 \times TP}{\text{样例总数} + TP - TN}$$

$$F_{\beta} = \frac{(1 + \beta^2) \times P \times R}{(\beta^2 \times P) + R}$$



性能度量

- 分类模型的性能度量
 - 多分类场景



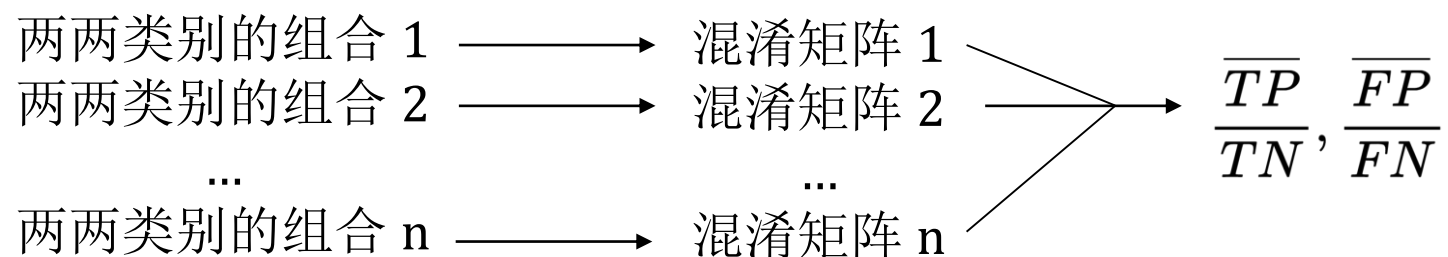
$$\text{macro} - P = \frac{1}{n} \sum_{i=1}^n P_i$$

$$\text{macro} - R = \frac{1}{n} \sum_{i=1}^n R_i$$

$$\text{macro} - F1 = \frac{2 \times \text{macro} - P \times \text{macro} - R}{\text{macro} - P + \text{macro} - R}$$

性能度量

- 分类模型的性能度量
 - 多分类场景



$$\text{micro-}P = \frac{\overline{TP}}{\overline{TP} + \overline{FP}}$$

$$\text{micro-}R = \frac{\overline{TP}}{\overline{TP} + \overline{FN}}$$

$$\text{micro-}F1 = \frac{2 \times \text{micro-}P \times \text{micro-}R}{\text{micro-}P + \text{micro-}R}$$

性能度量

- 分类模型的性能度量
 - 混淆矩阵

真实情况	预测结果	
	正例	反例
正例	TP (真正例)	FN (假反例)
反例	FP (假正例)	TN (真反例)

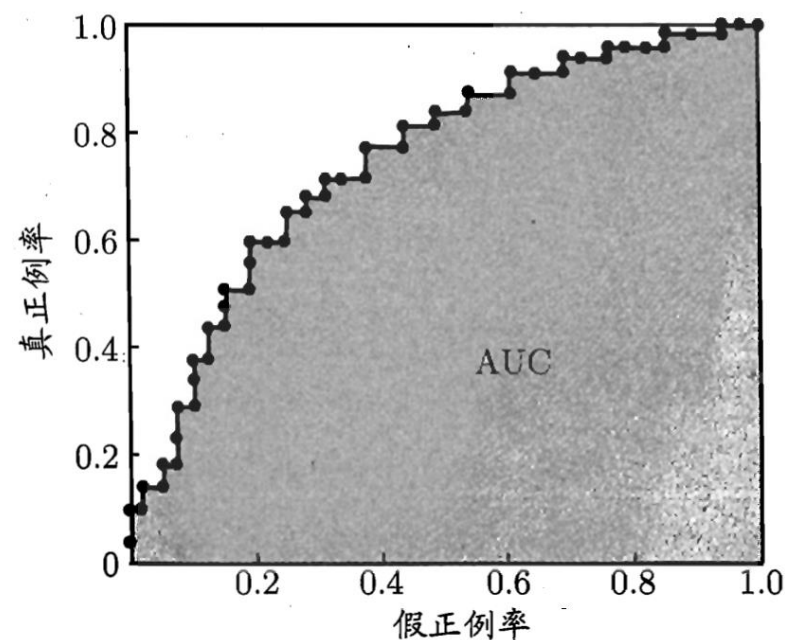
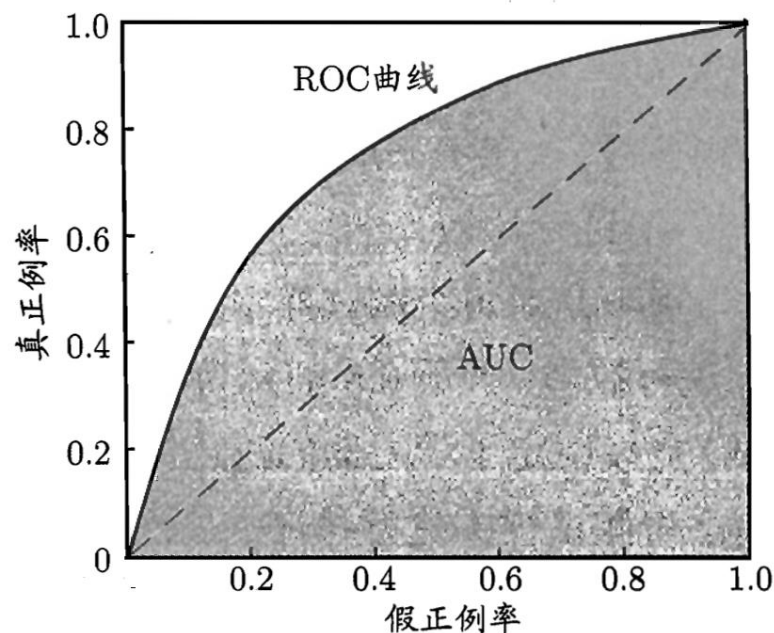
- “真正例率” TPR 与“假正例率” FPR

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{TN + FP}$$

性能度量

- 分类模型的性能度量
 - ROC 与 AUC

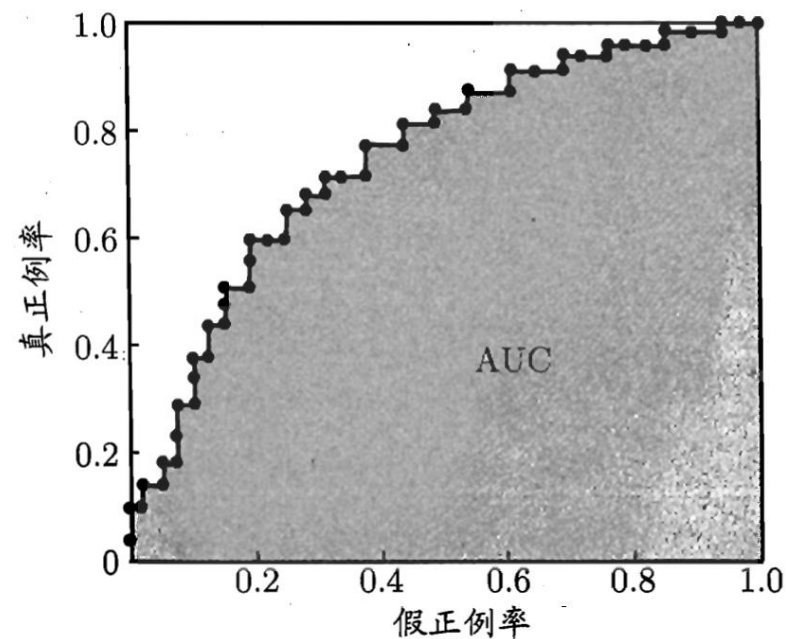
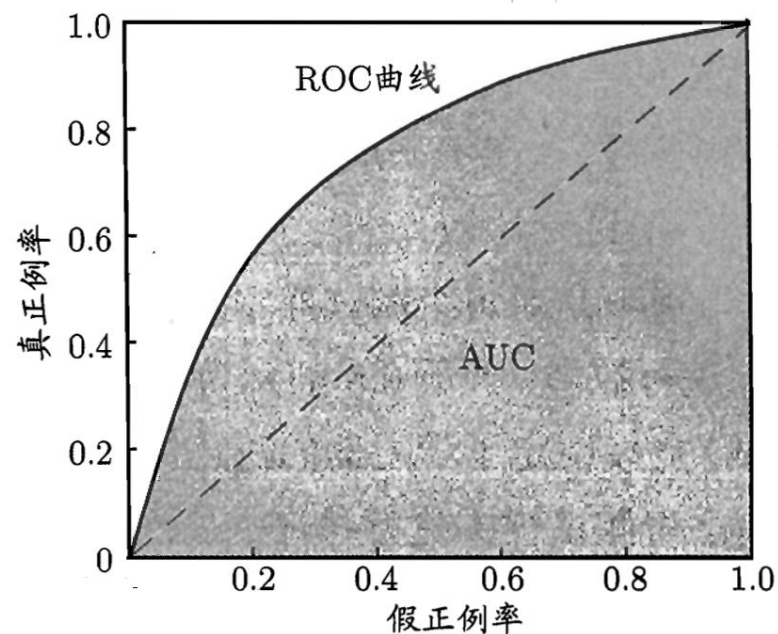


AUC = “ROC 曲线之下的面积” = 1 - “ROC 曲线之上的面积” = $1 - \ell_{\text{rank}}$

$$\ell_{\text{rank}} = \frac{1}{m^+ m^-} \sum_{\mathbf{x}^+ \in D^+} \sum_{\mathbf{x}^- \in D^-} \left(\mathbb{I}(f(\mathbf{x}^+) < f(\mathbf{x}^-)) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}^+) = f(\mathbf{x}^-)) \right)$$

性能度量

- 分类模型的性能度量
 - ROC 与 AUC



思考：多分类场景下的ROC与AUC

性能度量

- 分类模型的性能度量
 - 代价敏感场景

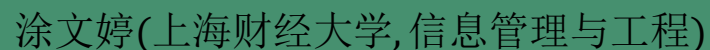
真实类别	预测类别	
	第 0 类	第 1 类
第 0 类	0	$cost_{10}$
第 1 类	$cost_{01}$	0

$$E(f; D; cost) = \frac{1}{m} \left(\sum_{x_i \in D^+} \mathbb{I}(f(x_i) \neq y_i) \times cost_{01} + \sum_{x_i \in D^-} \mathbb{I}(f(x_i) \neq y_i) \times cost_{10} \right)$$

- Q: 泛化错误率 $\epsilon_i^A \leq \epsilon_i^B$

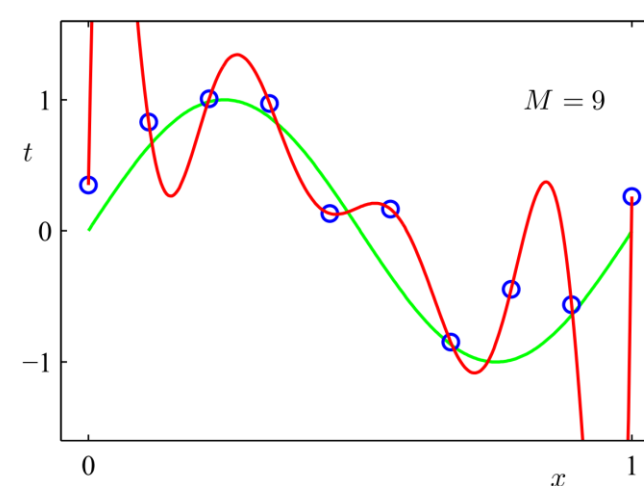
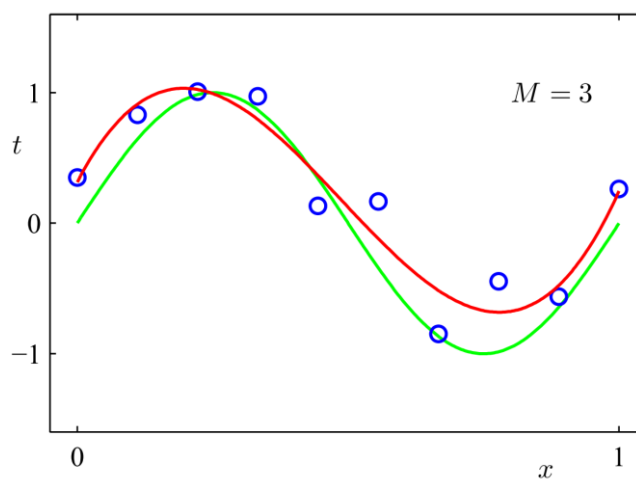
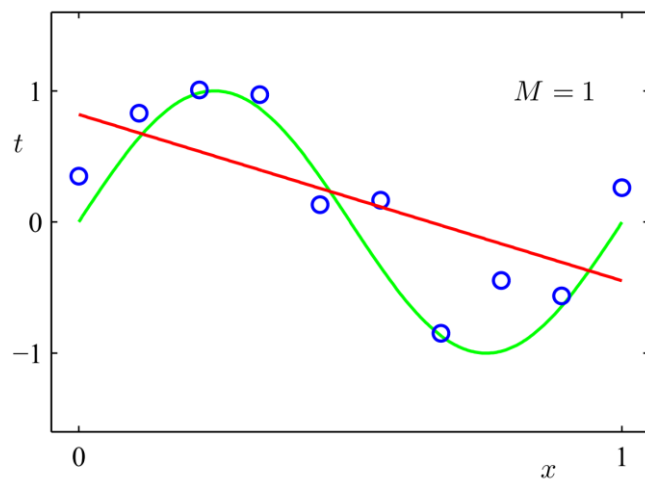
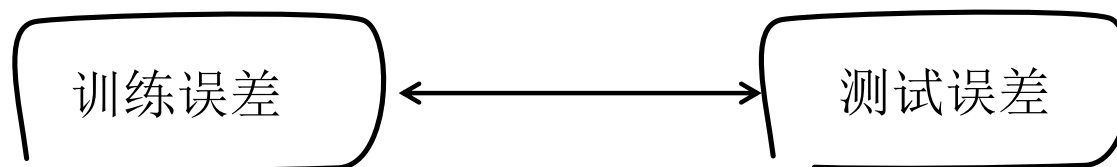
然后利用 t 检验: $\tau_t = \frac{\sqrt{k}\mu_{\leftarrow}}{\sigma_{\leftarrow}}$ 差值的均值 / 差值的方差

判定上述变量值是否小于临界值 $t_{\alpha/2, k-1}$ ，即尾部累积分布为 $\alpha/2$ 的临界值



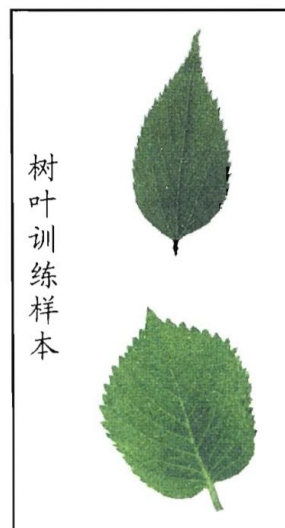
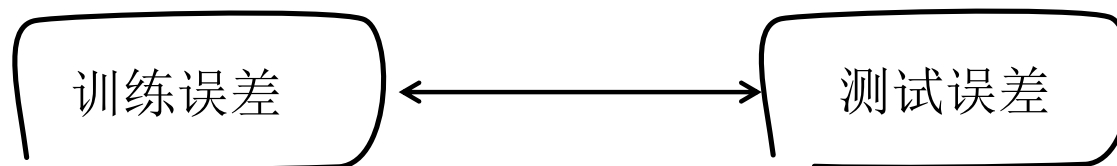
欠拟合 vs 过拟合

- 示例



欠拟合 vs 过拟合

- 示例



偏差与方差

• 偏差方差分解 (bias-variance decomposition)

$$\begin{aligned}
 E(f; D) &= \mathbb{E}_D[(f(\mathbf{x}; D) - y_D)^2] \\
 &= \mathbb{E}_D[(f(\mathbf{x}; D) - \bar{f}(\mathbf{x}))^2] + (\bar{f}(\mathbf{x}) - y)^2 + \mathbb{E}_D[(y_D - y)^2] \\
 &= \text{bias}^2(\mathbf{x}) + \text{var}(\mathbf{x}) + \varepsilon^2
 \end{aligned}$$

y_D : \mathbf{x} 在数据集中的标记

$f(\mathbf{x}; D)$: 训练集 D 上学得模型 f 在 \mathbf{x} 上的预测输出

$\bar{f}(\mathbf{x}) = \mathbb{E}_D[f(\mathbf{x}; D)]$:

以回归任务为例, 学习算法的期望预测

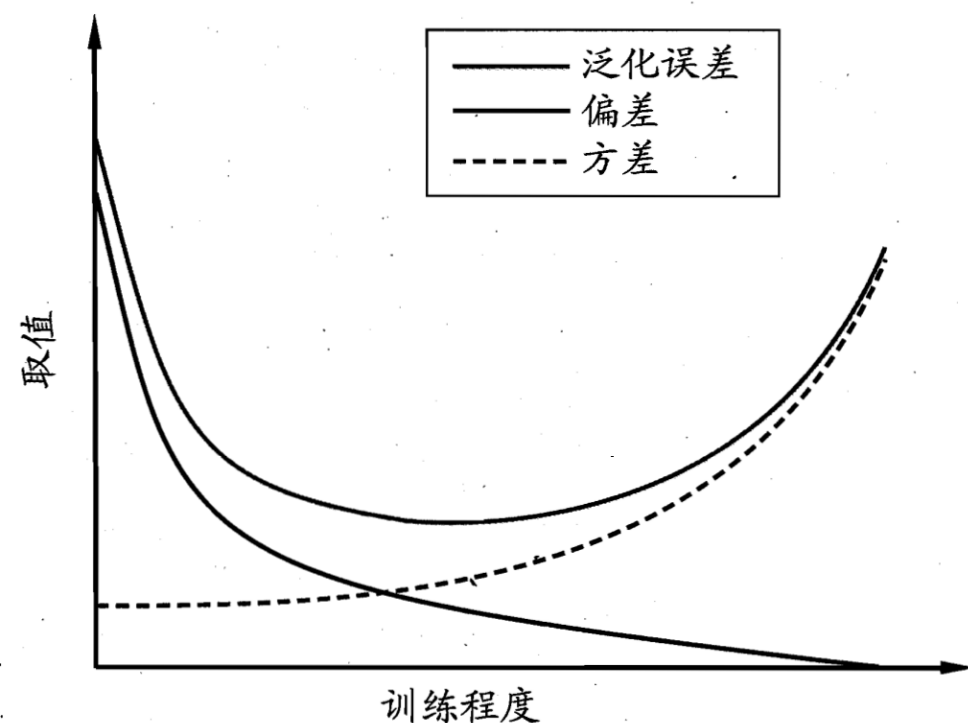
$\text{var}(\mathbf{x}) = \mathbb{E}_D[(f(\mathbf{x}; D) - \bar{f}(\mathbf{x}))^2]$:

使用样本数相同的不同训练集产生的方差为

$\varepsilon^2 = \mathbb{E}_D[(y_D - y)^2]$: 噪声

$\text{bias}^2(\mathbf{x}) = (\bar{f}(\mathbf{x}) - y)^2$:

期望输出与真实标记的差别称为偏差



正则化技术

- 正则化 *regularization* 技术

$$\min_f \Omega(f) + C \sum_{i=1}^m \ell(f(\mathbf{x}_i), y_i)$$

结构风险用于描述模型的某些性质

经验风险用于描述模型与训练数据的契合程度;

- 岭回归

$$\min_{\mathbf{w}, b} \|\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}\|_2^2 + \lambda \|\mathbf{w}\|_2^2, \quad \|\mathbf{w}\|_2^2 = \sum_{j=1}^d w_j^2$$

$$\min_{\mathbf{w}, b} (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}) + \lambda \|\mathbf{w}\|_2^2$$

$$\hat{\mathbf{w}}_{\text{RR}} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

岭回归与最大后验估计

- 岭回归的概率模型

> 最大似然估计

$$q(\mathbf{x}; \boldsymbol{\theta}) \rightarrow p(\mathbf{x})$$

$$\mathcal{D} = \{\mathbf{x}\}_{i=1}^n$$

Likelihood $p(\mathcal{D} | \boldsymbol{\theta})$

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\mathcal{D} | \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{i=1}^n q(\mathbf{x}_i; \boldsymbol{\theta})$$

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \log L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \left[\sum_{i=1}^n \log q(\mathbf{x}_i; \boldsymbol{\theta}) \right]$$

$$p(\mathbf{x}) = q(\mathbf{x}; \hat{\boldsymbol{\theta}}_{\text{MLE}})$$

岭回归与最大后验估计

- 岭回归的概率模型

- > 最大后验估计

$$q(\mathbf{x}; \boldsymbol{\theta}) \rightarrow p(\mathbf{x})$$

$$\mathcal{D} = \{\mathbf{x}\}_{i=1}^n$$

Likelihood $p(\mathcal{D} \mid \boldsymbol{\theta})$

Prior $p(\boldsymbol{\theta})$

Posterior $p(\boldsymbol{\theta} \mid \mathcal{D})$

$$\boldsymbol{\theta}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta} \mid \mathcal{D})$$

$$\boldsymbol{\theta}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left(\sum_{i=1}^n \log q(\mathbf{x}_i \mid \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \right)$$

$$p(\mathbf{x}) = q(\mathbf{x}; \hat{\boldsymbol{\theta}}_{\text{MAP}})$$

岭回归与最大后验估计

- 岭回归的概率模型

$$p(t | \mathbf{x})?$$

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

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}, \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

$$\mathbf{w}_{MAP}^* = \arg \max_{\mathbf{w}} \ln p(\mathbf{w} | \mathbf{t})$$

$$= \arg \max_{\mathbf{w}} \left(-\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const} \right)$$

$$\mathbf{w}_{MAP}^* = \mathbf{w}_{RR}^*$$

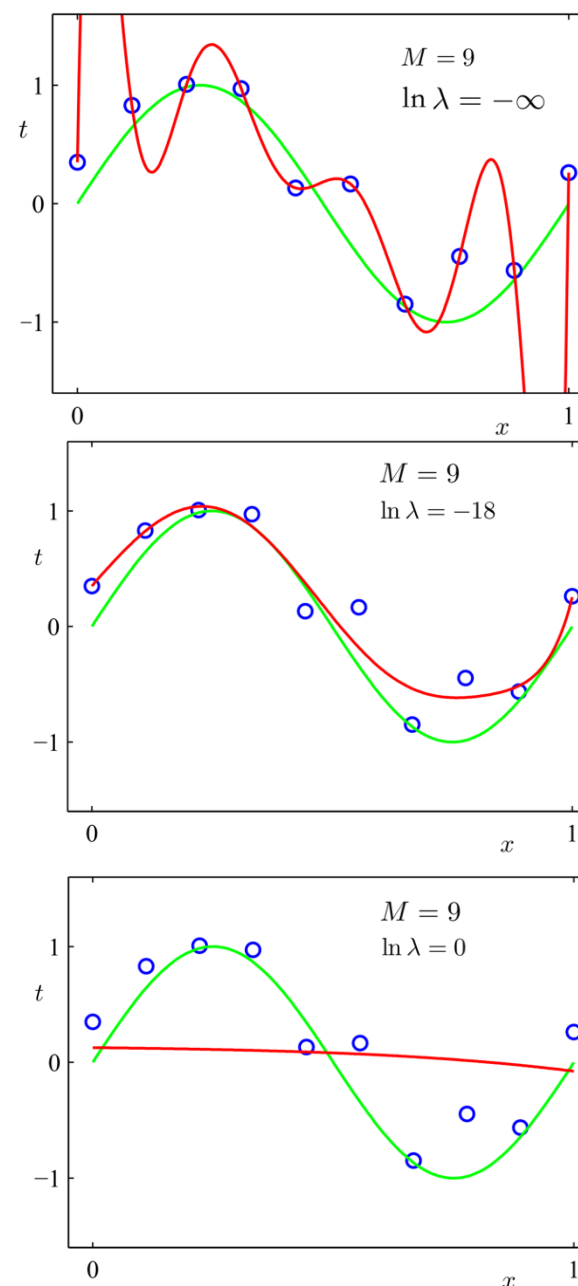
岭回归与LASSO

• 岭回归

$$\min_{\mathbf{w}, b} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$M = 9$

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01

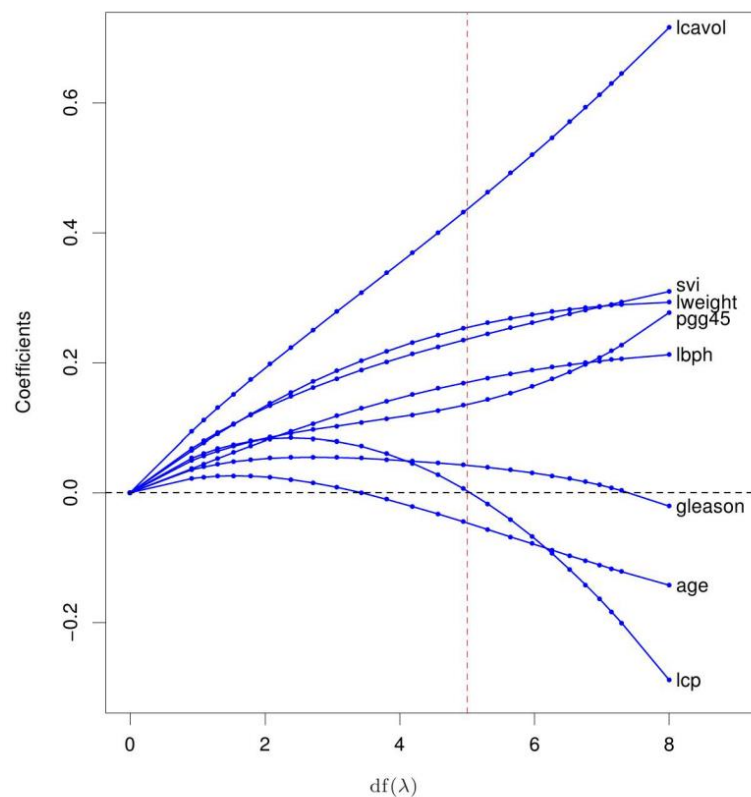


岭回归与LASSO

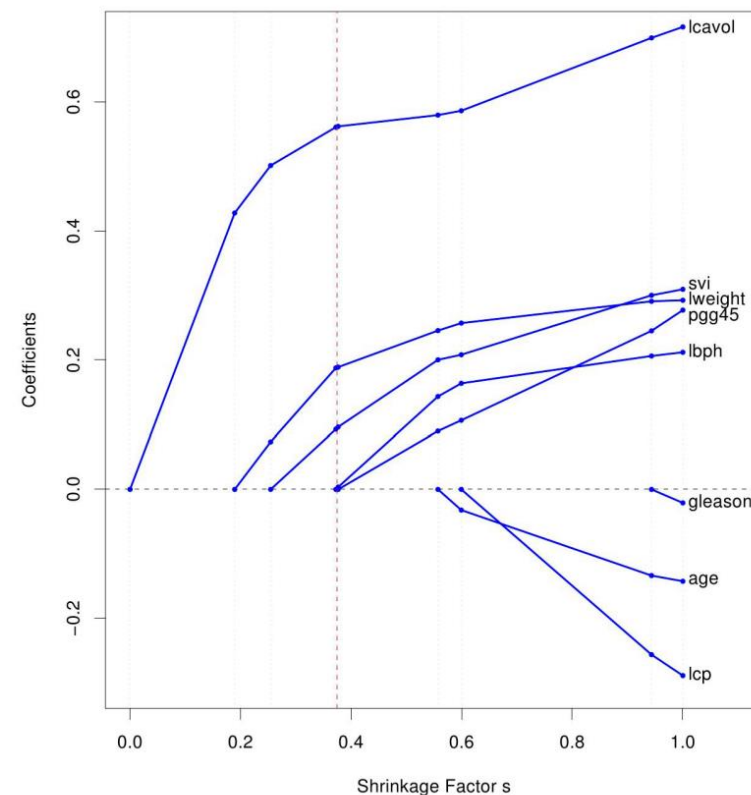
• LASSO

$$\min_{\mathbf{w}, b} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

$$\|\mathbf{w}\|_1 = \sum_{j=1}^d |w_j|$$



(a) $\|w\|_2$ penalty



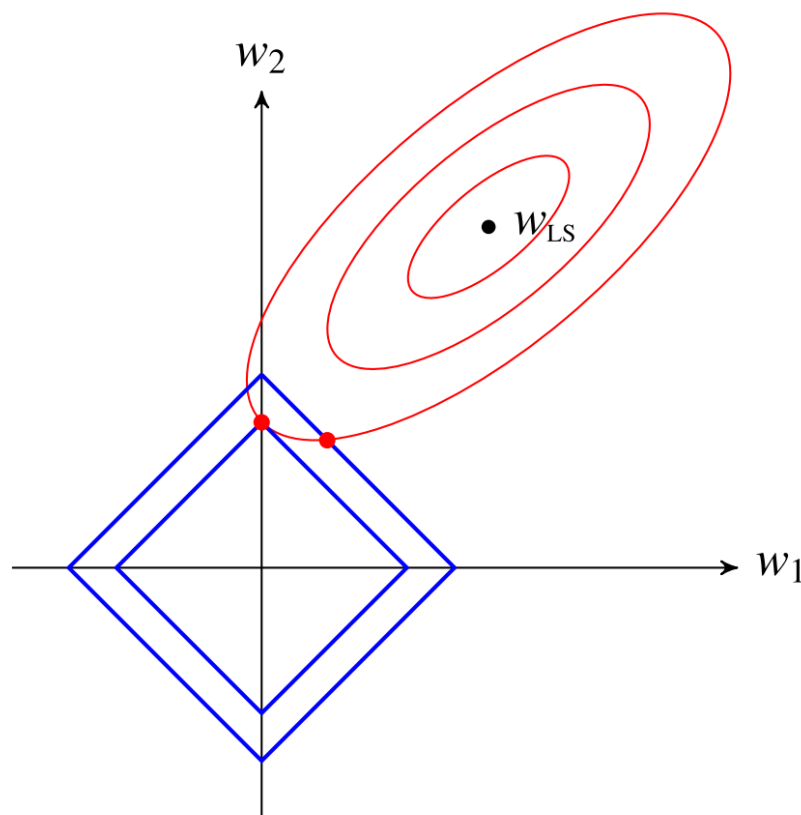
(b) $\|w\|_1$ penalty

岭回归与LASSO

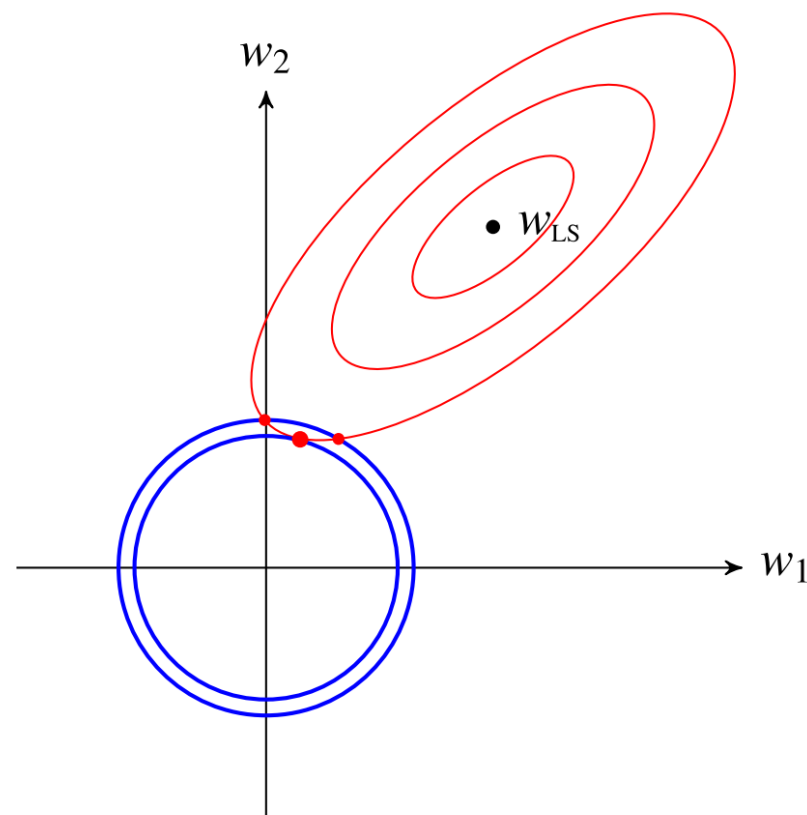
- LASSO

$$\min_{\mathbf{w}, b} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

$$\|\mathbf{w}\|_1 = \sum_{j=1}^d |w_j|$$



(a) $\|w\|_2$ penalty



(b) $\|w\|_1$ penalty

MACHINE LEARNING

实践

Practice



参考:

<https://scikit-learn.org/>

程序示例

- 通过sklearn在线文档可以获得编程练习

[3.1. Cross-validation: evaluating estimator performance](#)

[3.3. Metrics and scoring: quantifying the quality of predictions](#)