

MACHINE LEARNING

机器学习

Linear Models

线性模型

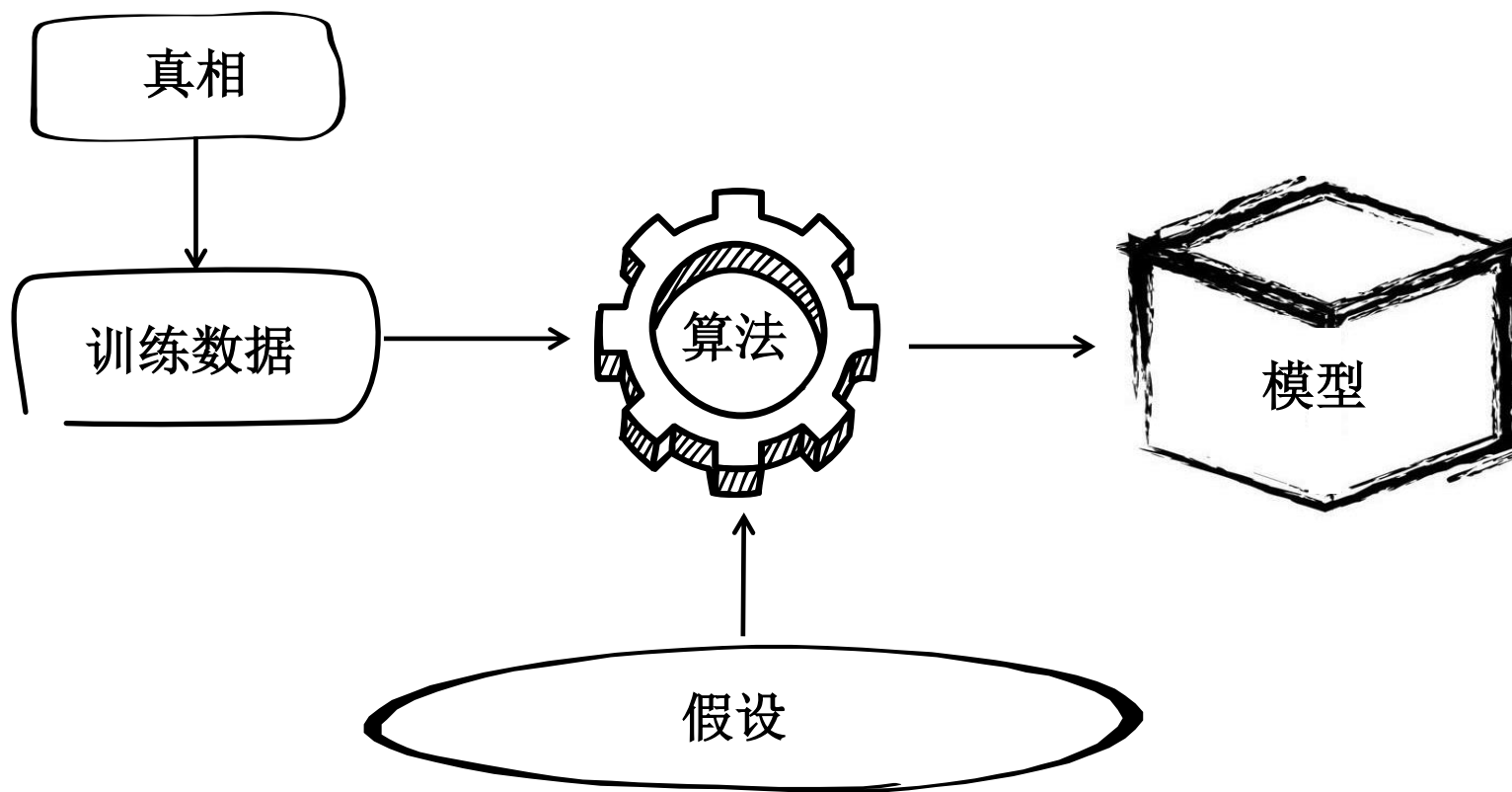


参考：
《机器学习》

Machine Learning Course
Copyright belongs to Wenting Tu.

定义

- 线性假设空间



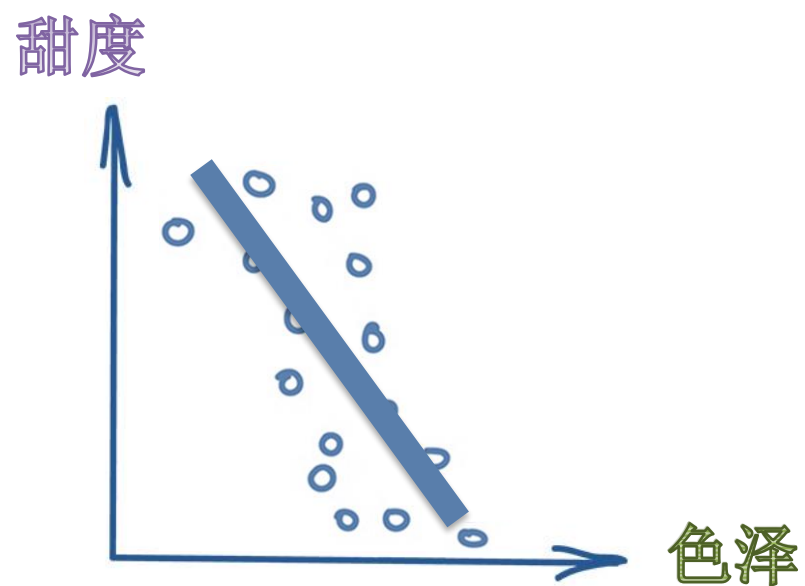
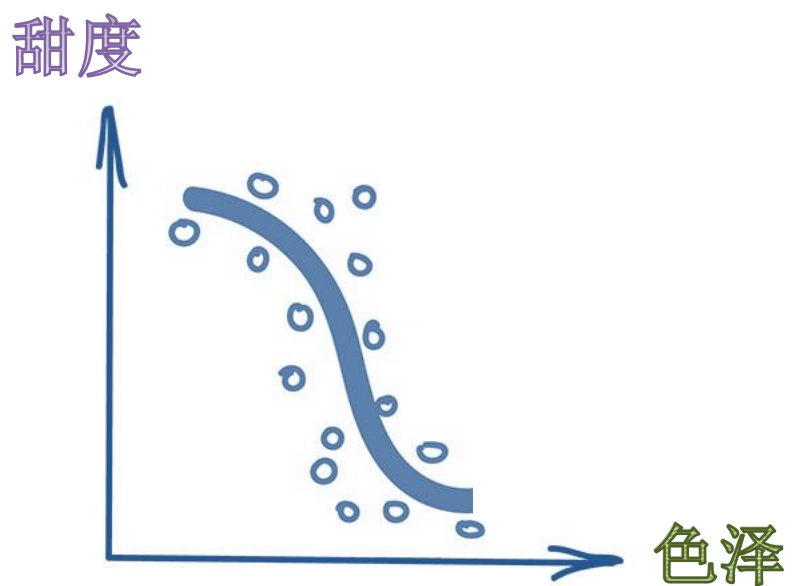
$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \cdots + w_dx_d + b$$

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

$$f_{\text{好瓜}}(\mathbf{x}) = 0.2 \cdot x_{\text{色泽}} + 0.5 \cdot x_{\text{根蒂}} + 0.3 \cdot x_{\text{敲声}} + 1$$

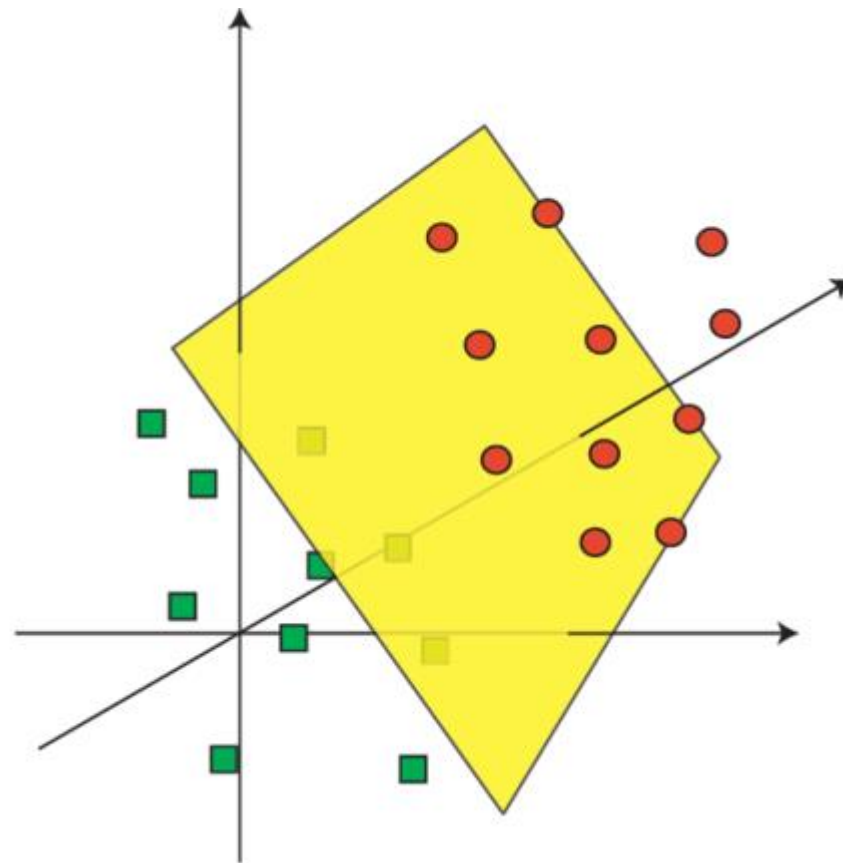
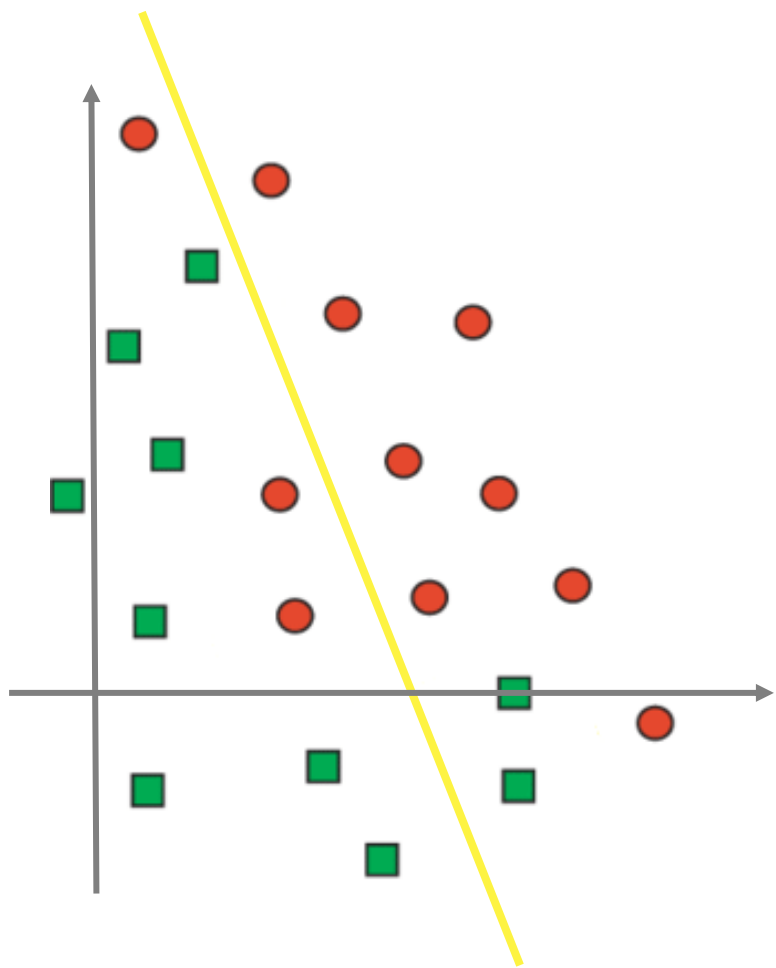
定义

- 线性拟合函数



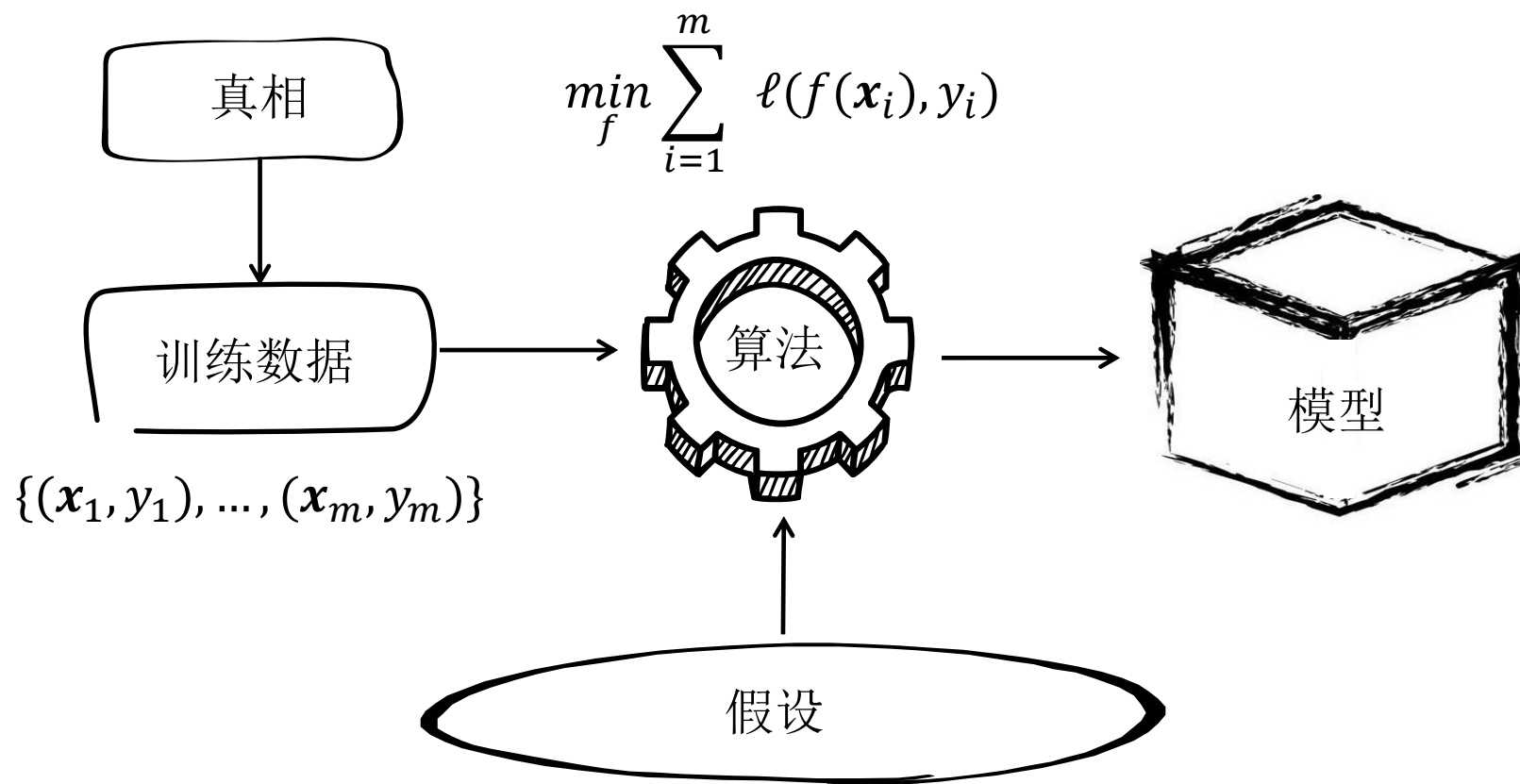
定义

- 线性分类超平面



线性回归

- 最小经验损失



$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$$
$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

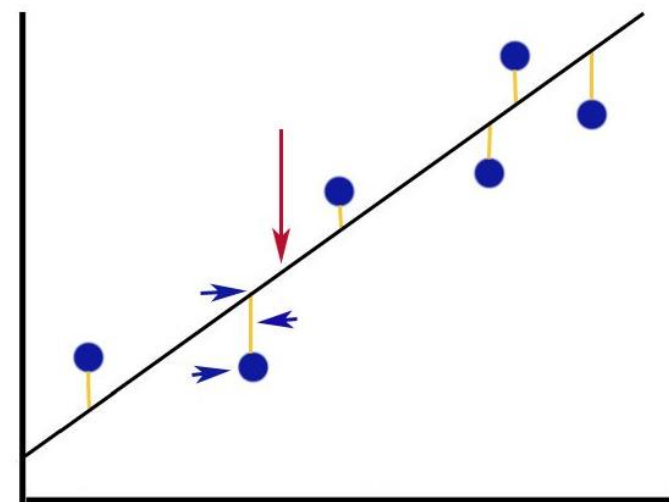
线性回归

- 损失函数

$$\ell(f(\mathbf{x}_i), y_i) = (f(x_i) - y_i)^2$$

- 经验误差

$$\begin{aligned} \min_f \sum_{i=1}^m \ell(f(\mathbf{x}_i), y_i) &= \arg \min_{(w,b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \\ &= \arg \min_{(w,b)} \sum_{i=1}^m (y_i - wx_i - b)^2 \end{aligned}$$



基于均方误差最小化来进行模型求解的方法称为“**最小二乘法**”(least square method).

在线性回归中，最小二乘法就是试图找到一条直线，使所有样本到直线上的欧氏距离之和最小。

线性回归

- 解析解

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix}$$

$$\mathbf{y} = (y_1; y_2; \dots; y_m)$$

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

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \rightarrow \hat{\mathbf{w}}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

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

线性回归

• 解析解

$\hat{\mathbf{w}}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, 此解析解假设了 $(\mathbf{X}^T \mathbf{X})^{-1}$ 存在, 这需要 $\mathbf{X}^T \mathbf{X}$ 是一个满秩矩阵 (full rank matrix)

若 $m < d+1$, 则 $(\mathbf{X}^T \mathbf{X})^{-1}$ 便不存在, 此时无法用解析解求解最小二乘法

• 梯度下降法

考虑无约束优化问题 $\min_{\mathbf{x}} f(\mathbf{x})$

若能构造一个序列 $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots$

满足 $f(\mathbf{x}^{t+1}) < f(\mathbf{x}^t), t=0,1,2, \dots$

则不断执行该过程即可收敛到局部极小点

根据泰勒展式有 $f(\mathbf{x} + \Delta \mathbf{x}) \simeq f(\mathbf{x}) + \Delta \mathbf{x}^T \nabla f(\mathbf{x})$

于是, 欲满足 $f(\mathbf{x} + \Delta \mathbf{x}) < f(\mathbf{x})$

可选择 $\Delta \mathbf{x} = -\eta \nabla f(\mathbf{x})$

其中步长 η 是一个小常数. 这就是梯度下降法

• 随机梯度下降法求解最小二乘法

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla \mathcal{L}_i$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta (y_i - \mathbf{w}^{(t)T} \mathbf{x}_i) \mathbf{x}_i$$

线性回归

- 对数线性回归

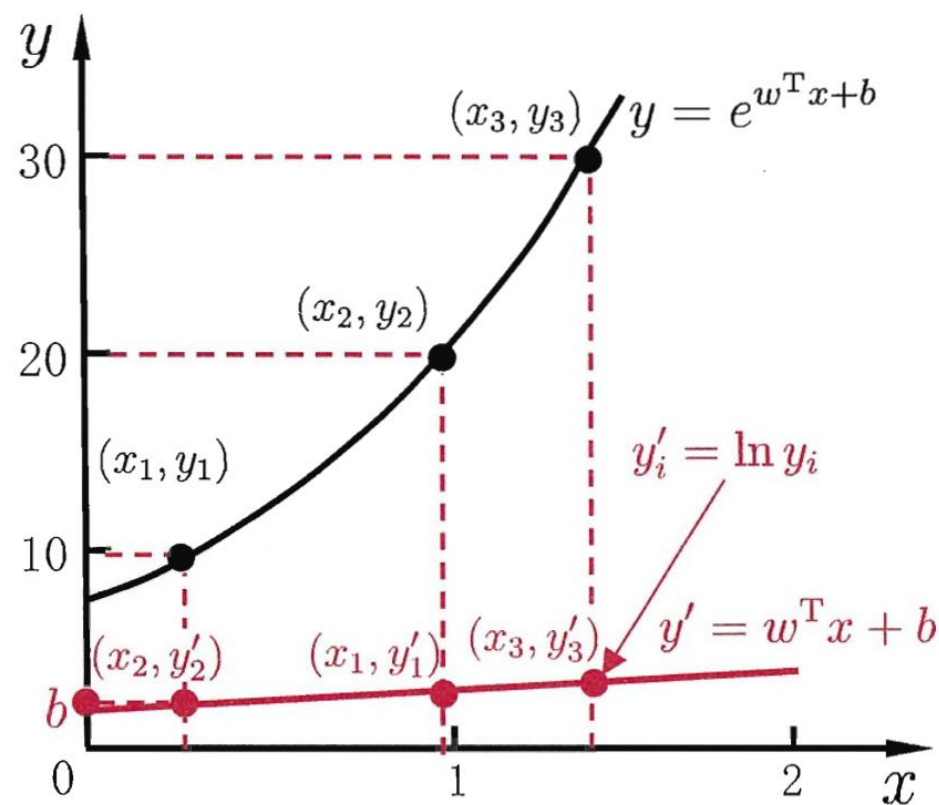
$$\ln y = \mathbf{w}^T \mathbf{x} + b$$

- 广义线性模型

$$y = g^{-1}(\mathbf{w}^T \mathbf{x} + b)$$



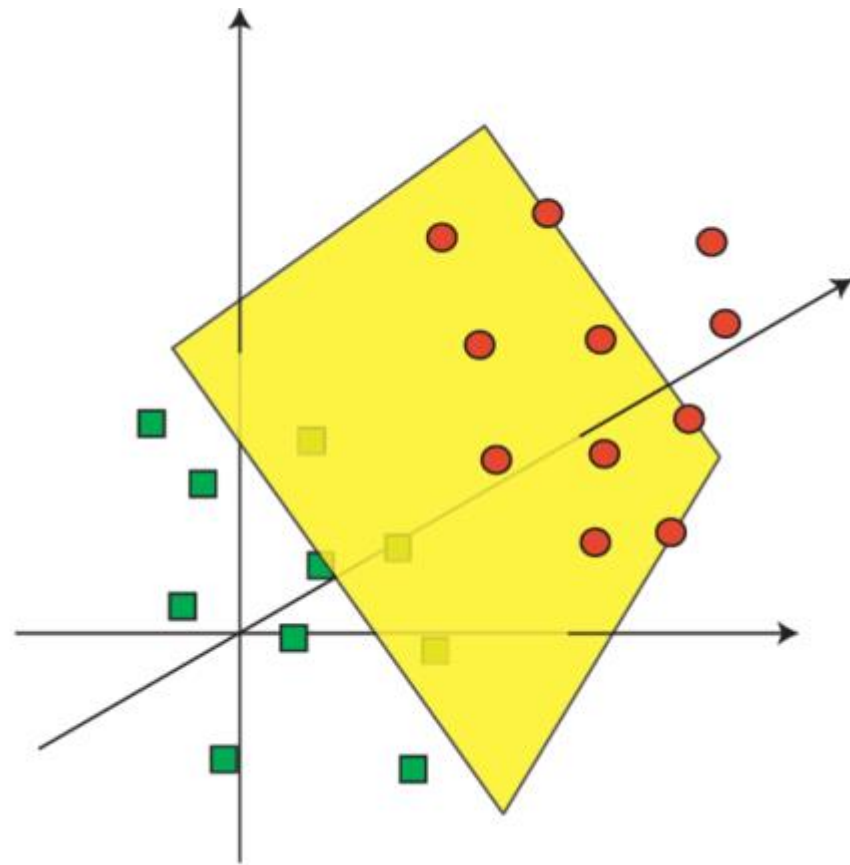
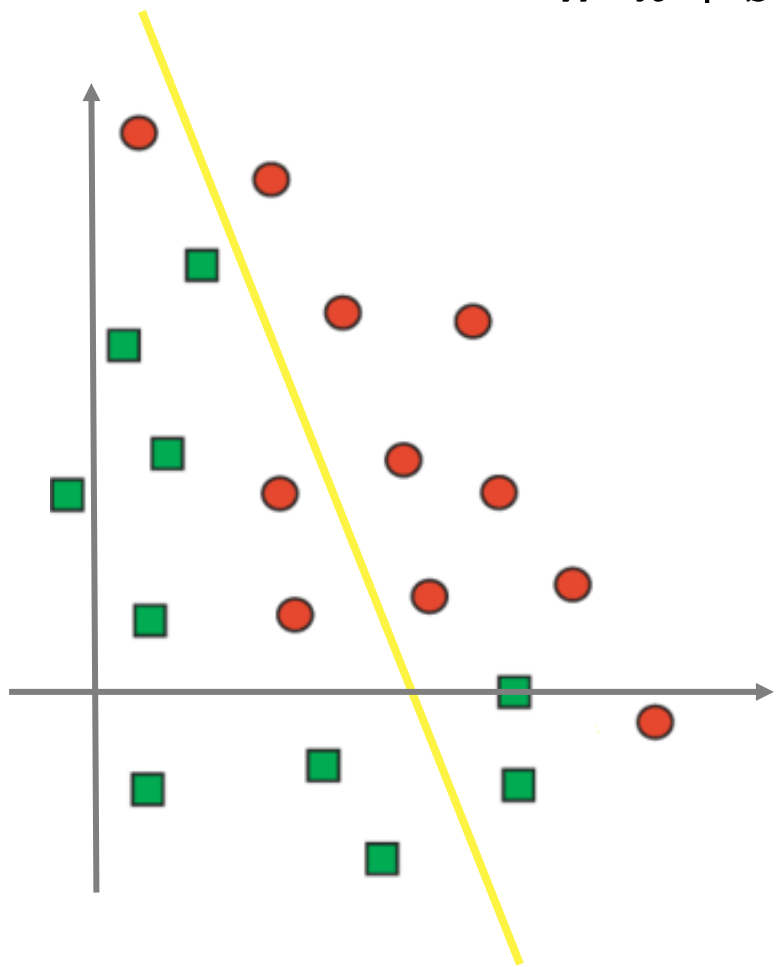
“联系函数”(link function)



线性分类

- 困难

$$\mathbf{w}^T \mathbf{x} + b \xrightarrow{?} y \in \{0, 1\}$$

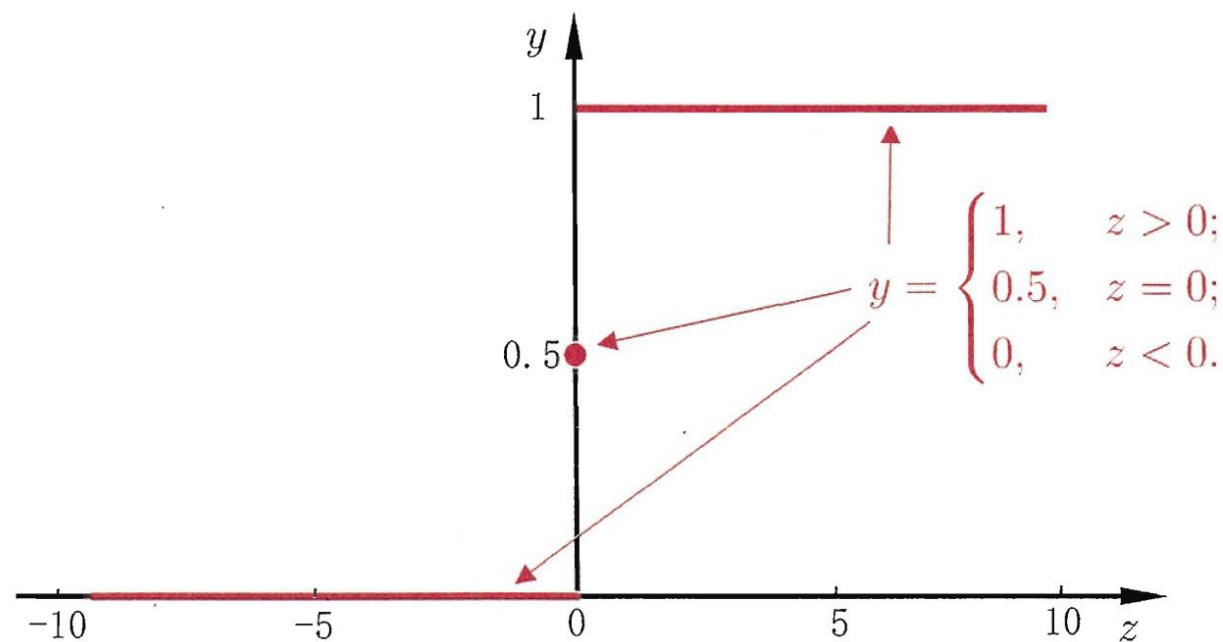


线性分类

- 阶跃函数来联系

$$z = \mathbf{w}^T \mathbf{x} + b$$

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

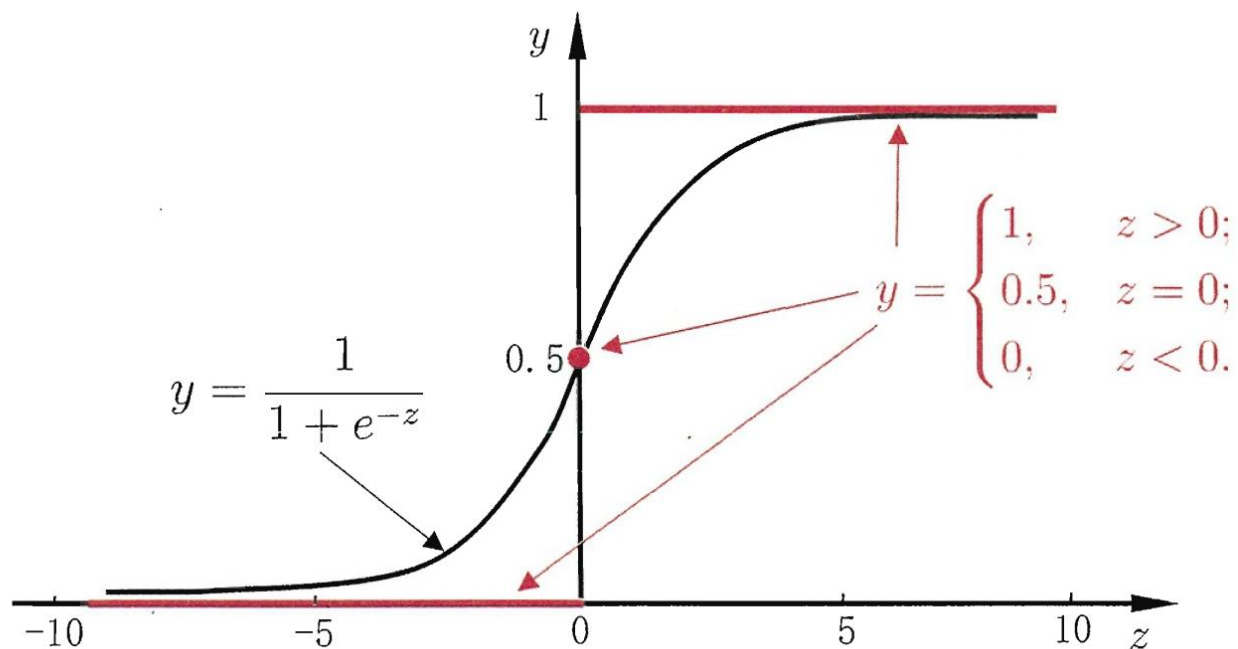


对数几率回归

- Sigmoid函数来联系

$$z = \mathbf{w}^T \mathbf{x} + b$$

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



对数几率回归

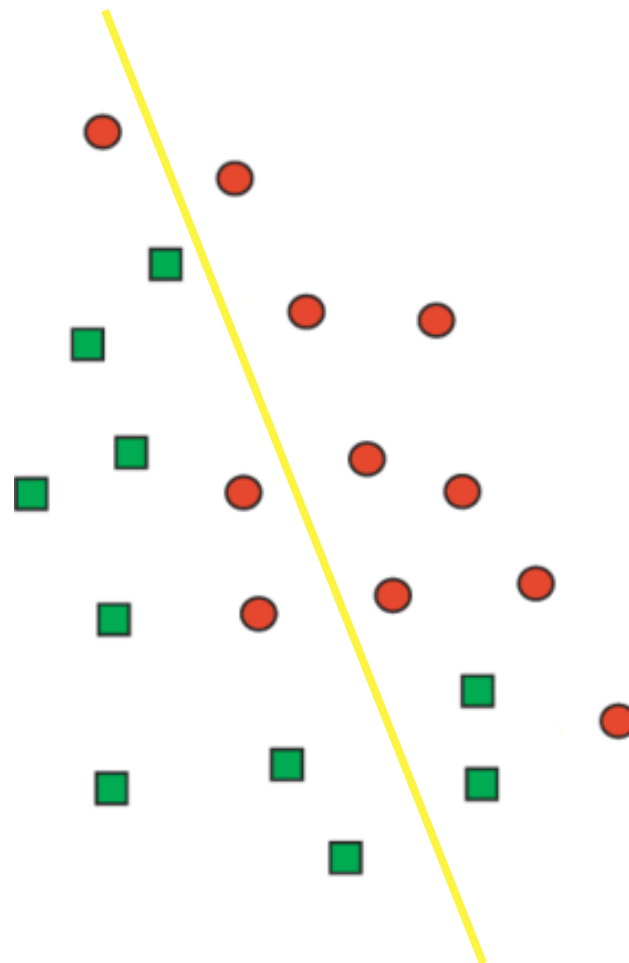
- 回归？分类？

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

$$\ln \frac{y}{1 - y} = \mathbf{w}^T \mathbf{x} + b$$

$$\ln \frac{p(y = 1 | \mathbf{x})}{p(y = 0 | \mathbf{x})}$$

几率
 对数几率
 回归



对数几率回归

- 最大似然法求解

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

$$p(y_i | \mathbf{x}_i; \mathbf{w}, b) = y_i p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta})$$

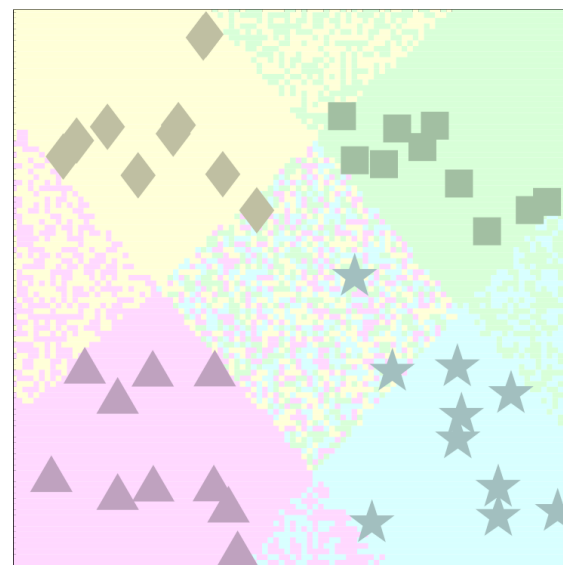
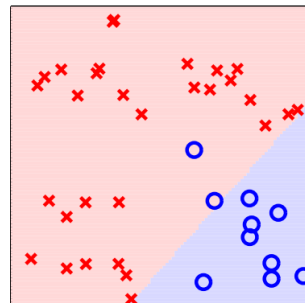
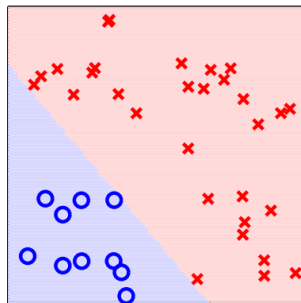
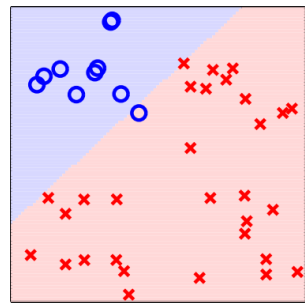
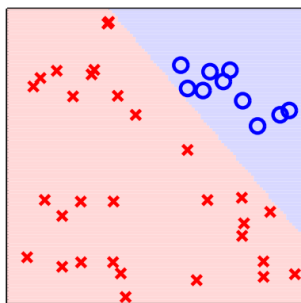
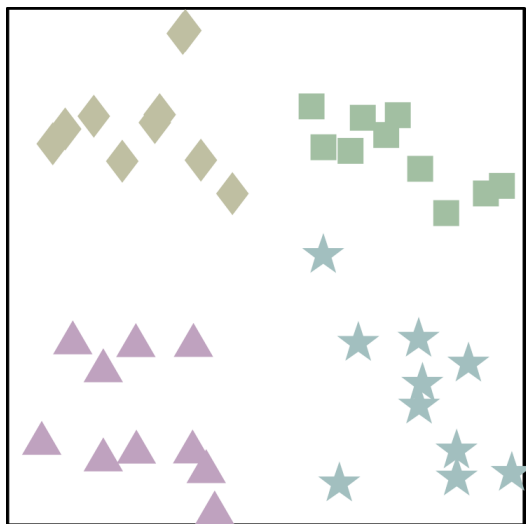
$$\begin{aligned} \ell(\boldsymbol{\beta}) &= \sum_{i=1}^m -\ln p(y_i | \mathbf{x}_i; \mathbf{w}, b) = \sum_{i=1}^m -\ln \left(y_i p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) \right) \\ &= \sum_{i=1}^m (-y_i \boldsymbol{\beta}^T \hat{\mathbf{x}}_i + \ln(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i})) \end{aligned}$$

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta})$$

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

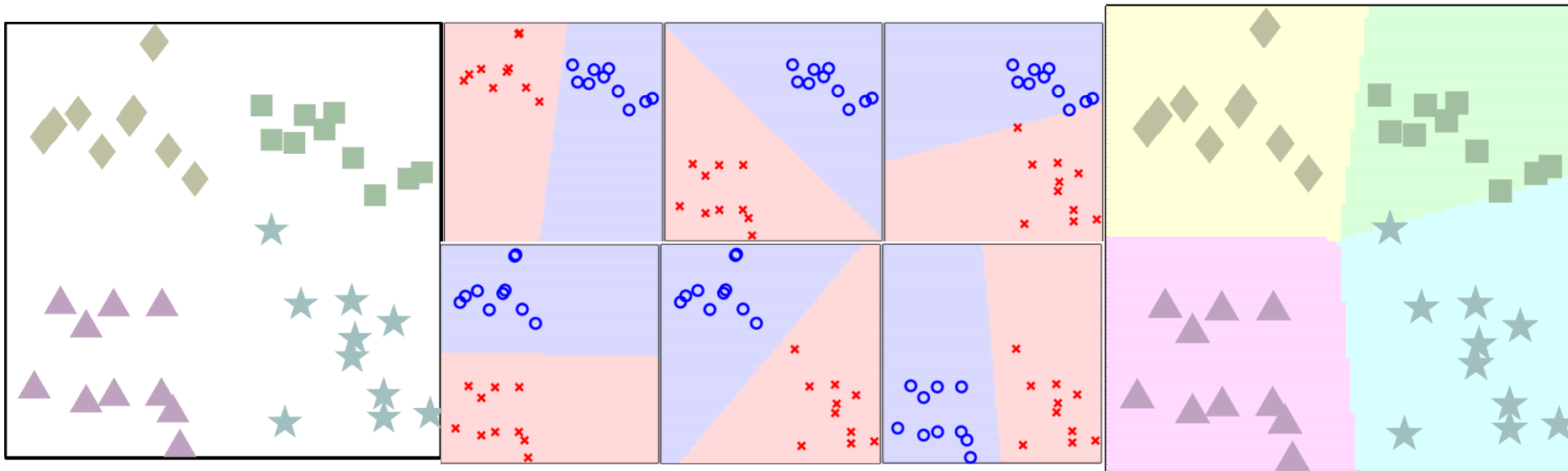
多分类

- 一对其余 One vs. Rest



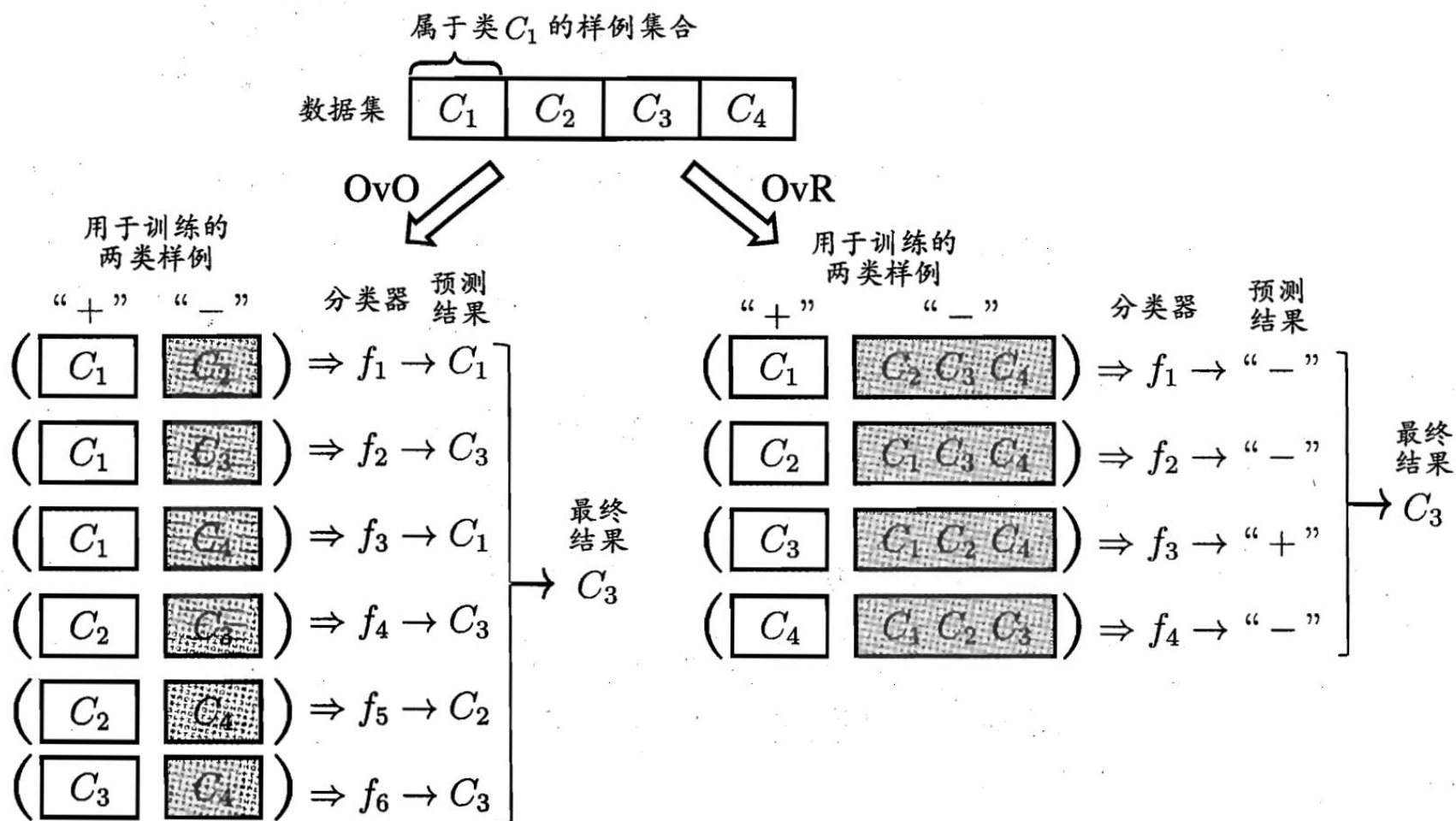
多分类

- 一对一 One vs. One



多分类

• 对比



类别不平衡

- 定义

类别不平衡 (class-imbalance) 就是指分类任务中不同类别的训练样例数目差别很大的情况.

- 策略

- 欠采样

- 过采样

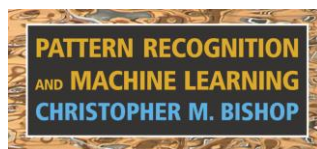
- 阈值移动

MACHINE LEARNING

机器学习

Linear Models

扩展



参考：
《PRML》

Machine Learning Course
Copyright belongs to Wenting Tu.

Softmax回归

- Softmax函数

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \text{softmax} \begin{bmatrix} W_{1,1}x_1 + W_{1,2}x_2 + W_{1,3}x_3 + b_1 \\ W_{2,1}x_1 + W_{2,2}x_2 + W_{2,3}x_3 + b_2 \\ W_{3,1}x_1 + W_{3,2}x_2 + W_{3,3}x_3 + b_3 \end{bmatrix} \begin{matrix} \leftarrow \mathbf{w}_1^\top \mathbf{x} \\ \leftarrow \mathbf{w}_2^\top \mathbf{x} \\ \leftarrow \mathbf{w}_3^\top \mathbf{x} \end{matrix}$$

$$p(y = c \mid \mathbf{x}) = \text{softmax}(\mathbf{w}_c^\top \mathbf{x}) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{c=1}^C \exp(\mathbf{w}_c^\top \mathbf{x})}$$

Softmax回归

- 交叉熵损失函数

$$\mathcal{L}(\mathbf{y}, f(\mathbf{x}, \theta)) = - \sum_{c=1}^C y_c \log f_c(\mathbf{x}, \theta) = -\log f_{\mathbf{y}}(\mathbf{x}, \theta)$$

y_c 是一个 C 维的向量，用来标注样本的多类标签。假设样本的标签为 c ，那么它只有 c 维是1，其余维都为0

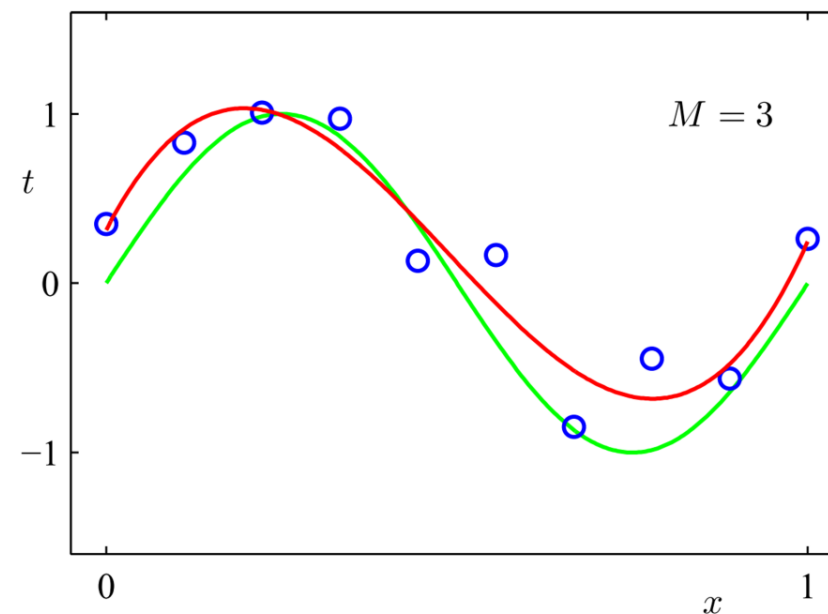
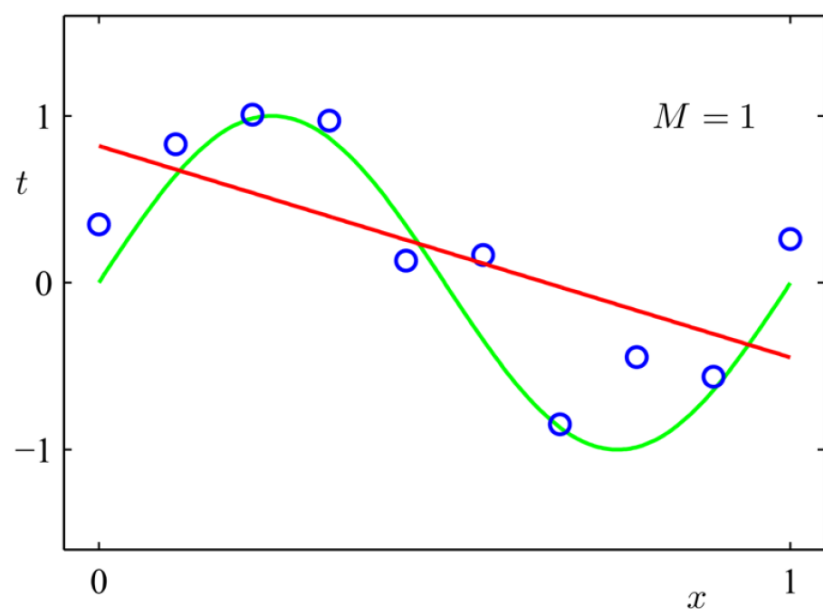
$f_c(\mathbf{x}, \theta)$ 为模型（参数为 θ ）判断的样本为第 c 类的概率

交叉熵角度：当我们将 \mathbf{y} 看做是样本标签的真实概率分布， $f(\mathbf{x}, \theta)$ 看作是类别标签的条件概率分布（模型估计的），那么 $\mathcal{L}(\mathbf{y}, f(\mathbf{x}, \theta))$ 就对应于信息论里面交叉熵的概念，是一种衡量两个分布差距的度量。

最大似然估计角度： $\log f_{\mathbf{y}}(\mathbf{x}, \theta)$ 实际上对应于真实类别 \mathbf{y} 的对数似然函数

线性基函数模型 *Linear Basis Function Models*

• 示例



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 + \cdots + w_Mx^M$$

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 + \cdots + w_Mx^M$$

线性基函数模型 *Linear Basis Function Models*

- 定义

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \quad \phi_0(\mathbf{x}) = 1$$

$$\phi_j(x) = x^j$$

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

$$\phi_j(x) = \sigma \left(\frac{x - \mu_j}{s} \right) \quad \sigma_a = \frac{1}{1 + \exp(-a)}$$

线性基函数模型 *Linear Basis Function Models*

- 定义

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \quad \phi_0(\mathbf{x}) = 1$$

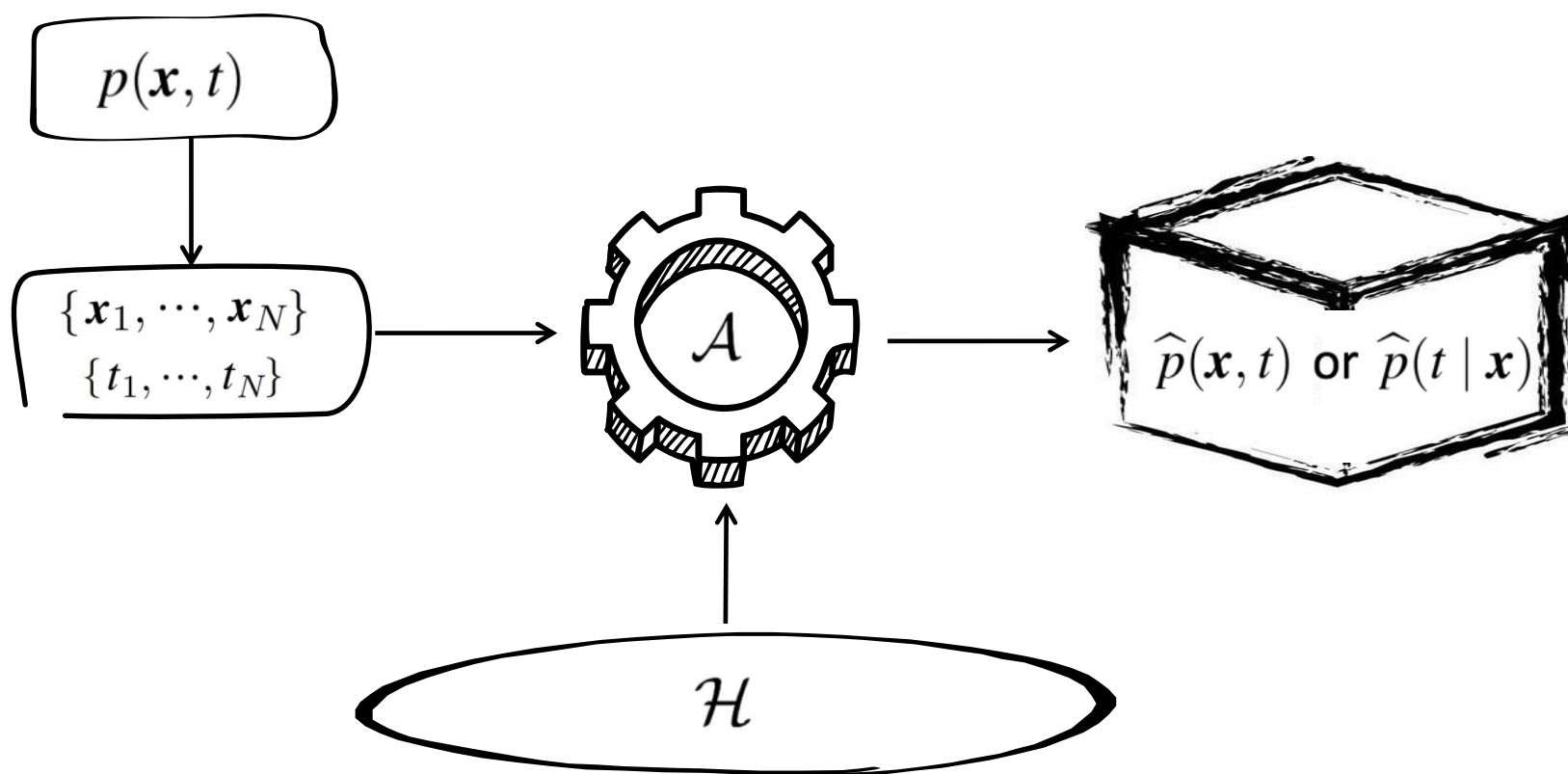
$$\phi_j(x) = x^j$$

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

$$\phi_j(x) = \sigma \left(\frac{x - \mu_j}{s} \right) \quad \sigma_a = \frac{1}{1 + \exp(-a)}$$

最大似然法与最小二乘法

- 回归任务的概率体系框架



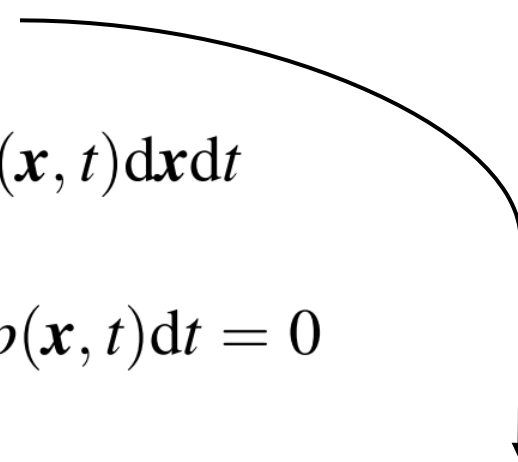
最大似然法与最小二乘法

- 关于回归任务的决策理论

$$L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$$

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

$$\frac{\partial \mathbb{E}[L]}{\partial y(\mathbf{x})} = 2 \int \{y(\mathbf{x}) - t\} p(\mathbf{x}, t) dt = 0$$

$$y(\mathbf{x}) = \frac{\int tp(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int tp(t | \mathbf{x}) dt = \mathbb{E}_t[t | \mathbf{x}]$$


最大似然法与最小二乘法

- 最小二乘法与最大似然估计

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

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

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \quad \mathbf{t} = \{t_1, \dots, t_N\}$$

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

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

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

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

$$\mathbf{w}_{MLE}^* = \mathbf{w}_{LS}^*$$