

# 机器学习

## Machine Learning

### 第三讲：回归的基本学习算法

# 1. 线性回归模型

特征向量（输入向量）

$$\mathbf{x} = [x_1, x_2, \cdots, x_K]^T$$

扩充特征向量

$$\bar{\mathbf{x}} = [1, x_1, x_2, \cdots, x_K]^T \quad \text{即 } x_0 = 1$$

权系数向量

$$\mathbf{w} = [w_0, w_1, w_2, \cdots, w_K]^T$$

则线性回归模型为

$$\hat{y}(\mathbf{x}, \mathbf{w}) = \sum_{k=0}^K w_k x_k = \mathbf{w}^T \bar{\mathbf{x}}$$

## 线性回归模型（续）

给出训练序列

$$\begin{aligned}\mathbf{D} &= \left\{ (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \cdots, (\mathbf{x}_N, y_N) \right\} \\ &= \left\{ (\mathbf{x}_n, y_n) \right\}_{n=1}^N\end{aligned}$$

对于给出的损失函数最小化，得到回归参数  $\mathbf{w}$

对于给出的新特征向量  $\mathbf{x}$ ，得到预测值

$$\hat{y}(\mathbf{x}, \mathbf{w}) = \sum_{k=0}^K w_k x_k = \mathbf{w}^T \bar{\mathbf{x}}$$

## 2. 扩充：线性基函数回归模型

特点：对参数向量线性，对特征向量非线性

对特征向量映射基函数：

$$\mathbf{x} = [x_1, x_2, \dots, x_K]^T \Rightarrow$$

$$\boldsymbol{\phi}(\mathbf{x}) = [\phi_0(\mathbf{x}), \phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x})]^T$$

线性基函数模型：

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

特例：线性回归是线性基函数回归的特例： $\boldsymbol{\phi}(\mathbf{x}) = \bar{\mathbf{x}}$

## 基函数集的例子-1

$$\mathbf{x} = [x_1, x_2, x_3]^T \Rightarrow$$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1x_2, x_1x_3, x_2x_3]^T$$

对应线性回归模型

$$\hat{y}(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + w_2x_2 + w_3x_3$$

对应线性基函数回归模型

$$\begin{aligned} \hat{y}(\mathbf{x}, \mathbf{w}) = & w_0 + w_1x_1 + w_2x_2 + w_3x_3 + w_4x_1^2 + \\ & + w_5x_2^2 + w_6x_3^2 + w_7x_1x_2 + w_8x_1x_3 + w_9x_2x_3 \end{aligned}$$

### 3. 基本线性回归模型的学习

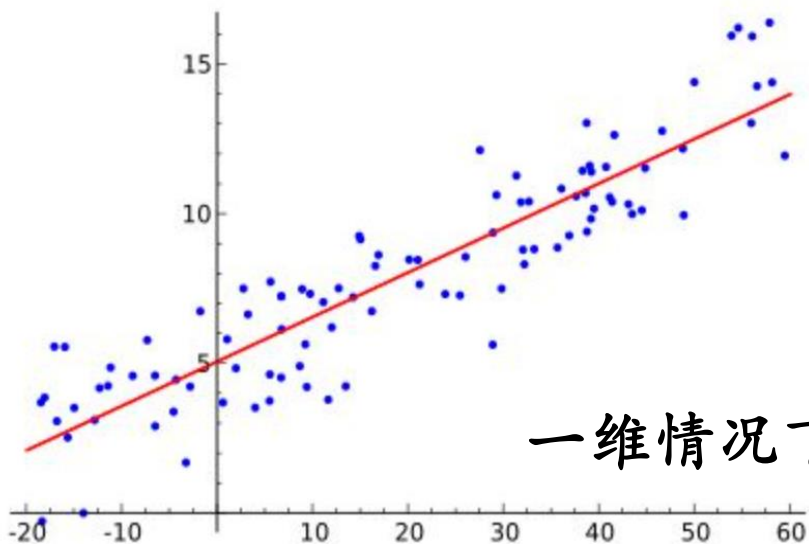
独立同分布条件 (I.I.D) 的训练数据集

$$\mathbf{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$$

对每个样本，模型与标注之间存在误差  $\varepsilon_i$

$$y_i = \hat{y}(\mathbf{x}_i, \mathbf{w}) + \varepsilon_i = \mathbf{w}^T \bar{\mathbf{x}}_i + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$$



一维情况下，线性回归表示的例子

## 通过最大似然原理导出线性回归的学习算法

$y_i$  的概率密度函数

$$p_y(y_i | \mathbf{w}) = \frac{1}{(2\pi\sigma_\varepsilon^2)^{\frac{1}{2}}} \exp \left[ -\frac{1}{2\sigma_\varepsilon^2} (y_i - \hat{y}(\mathbf{x}_i, \mathbf{w}))^2 \right]$$

所有样本的标注值表示为向量

$$\mathbf{y} = [y_1, y_2, y_3, \dots, y_N]^T$$

由样本集的I.I.D性得似然函数

$$\begin{aligned} p_y(\mathbf{y} | \mathbf{w}) &= \prod_{i=1}^N p_y(y_i | \mathbf{w}) \\ &= \frac{1}{(2\pi\sigma_\varepsilon^2)^{\frac{N}{2}}} \exp \left[ -\frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^N (y_i - \hat{y}(\mathbf{x}_i, \mathbf{w}))^2 \right] \end{aligned}$$

## ML导出线性回归的学习算法（续）

取对数似然为

$$\log p_y(\mathbf{y}|\mathbf{w}) = -\frac{N}{2}\log(2\pi\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^N (y_i - \hat{y}(\mathbf{x}_i, \mathbf{w}))^2$$

对数似然函数最大，对应平方误差和最小

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}(\mathbf{x}_i, \mathbf{w}))^2$$

重写误差平方和

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \varepsilon_i^2 = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \bar{\mathbf{x}}_i)^2 \\ &= \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &= \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \end{aligned}$$

其中

$$\mathbf{X} = \begin{bmatrix} \bar{\mathbf{x}}_1^T \\ \bar{\mathbf{x}}_2^T \\ \vdots \\ \bar{\mathbf{x}}_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1K} \\ 1 & x_{21} & \cdots & x_{2K} \\ \vdots & \vdots & \cdots & \vdots \\ 1 & x_{N1} & \cdots & x_{NK} \end{bmatrix}$$



## ML导出线性回归的学习算法（续）

求解系数向量

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{0}$$

系数向量满足方程

$$\mathbf{X}^T \mathbf{X} \mathbf{w}_{ML} = \mathbf{X}^T \mathbf{y}$$

若 $\mathbf{X}^T \mathbf{X}$ 可逆，解为（线性回归的最小二乘（LS）解）

$$\mathbf{w}_{ML} = \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

权系数向量得到后，线性回归函数确定为

$$\hat{y}(\mathbf{x}, \mathbf{w}) = \mathbf{w}_{ML}^T \bar{\mathbf{x}}$$

# 线性回归的几何解释

线性回归在训练集上的输出向量

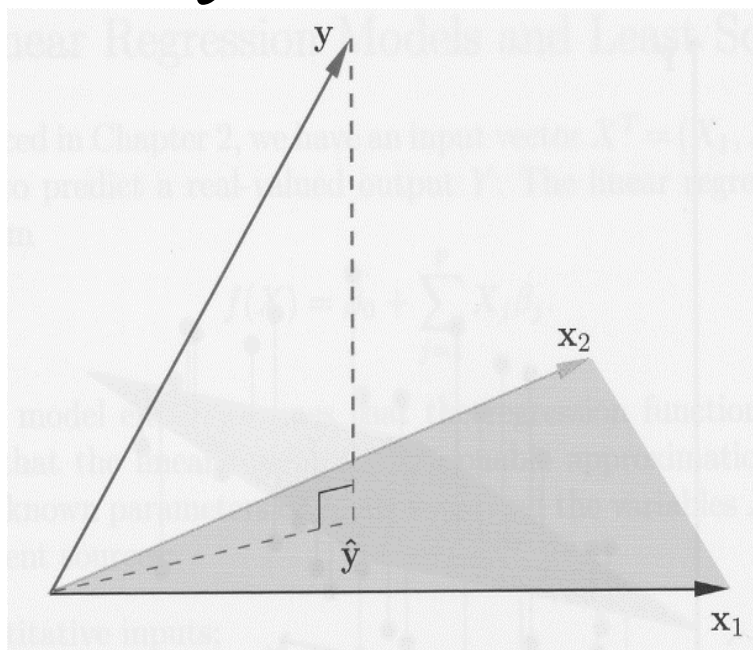
$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}_{ML} = \begin{bmatrix} \tilde{\mathbf{x}}_0 & \tilde{\mathbf{x}}_1 & \cdots & \tilde{\mathbf{x}}_K \end{bmatrix} \mathbf{w}_{ML} = \sum_{i=0}^K w_{ML,i} \tilde{\mathbf{x}}_i$$

线性回归对每一个标注值的逼近误差向量

$$\boldsymbol{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{P}\mathbf{y} = \mathbf{P}^\perp \mathbf{y}$$

回归误差正交性

$$\boldsymbol{\varepsilon}^T \hat{\mathbf{y}} = 0$$



## 4. 线性回归模型的递推学习

对数据集计算平均梯度

$$\frac{1}{N} \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^{(k)}} = -\frac{1}{N} \sum_{i=1}^N \left( y_i - \mathbf{w}^{(k)\top} \bar{\mathbf{x}}_i \right) \bar{\mathbf{x}}_i$$

从初始猜测值  $\mathbf{w}^{(0)}$  按梯度下降算法更新

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \mathbf{w}^{(k)} - \eta_k \frac{1}{N} \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^{(k)}} \\ &= \mathbf{w}^{(k)} + \frac{\eta_k}{N} \sum_{i=1}^N \left( y_i - \mathbf{w}^{(k)\top} \bar{\mathbf{x}}_i \right) \bar{\mathbf{x}}_i \end{aligned}$$

# 随机梯度下降算法

(stochastic gradient descent, SGD)

误差和可分解为

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left( y_i - \mathbf{w}^T \bar{\mathbf{x}}_i \right)^2 = \sum_{i=1}^N J_i(\mathbf{w})$$

一个样本的梯度

$$\left. \frac{\partial J_i(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^{(k)}} = - \left( y_i - \mathbf{w}^{(k)T} \bar{\mathbf{x}}_i \right) \bar{\mathbf{x}}_i$$

利用一个样本梯度对权系数向量的更新

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \mathbf{w}^{(k)} - \eta \left. \frac{\partial J_i(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^{(k)}} \\ &= \mathbf{w}^{(k)} + \eta_k \left( y_i - \mathbf{w}^{(k)T} \bar{\mathbf{x}}_i \right) \bar{\mathbf{x}}_i \end{aligned}$$

梯度是随机的，  
每次迭代样本  
是随机选取的

# 小批量SGD算法

从数据集随机抽取一小批量样本

$$\mathbf{D}_{k+1} = \left\{ (\mathbf{x}_m, y_m) \right\}_{m=1}^{N_1}$$

小批量SGD算法如下

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \eta \frac{1}{N_1} \sum_{m=1}^{N_1} \left( y_m - \mathbf{w}^{(k)\top} \bar{\mathbf{x}}_m \right) \bar{\mathbf{x}}_m$$

一般收敛性条件

$$\sum_{k=1}^{\infty} \eta_k = \infty$$

$$\sum_{k=1}^{\infty} \eta_k^2 < \infty$$

## 正则化目标函数 与贝叶斯框架的等价性！

### 5. 正则化线性回归

正则化：在目标函数中增加约束参数向量的量

一种常用选择为参数向量的范数平方约束

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_{i=1}^N \left( y_i - \mathbf{w}^T \bar{\mathbf{x}}_i \right)^2 + \frac{\lambda}{2} \sum_{i=1}^K w_i^2 \\ &= \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

参数向量的正则化LS解为

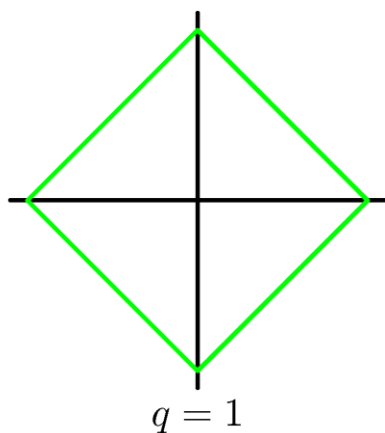
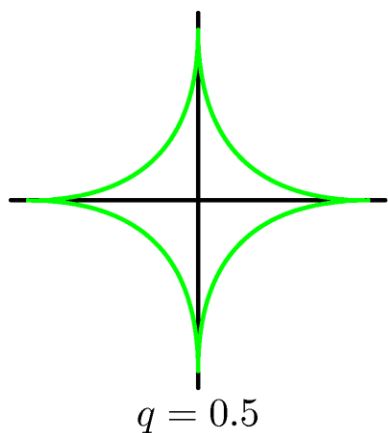
$$\mathbf{w}_R = \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

# 更一般的正则化

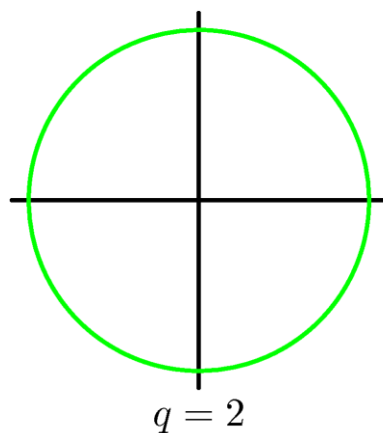
注：贝叶斯框架下，一个 $q$ 取值  
对应参数向量的一种不同的先验  
概率假设

对更一般的 $q$ 值，定义正则化目标函数

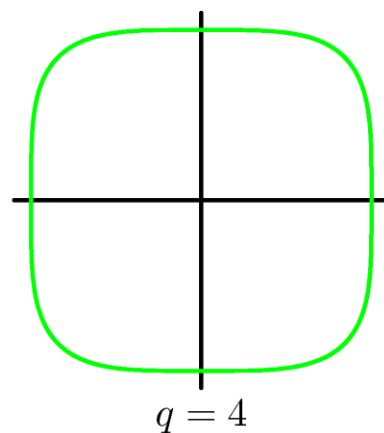
$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left( y_i - \mathbf{w}^T \bar{\mathbf{x}}_i \right)^2 + \frac{\lambda}{2} \sum_{i=1}^K |w_i|^q$$



Lasso

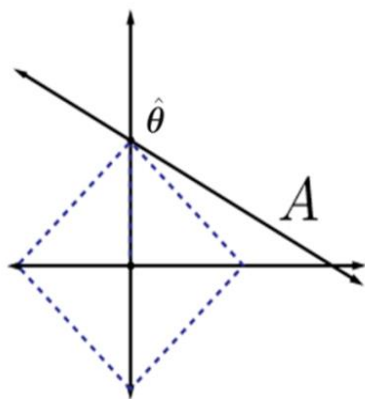


Quadratic

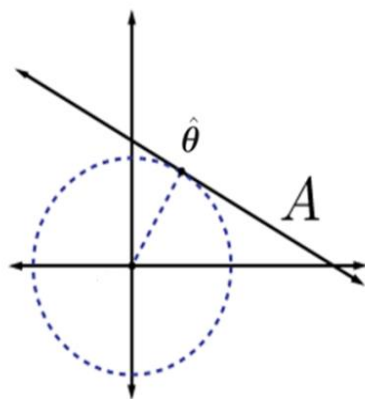


# 更一般的正则化

$q=1$ , Lasso解趋向于稀疏化



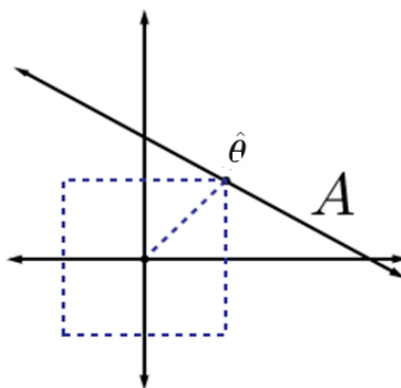
$p = 1$



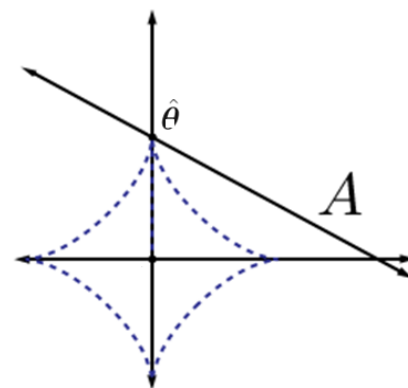
$p = 2$

$q < 1$ : 稀疏化

$q > 1$ : 非稀疏化



$p = \infty$



$p = \frac{1}{2}$



## 6. 线性基函数回归

输入向量映射:  $\mathbf{x} \mapsto \boldsymbol{\phi}(\mathbf{x})$

$$\boldsymbol{\phi}(\mathbf{x}) = [\phi_0(\mathbf{x}), \phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})]^T$$

线性基函数回归模型

$$\hat{y}(\boldsymbol{\phi}, \mathbf{w}) = \sum_{k=0}^M w_k \phi_k(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

线性基函数回归系数向量的解为

$$\mathbf{w}_{ML} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{y}$$

基函数数据矩阵为

$$\begin{aligned}\boldsymbol{\Phi} &= \begin{bmatrix} \boldsymbol{\phi}^T(\mathbf{x}_1) \\ \boldsymbol{\phi}^T(\mathbf{x}_2) \\ \vdots \\ \boldsymbol{\phi}^T(\mathbf{x}_N) \end{bmatrix} \\ &= \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_M(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_M(\mathbf{x}_2) \\ \vdots & \vdots & \cdots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_M(\mathbf{x}_N) \end{bmatrix}\end{aligned}$$

三维向量例子  $\mathbf{x}_n = [x_{n,1}, x_{n,2}, x_{n,3}]^T$

$$\boldsymbol{\phi}(\mathbf{x}_n) = [\phi_0(\mathbf{x}_n), \phi_1(\mathbf{x}_n), \cdots, \phi_9(\mathbf{x}_n)]^T$$

$$= [1, x_{n,1}, x_{n,2}, x_{n,3}, x_{n,1}^2, x_{n,2}^2, x_{n,3}^2, x_{n,1}x_{n,2}, x_{n,2}x_{n,3}, x_{n,1}x_{n,3}]^T$$

# 数值实例

内在输入输出模型

$$f(x) = \frac{1}{1 + \exp(-5x)}$$

采样样本方式

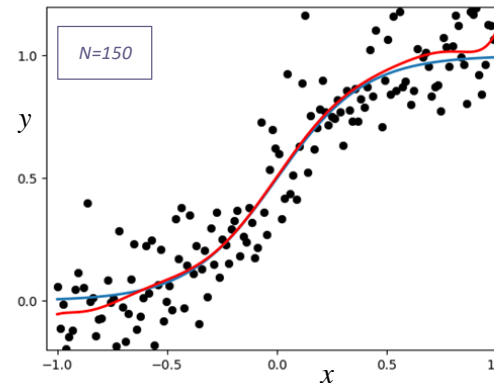
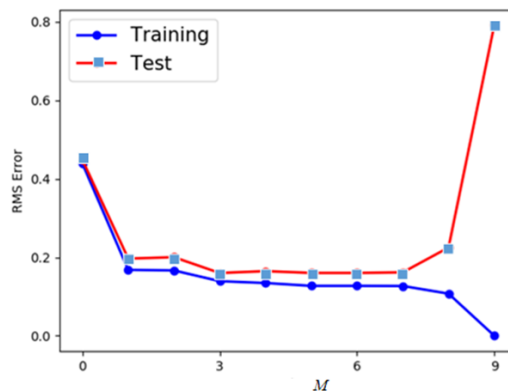
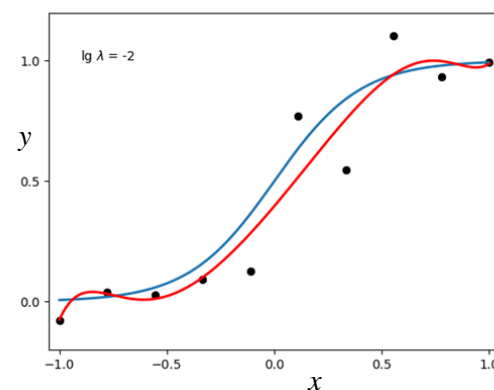
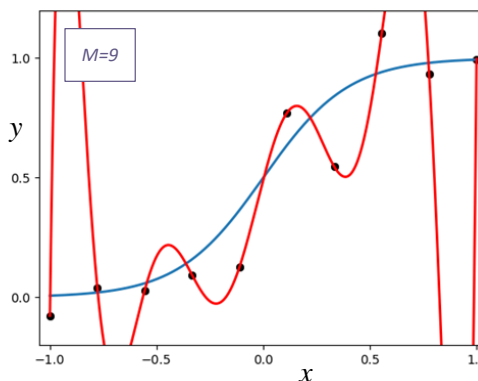
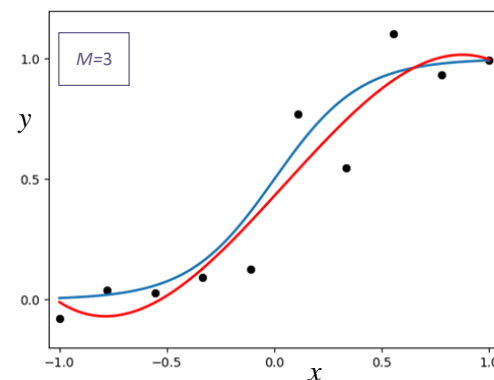
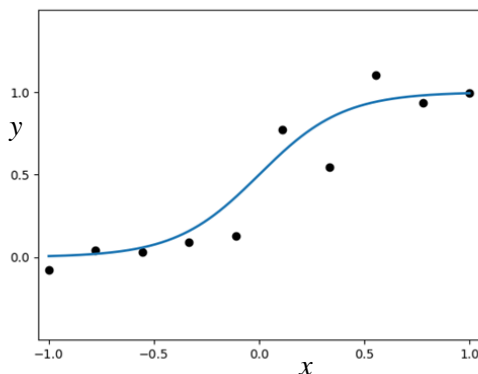
$$y_n = f(x_n) + \varepsilon_n$$

$$\varepsilon_n \sim N(0, 0.15^2)$$

基函数回归

$$\phi(x_n) = [1, x_n, x_n^2, \dots, x_n^M]^T$$

右图：不同参数的实验结果  
注意过拟合和正则化



## 7. 机器学习模型的误差分解

以回归模型作为讨论对象，考虑平方误差

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

模型的误差期望（泛化误差）

$$E(L) = \int \int (\hat{y}(\mathbf{x}) - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

最优回归模型是  $h(\mathbf{x}) = E(y|\mathbf{x})$

误差分解

$$\begin{aligned} E(L) &= \int \int (\hat{y}(\mathbf{x}) - h(\mathbf{x}) + h(\mathbf{x}) - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy \\ &= \int (\hat{y}(\mathbf{x}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \int \int (E(y|\mathbf{x}) - y)^2 p(\mathbf{x}, y) d\mathbf{x} dy \end{aligned}$$

# 模型的误差分解

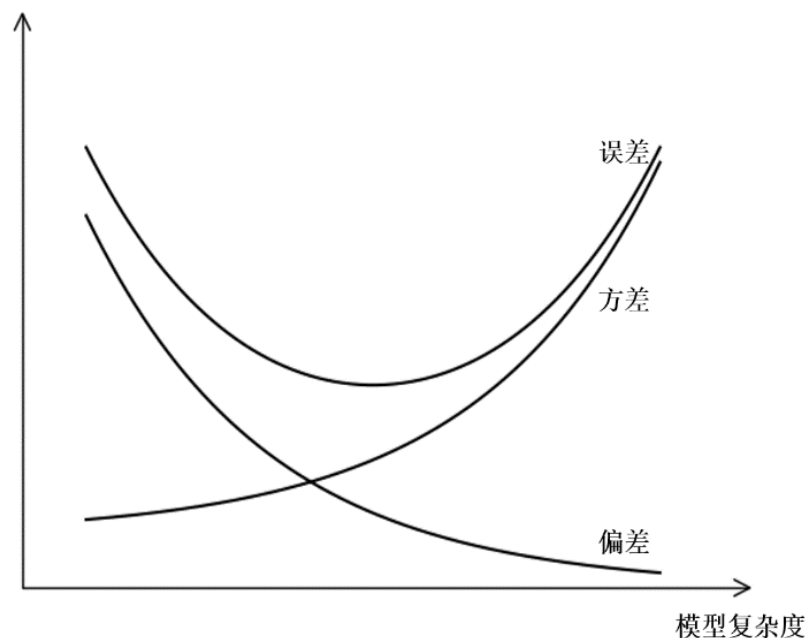
$$\begin{aligned} E(L) &= \\ &= \int \left[ E_D \left( \hat{y}(\mathbf{x}; \mathbf{D}) \right) - h(\mathbf{x}) \right]^2 p(\mathbf{x}) d\mathbf{x} \\ &\quad + \int E_D \left\{ \left[ \hat{y}(\mathbf{x}; \mathbf{D}) - E_D \left( \hat{y}(\mathbf{x}; \mathbf{D}) \right) \right]^2 \right\} p(\mathbf{x}) d\mathbf{x} \\ &\quad + \int \int \left( E(y|\mathbf{x}) - y \right)^2 p(\mathbf{x}, y) d\mathbf{x} dy \\ &= (bias)^2 + \text{方差} + \text{固有误差} \end{aligned}$$

泛化误差由三部分组成：

偏、方差和固有误差

模型简单，方差小，偏大；

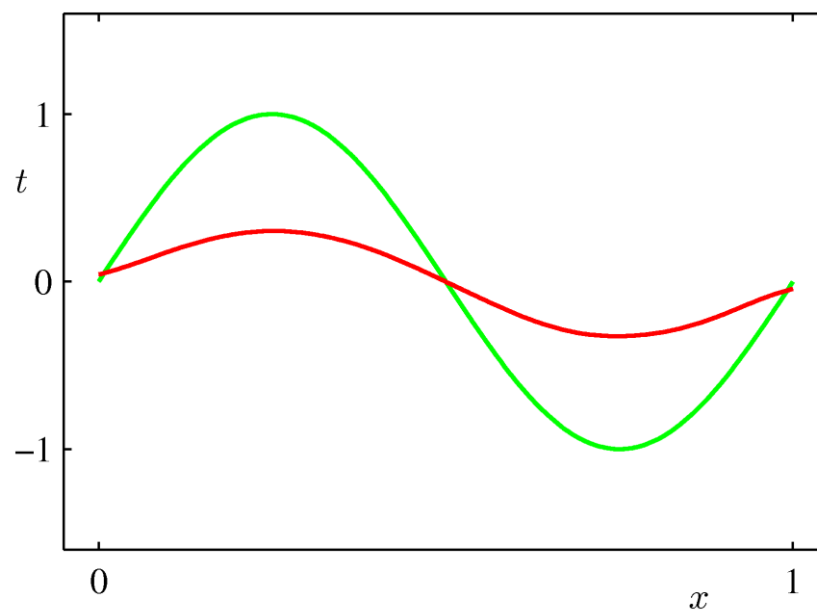
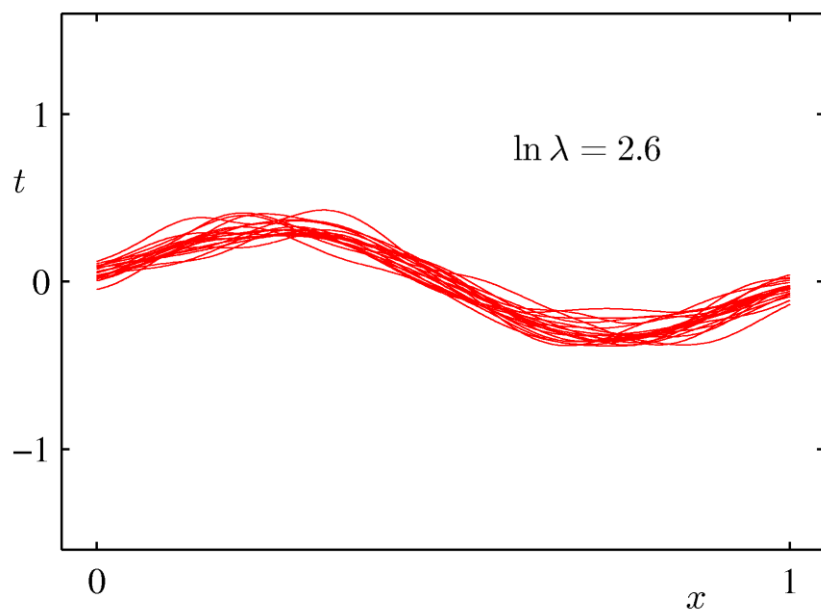
模型复杂，方差大，偏小



# 误差分解的数值实例

---

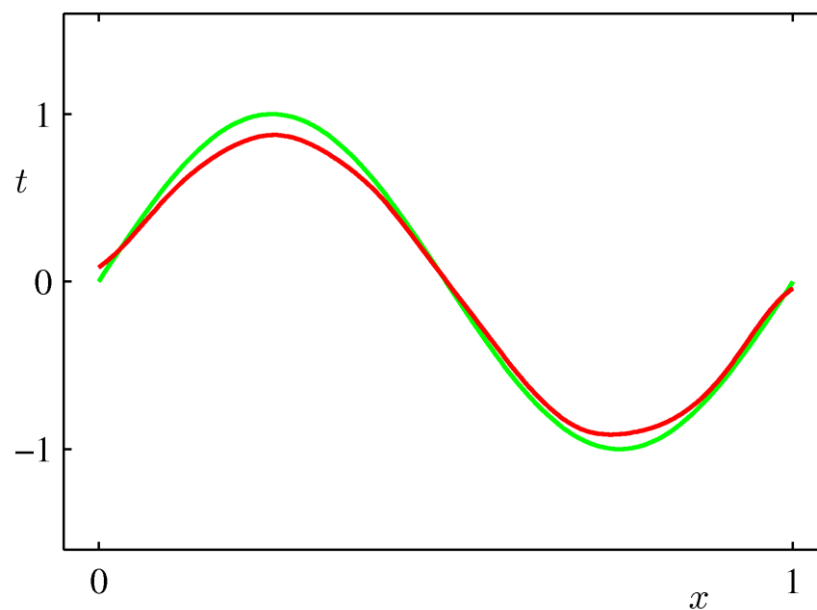
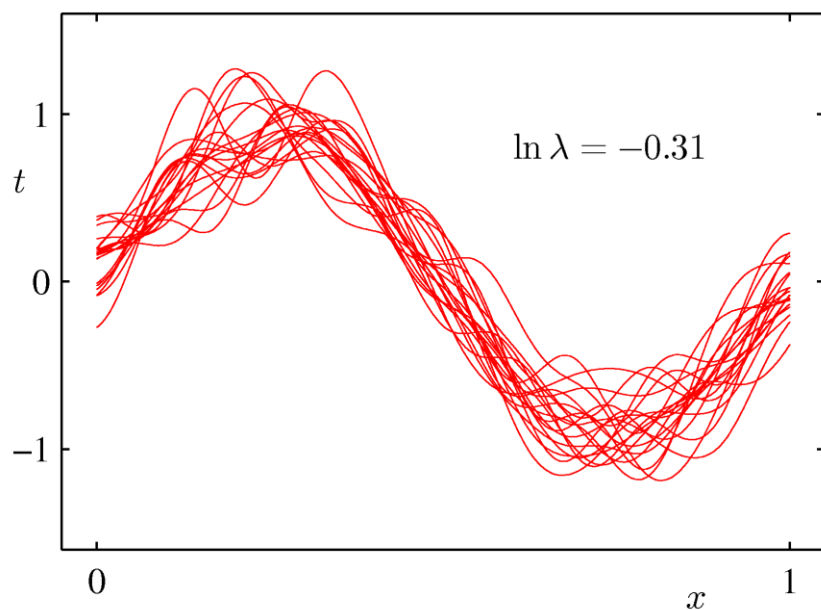
例: 25 个数据集, 在基函数回归下, 变化正则化参数  $\lambda$ . 模型误差的方差和偏折中示例



# 误差分解的数值实例

---

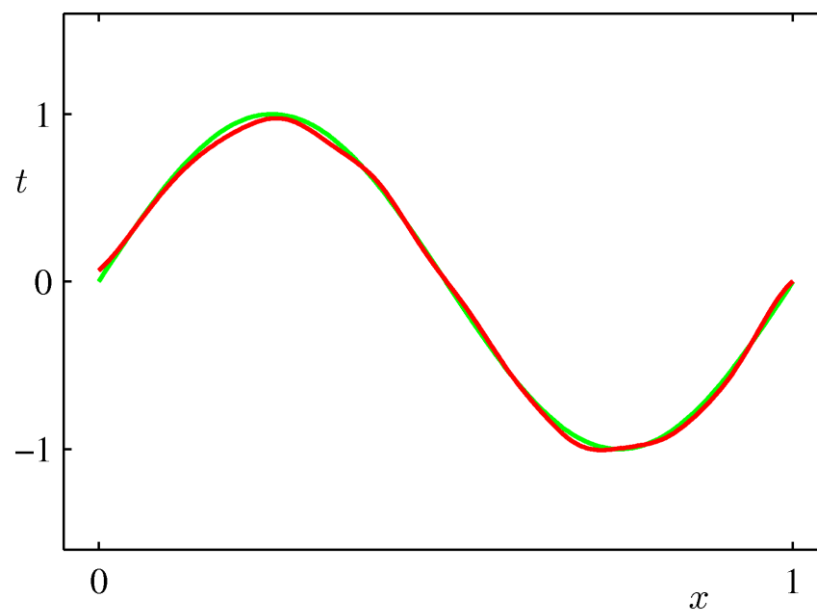
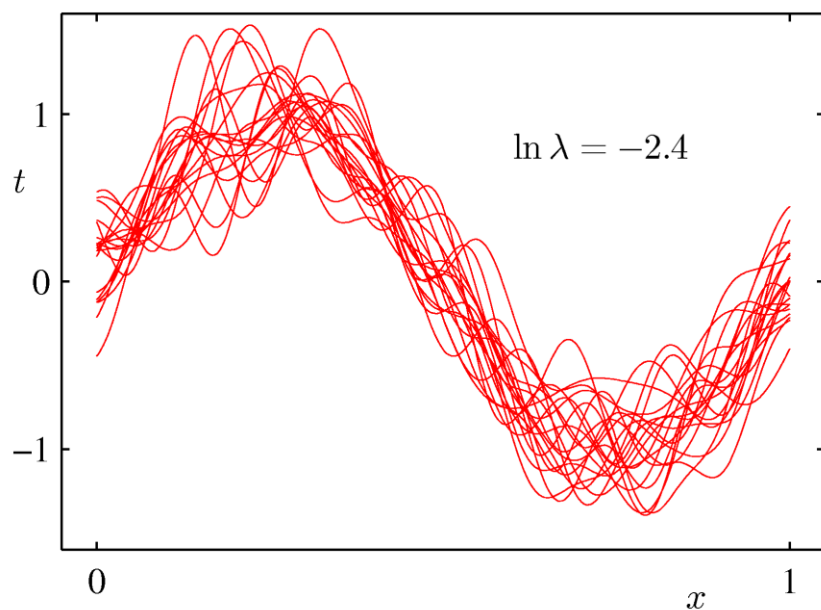
例: 25 个数据集, 在基函数回归下, 变化正则化参数  $\lambda$ . 模型误差的方差和偏折中示例



# 误差分解的数值实例

---

例: 25 个数据集, 在基函数回归下, 变化正则化参数  $\lambda$ . 模型误差的方差和偏折中示例





# 误差分解实例

$f(x)$ 无法直接观测到，采样过程为

$$y = f(x) + v$$

采样数据构成I.I.D数据集  $\{x_i, y_i\}_{i=1}^N$

采用K近邻回归算法训练模型为

(非参数模型，且K越小对应模型越复杂)

$$\hat{y} = \hat{f}(x) = \frac{1}{K} \sum_{l=1}^K y_{(l)}$$

该问题的统计最优模型为

$$h(x) = E(y|x) = f(x)$$

直接求得误差分解为

$$\begin{aligned} E(L) &= E \left\{ (\hat{y} - h(x))^2 \right\} + \sigma_v^2 \\ &= E \left\{ \left( \frac{1}{K} \left( \sum_{l=1}^K f(x_{(l)}) + v_l \right) - f(x) \right)^2 \right\} + \sigma_v^2 \\ &= E \left\{ \left[ \frac{1}{K} \sum_{l=1}^K f(x_{(l)}) - f(x) \right]^2 \right\} + E \left\{ \left( \frac{1}{K} \sum_{l=1}^K v_l \right)^2 \right\} + \sigma_v^2 \\ &= \left[ \frac{1}{K} \sum_{l=1}^K f(x_{(l)}) - f(x) \right]^2 + \frac{\sigma_v^2}{K} + \sigma_v^2 \end{aligned}$$

第一项是偏，K越大模型越简单，偏越大；第二项是方差，K越大，模型越简单方差越小

- 注释** 测试误差和泛化误差曲线呈现一种类似“U”曲线。对于传统的单一机器学习模型“U”曲线具有一般性；
- 但在深度学习中，当深度网络复杂度达到一定规模后，测试误差的表现更加复杂；
- 对于集成学习中一些方法，如随机森林和提升算法，测试误差一般也并没有呈现出“U”曲线，换言之，集成学习更不易出现过拟合问题。
- 机器学习是仍在快速发展中的领域，在发展中一些传统结论，可能被不断补充和修改。