

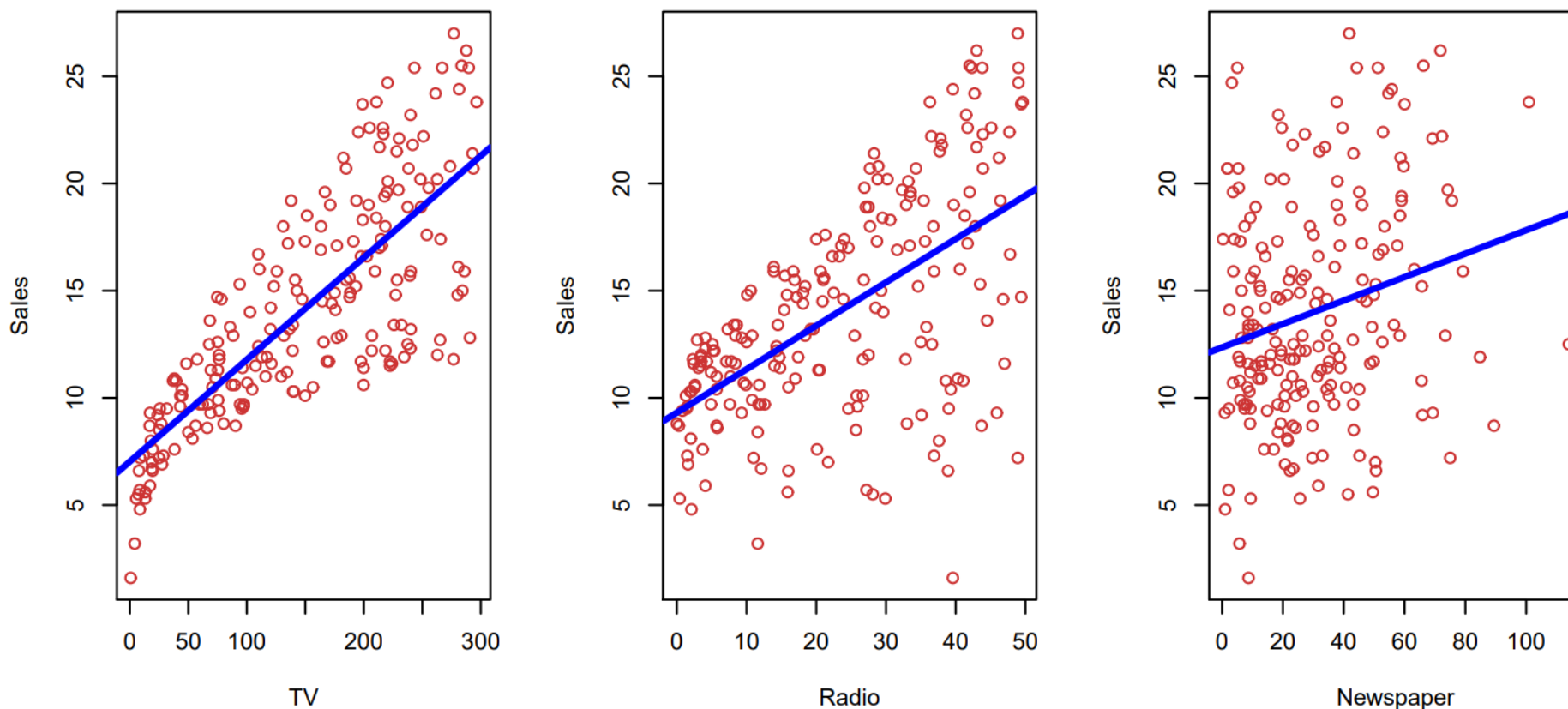


第二章 统计学习



1.什么是统计学习?

问题：是否可以用这三种投入来预测销售的结果？



散点图：TV, Radio and Newspaper的投入 vs 销售额

蓝线：线性回归拟合三种结果（linear-regression）

$$\text{Sales} \approx f(\text{TV}, \text{Radio}, \text{Newspaper})$$

标记方法

X: 输入变量, 特征变量, 属性变量, 预测变量, 自变量, 变量。如: 电视X1, 收音机X2, 使用输入向量集合

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$$

标记方法

Y: 输出变量, 响应变量, 目标变量, 因变量, 输出结果等。

如: 希望能够预测的销售额, 用Y来指代。

模型: 预测模型, 分类模型等。

$$Y = f(X) + \varepsilon$$

ε : 随机误差项, 与X独立, 均值为0.

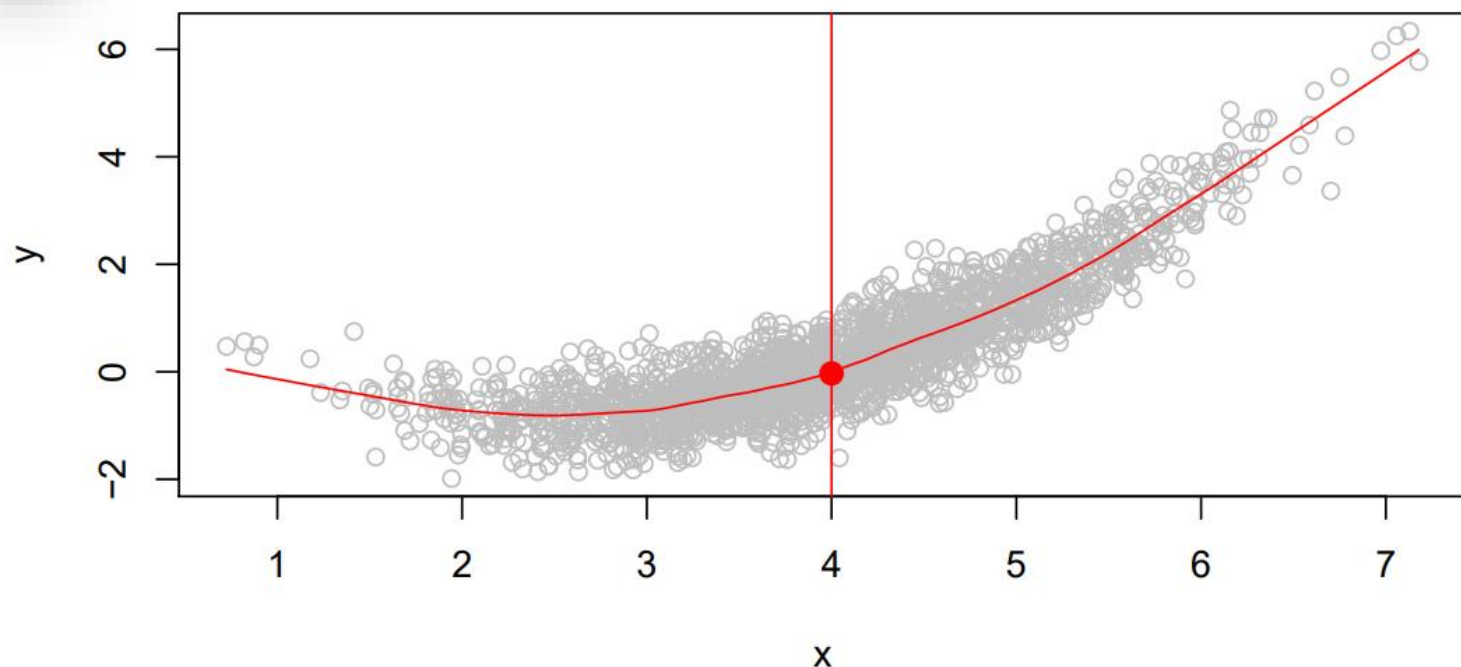
f : 函数, 固定但是未知, f 表达了X提供给Y的系统信息

估计 f 的主要原因:

- 有一个好的 f ，可以在新点 $X = x$ 上预测 Y 。
- 可以理解 $X = (X_1, X_2, \dots, X_p)$ 在解释 Y 时哪些是重要元素，而哪些是不相关的。
- 根据 f 的复杂性，能够理解 X 的每个分量 X_j 如何影响 Y 的。

预测

推断



有理想的 $f(X)$ 吗?

$f(X)$ 在任意选定的 X 值处的最佳值是多少, 比如 $X = 4$? 在 $X = 4$ 处可以有很多 Y 值。一个好的值是:

$$f(4) = E(Y | X = 4)$$

- $E(Y | X = 4)$ 表示给定 $X = 4$ 的 Y 的期望值(平均值)。
- 理想 $f(x) = E(Y | X = x)$ 称为回归函数。

对向量 X 也有定义，如：

$$f(x) = f(x_1, x_2, x_3) = E(Y | X_1 = x_1, X_2 = x_2, X_3 = x_3)$$

对于最优和最理想 Y 的均方预测误差：

$$f(x) = E(Y | X = x)$$

是使函数 g 在所有 $X=x$ 点上的 $E[(Y - g(X))^2 | X = x]$ 最小化的函数。

对于 $f(x)$ 的任何估计 $\hat{f}(x)$ 有:

$$E[(Y - \hat{f}(X))^2 | X = x] = [f(x) - \hat{f}(x)]^2 + \text{Var}(\epsilon)$$

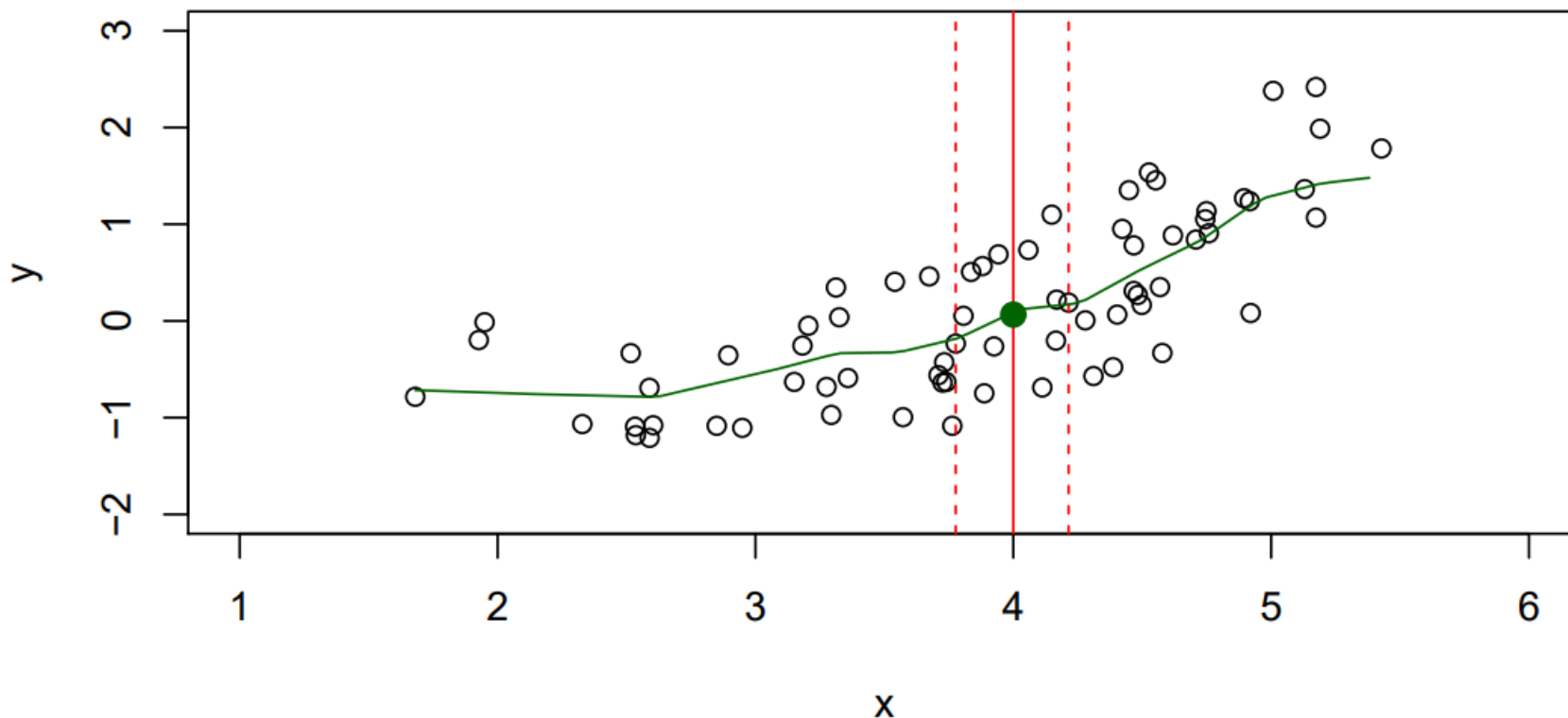
可约误差
reducible error

不可约误差
irreducible
error

对于 $f(x)$ 的任何估计 $\hat{f}(x)$ 有:

$$E[(Y - \hat{f}(X))^2 | X = x] = [f(x) - \hat{f}(x)]^2 + \text{Var}(\epsilon)$$

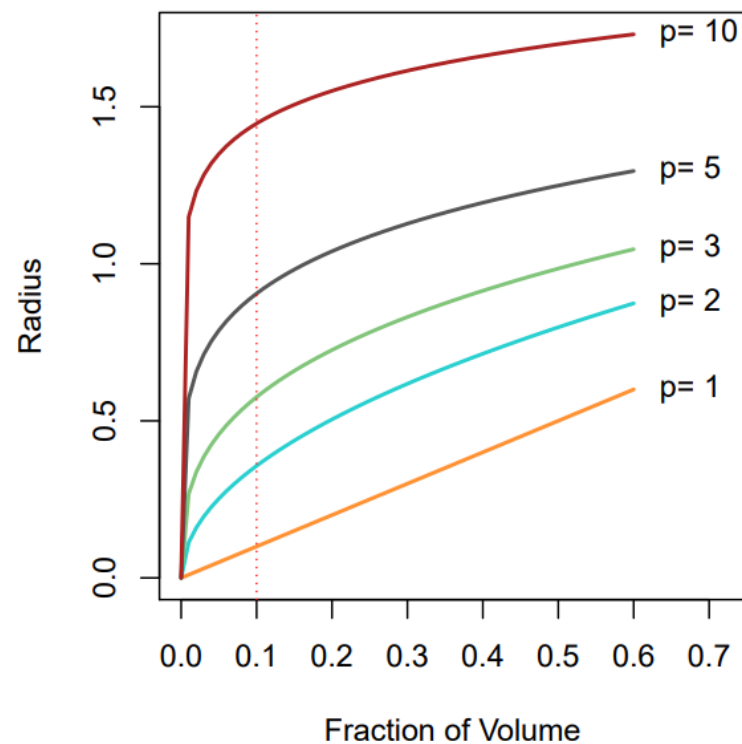
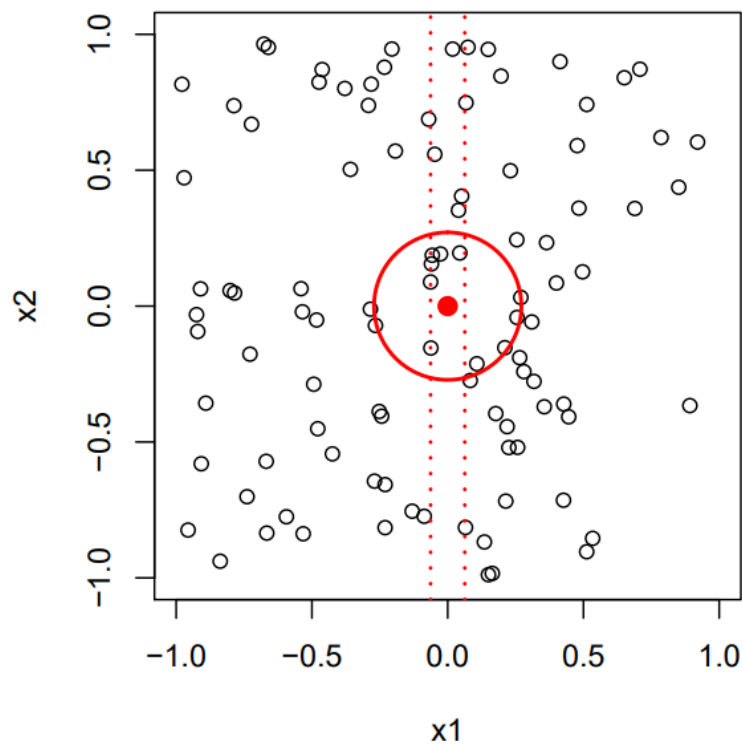
- $E(Y - \hat{Y})^2$ 代表预测量和实际值 Y 的均方误差或期望平方误差值， $\text{Var}(\epsilon)$ 表示误差项 ϵ 的方差。
- $\epsilon = Y - f(x)$ 是不可减少的误差——即使知道 $f(x)$ ，在预测时仍然会出错，因为在每个 $X = x$ 处，可能的 Y 值通常都有一个分布。



若 $X = 4$ 的数据点很少，怎么办？

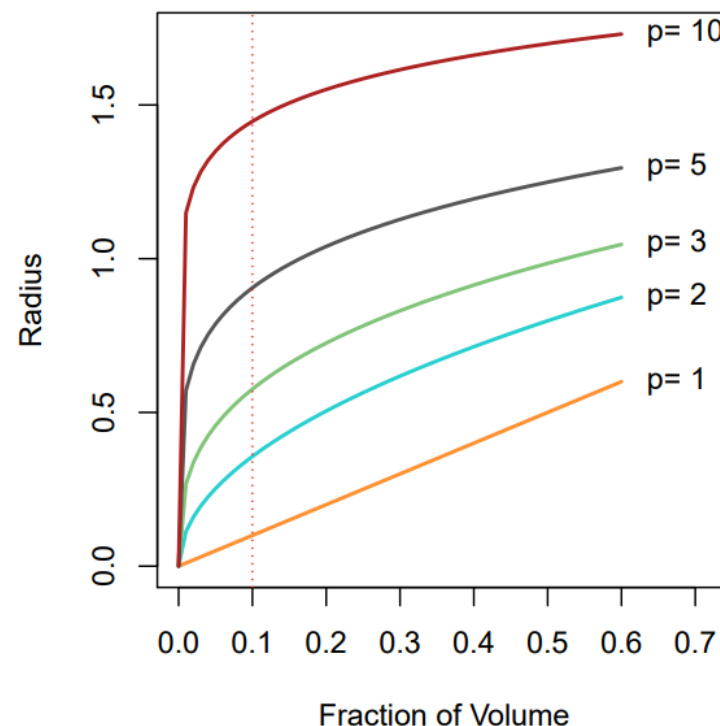
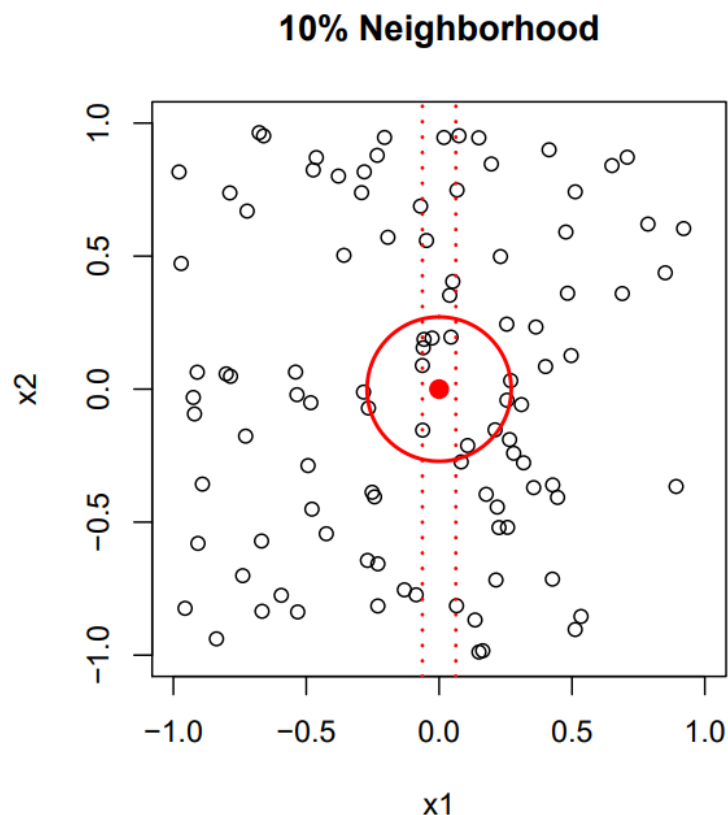
- 放宽定义，设： $\hat{f}(x) = \text{Ave}(Y|X \in \mathcal{N}(x))$
- 其中 $\mathcal{N}(x)$ 是 x 的邻域。

10% Neighborhood



- 对于小的 p ，最近邻平均是非常好的，即 $p \leq 4$ 和大的 N 。
- 当 p 很大时，最近邻方法可能很糟糕。

原因:维数灾难。



在高维空间中，最近的邻居往往离我们很远。

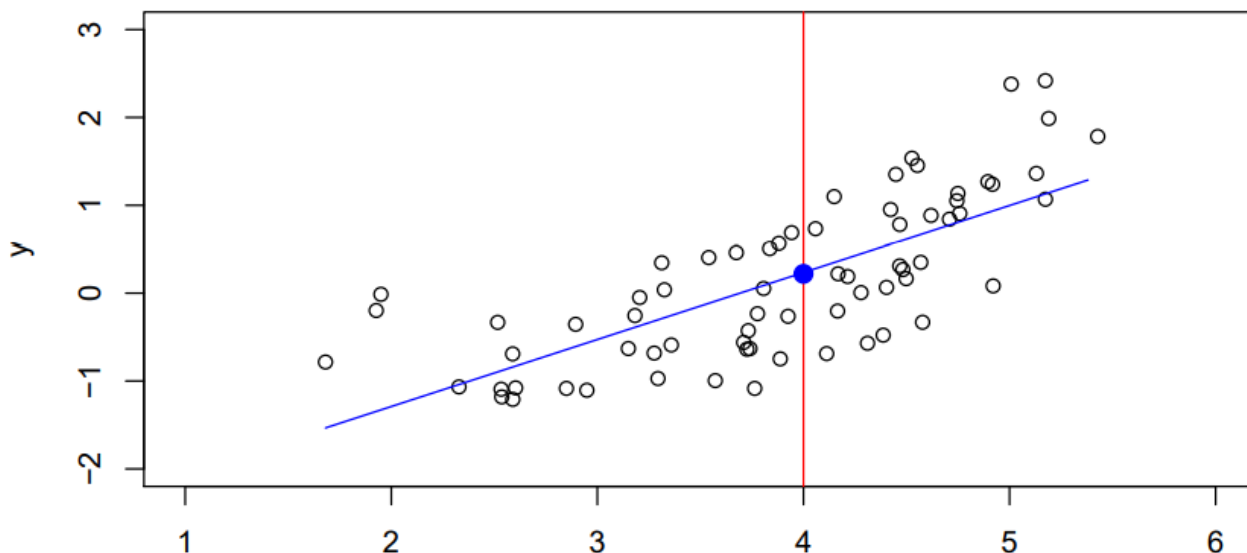
- 需要得到 y_i 的 N 值的一个合理的部分来取平均，以降低方差。
- 高维中10%的邻域不再是局部的，因此不能通过局部平均来估计 $E(Y | X = X)$ 。

线性模型是参数模型的一个重要例子：

$$f_L(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \beta_p X_p$$

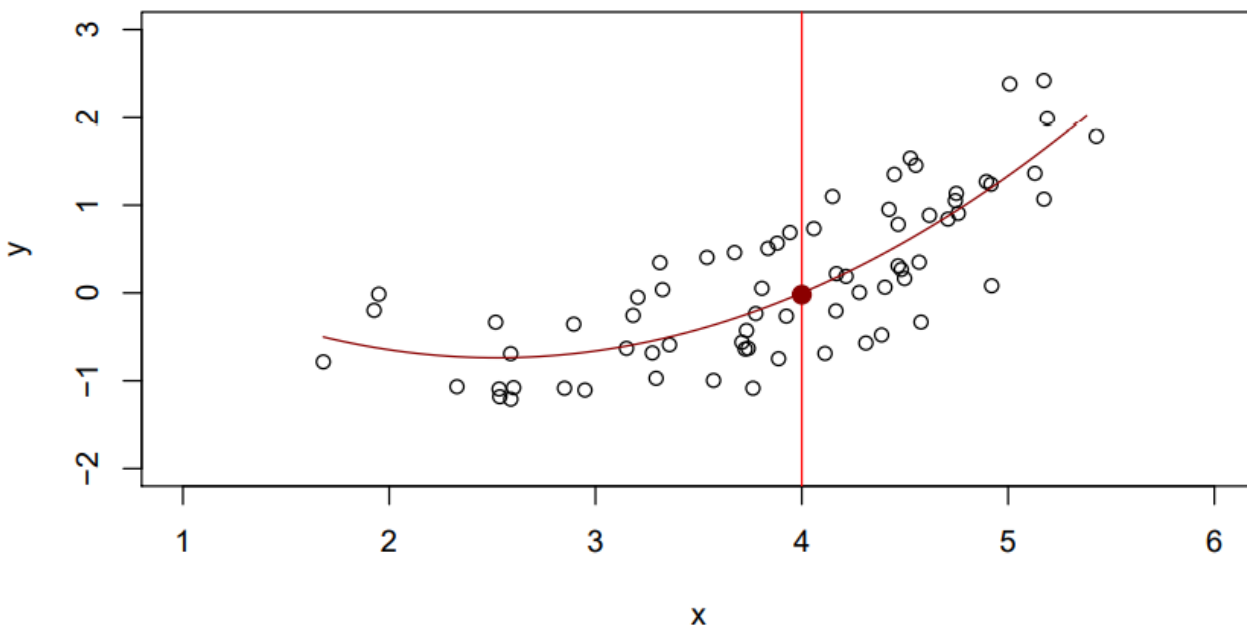
线性模型用 $p + 1$ 参数表示 $\beta_0, \beta_1, \dots, \beta_p$

- 参数方法：估计参数，采用训练数据来拟合模型。
- 虽然线性模型常常都不是正确的，但它通常是未知真函数 $f(X)$ 的一个良好的、可解释的逼近。



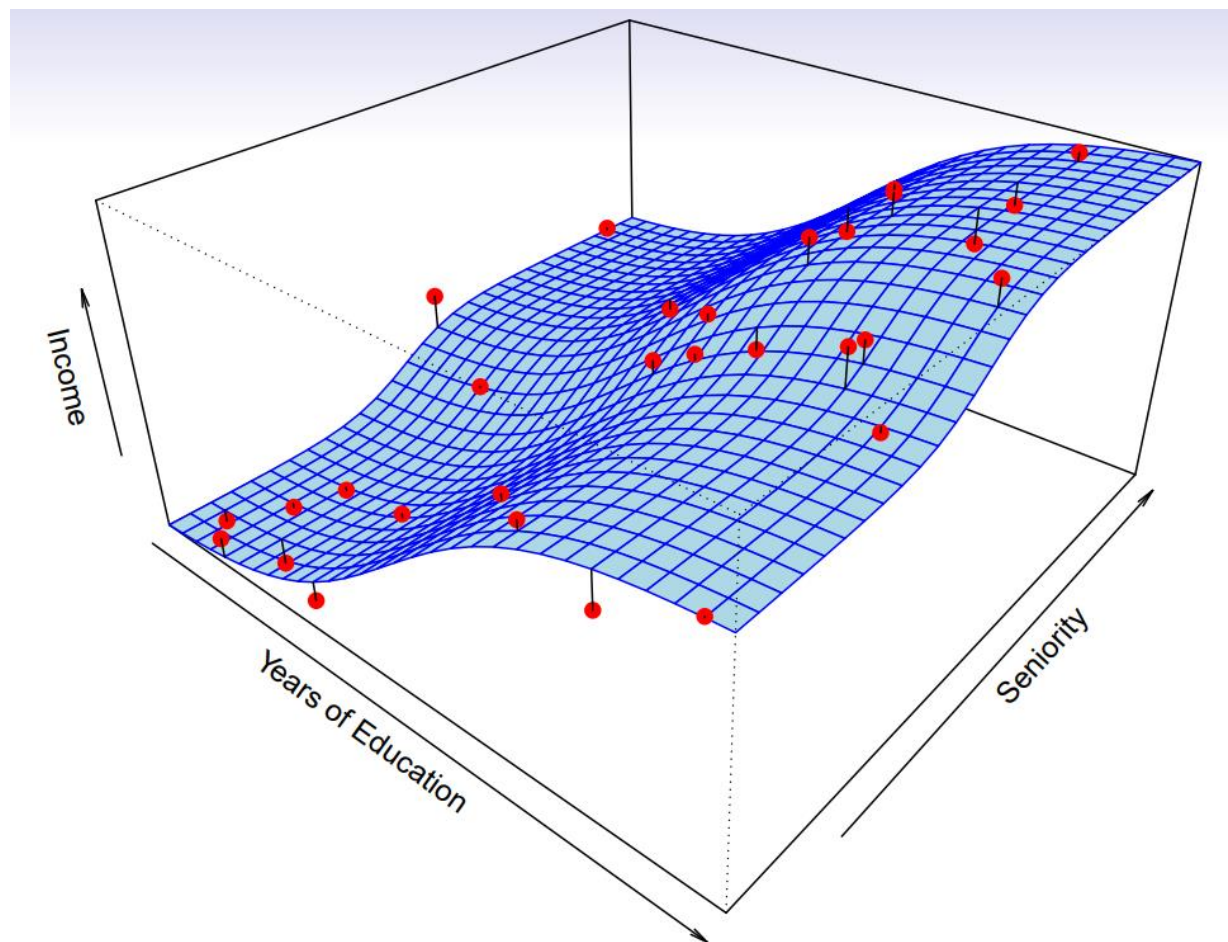
线性模型

$$\hat{f}_L(X) = \hat{\beta}_0 + \hat{\beta}_1 X$$



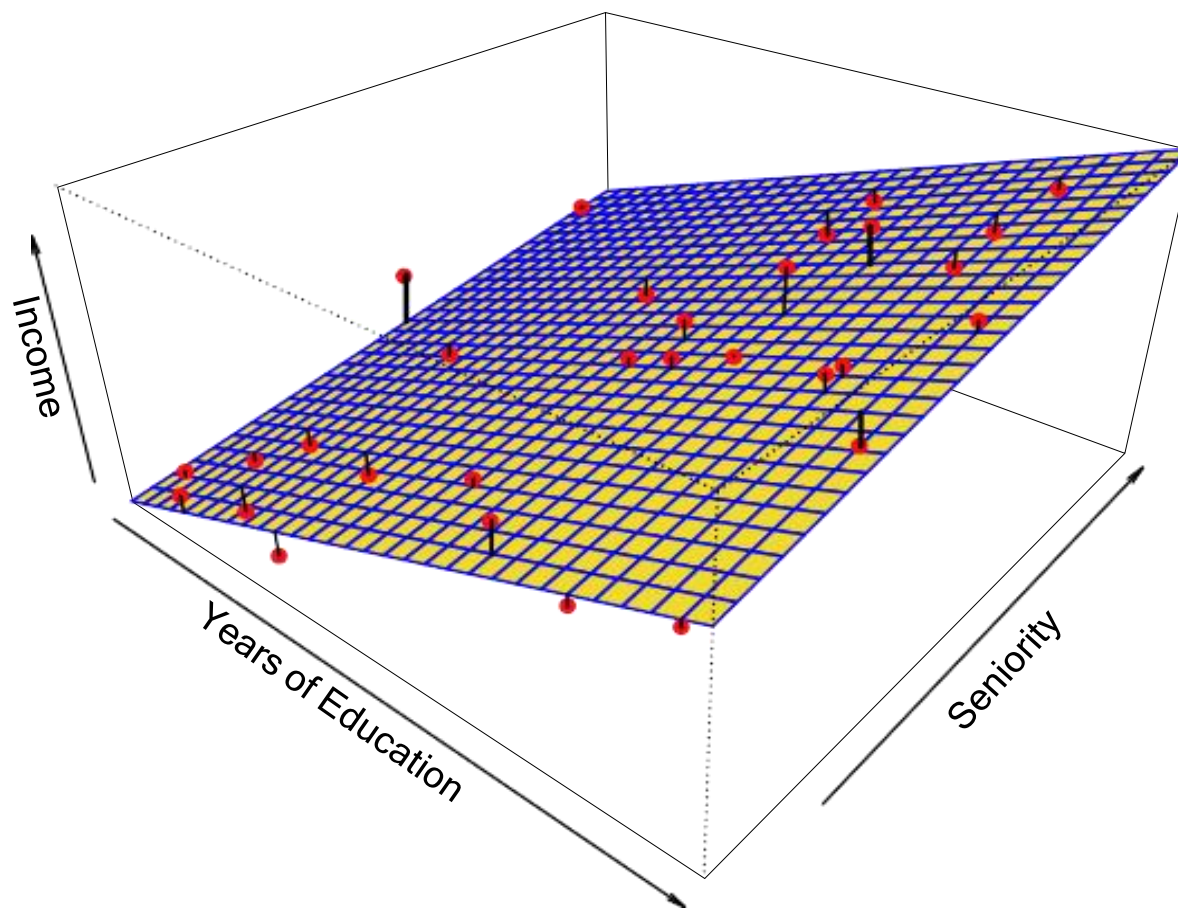
二次方程

$$\hat{f}_Q(X) = \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2$$



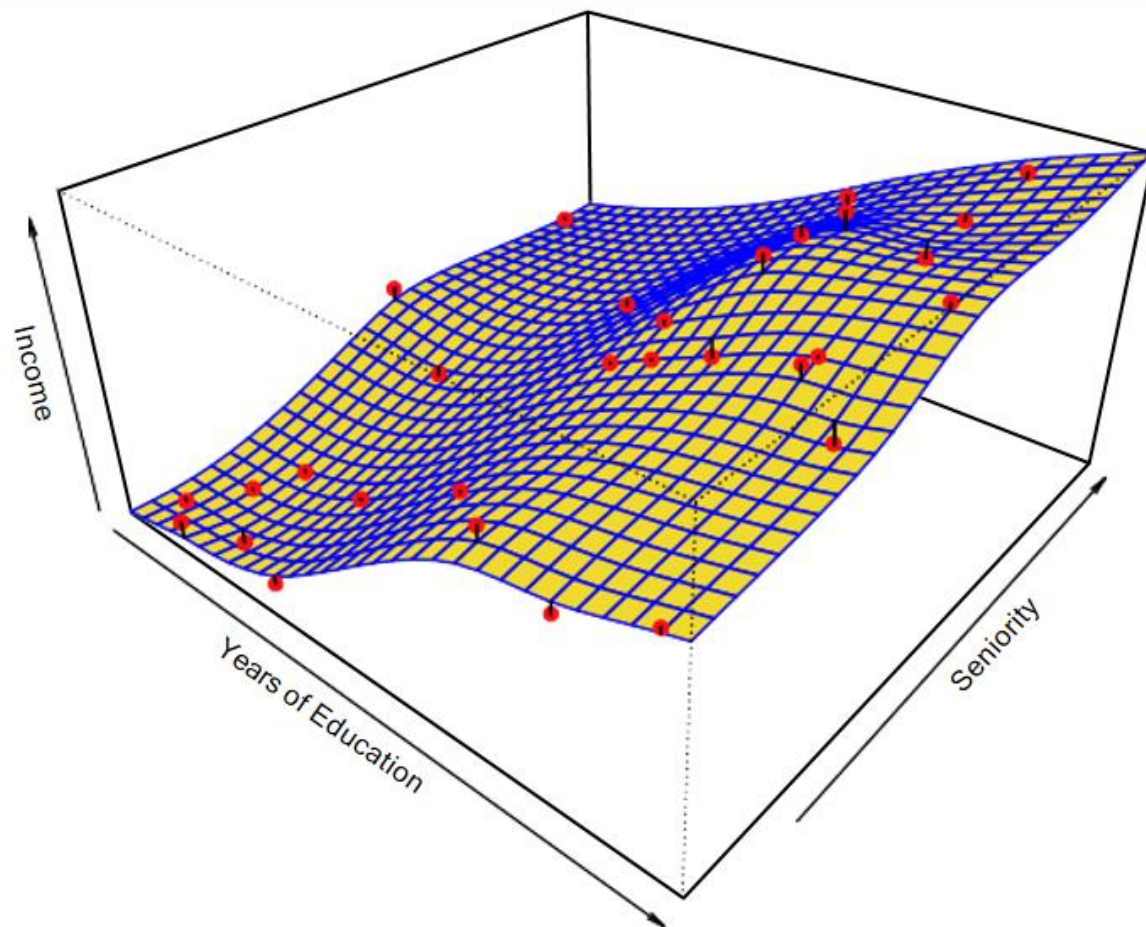
红点是模型所得的模拟值， f 是蓝色的曲面

$$\text{income} = f(\text{education}, \text{seniority}) + \epsilon$$



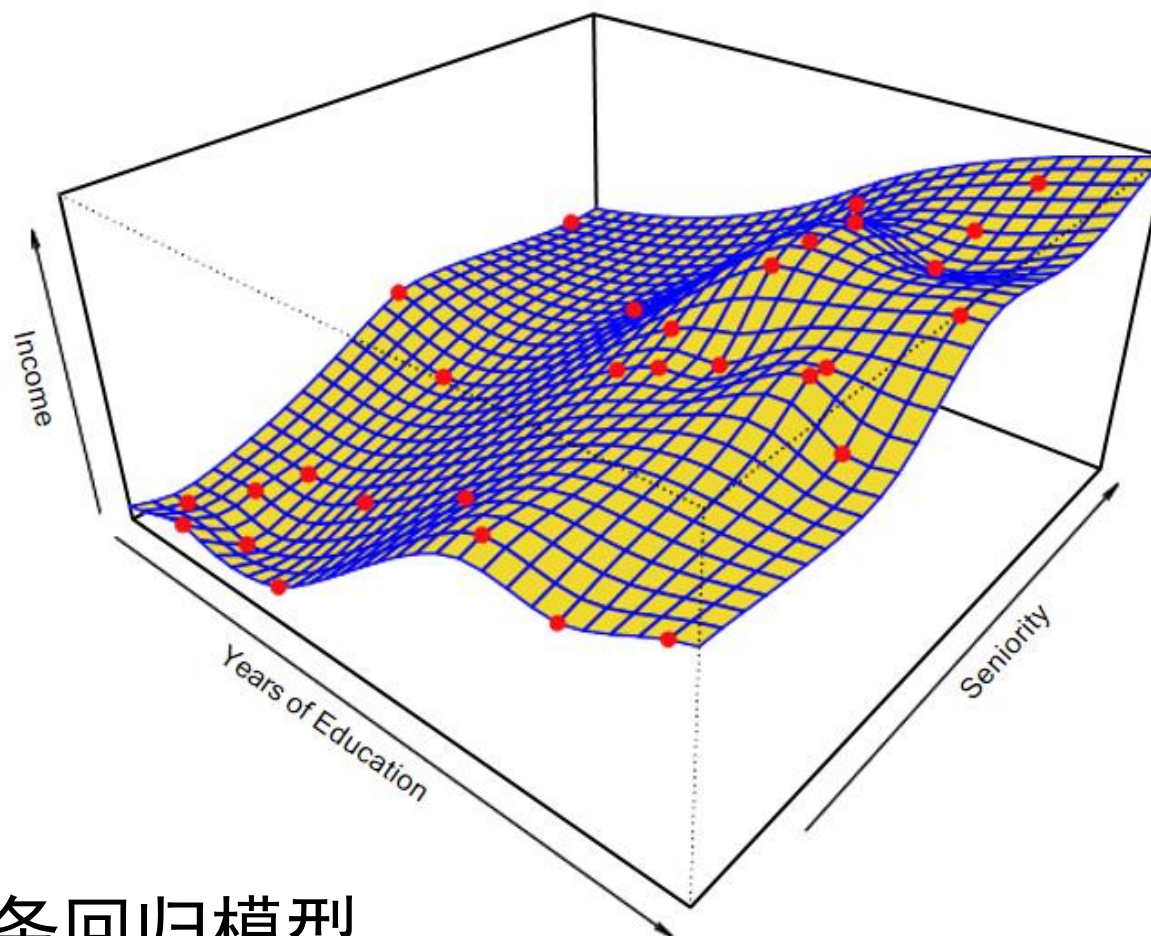
进行最小二乘线性回归模型拟合

$$\hat{f}_L(\text{education}, \text{seniority}) = \hat{\beta}_0 + \hat{\beta}_1 \times \text{education} + \hat{\beta}_2 \times \text{seniority} \quad 17$$



使用光滑薄板样条的技术拟合

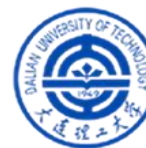
$$\hat{f}_S(\text{education}, \text{seniority})$$



更灵活的样条回归模型

$$\hat{f}_S(\text{education}, \text{seniority})$$

这里拟合模型在训练数据上没有误差!也被称为过拟合。



- 预测精度VS可解释性。

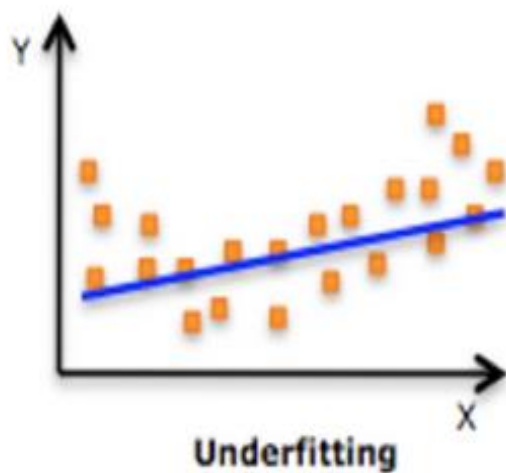
线性模型易于解释;薄板样条则不然。

- 好的拟合VS欠拟合与过拟合。

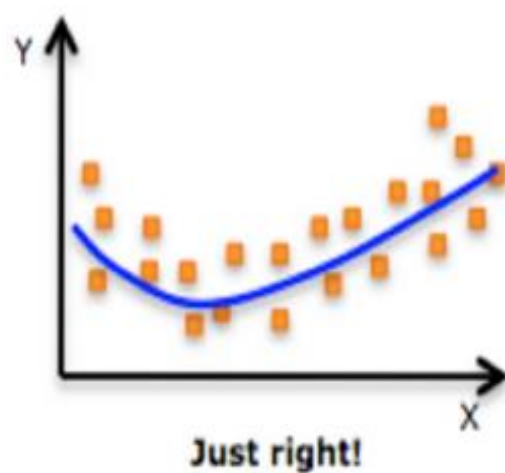
怎么知道什么时候合适?

- 简单VS黑盒 (black-box) 。

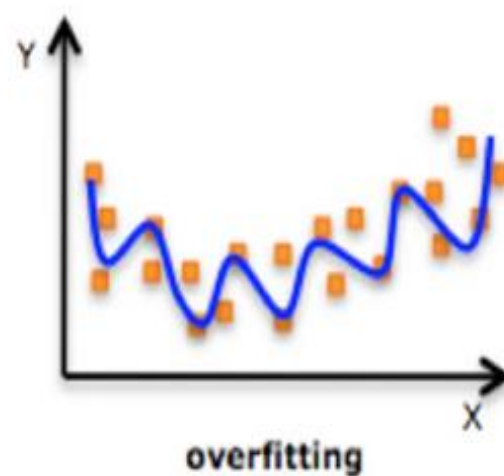
通常更喜欢包含更少变量的简单模型，而不是包含所有变量的黑箱预测器。



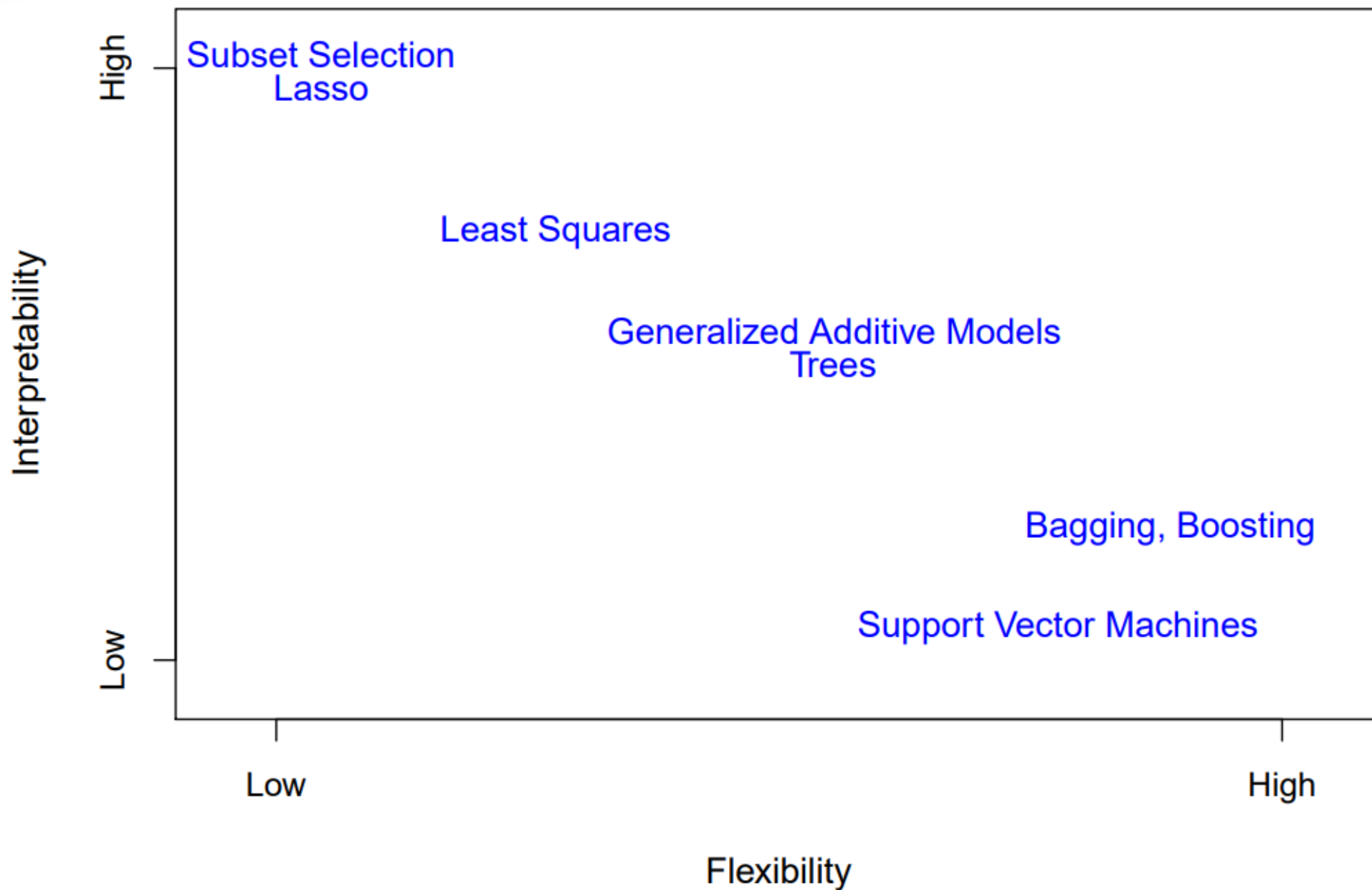
$$\theta_0 + \theta_1 x$$



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

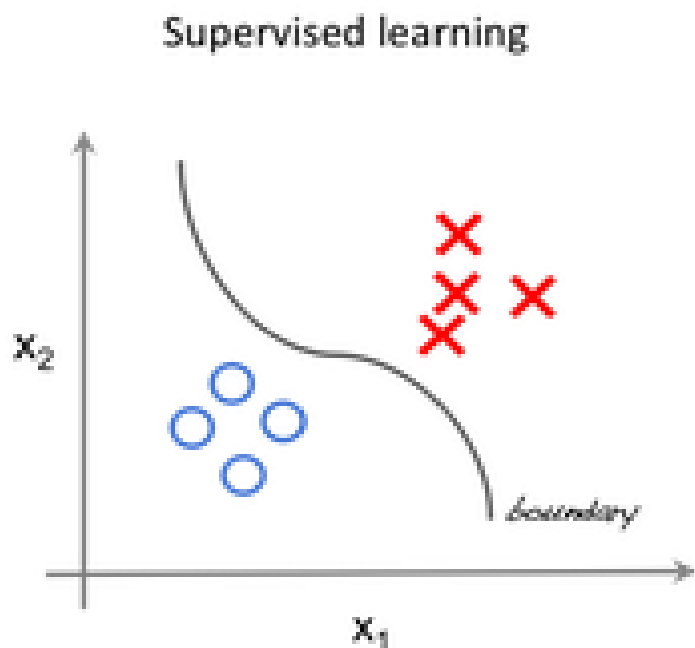


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$



指导学习 (supervised learning)

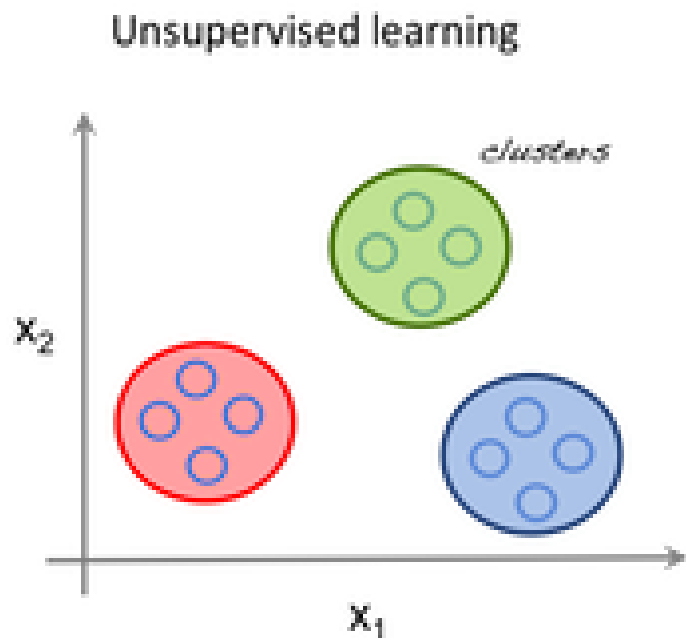
对每一个预测变量观测值 $x_i (i = 1, \dots, n)$ 都有相应的响应变量的观测 y_i 。



- 线性回归
- 逻辑斯谛回归
- 广义可加模型 (GAM)
- 支持向量机 (SVM)
-

无指导学习 (unsupervised learning)

只有预测变量的观测值 $x_i (i = 1, \dots, n)$ ，这些向量没有相应的响应变量 y_i 与之对应。



典型统计学习工具：
聚类分析
(cluster analysis)

半指导学习 (semi supervised learning)

- 假设有 n 个观测，其中 m ($m < n$) 个观测点可同时观测到预测变量与响应变量，而其余 $n-m$ 个观测点，只能观测到预测变量但无法观测到响应变量
- 这类统计方法就希望能够既用到 m 个观测点的预测变量和响应变量的信息，同时又包含了 $n-m$ 个不能获取响应变量观测值的信息。

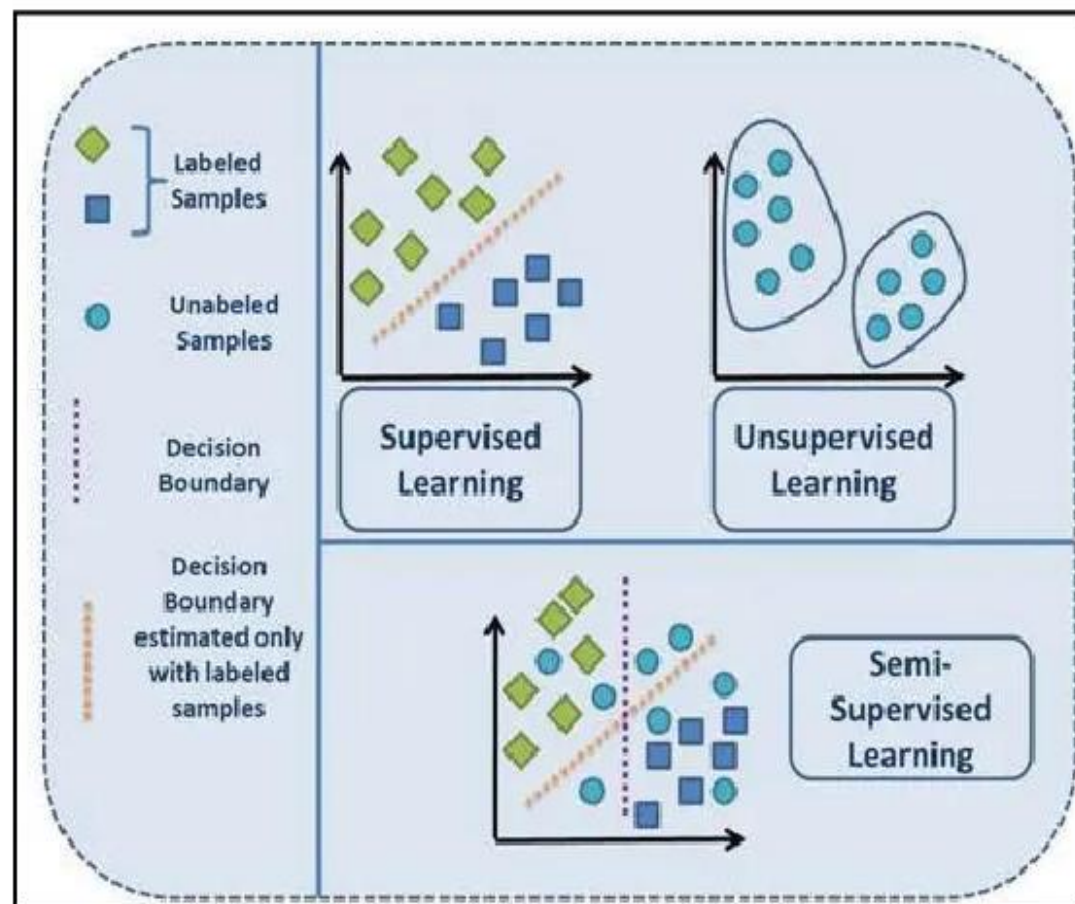


Fig. 1 Semi-Supervised Learning

变量

```
graph LR; A[变量] --> B[定量]; A --> C["定性 (分类变量)"]
```

定量

呈**数值型**，如：年龄身高收入房屋的价值以及股票的价格

定性
(分类变量)

取K个不同类的其中一个值，如：性别、品牌等

回归问题分析问题：将响应变量为定量的问题。

分类问题：具有定性响应变量的问题

回归与分类问题并不绝对！



2. 评价模型精度

在回归中最常用的评价准则：

均方误差 (mean squared error, MSE)

假设我们对训练数据 $\text{Tr} = \{x_i, y_i\}_1^N$ 拟合一个模型 $\hat{f}(x)$:

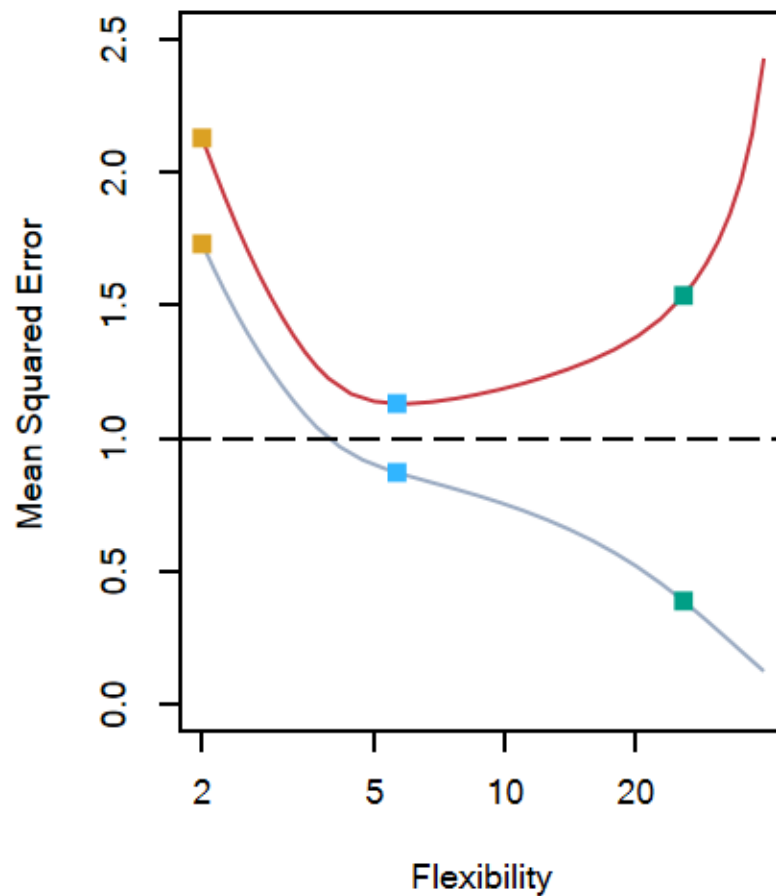
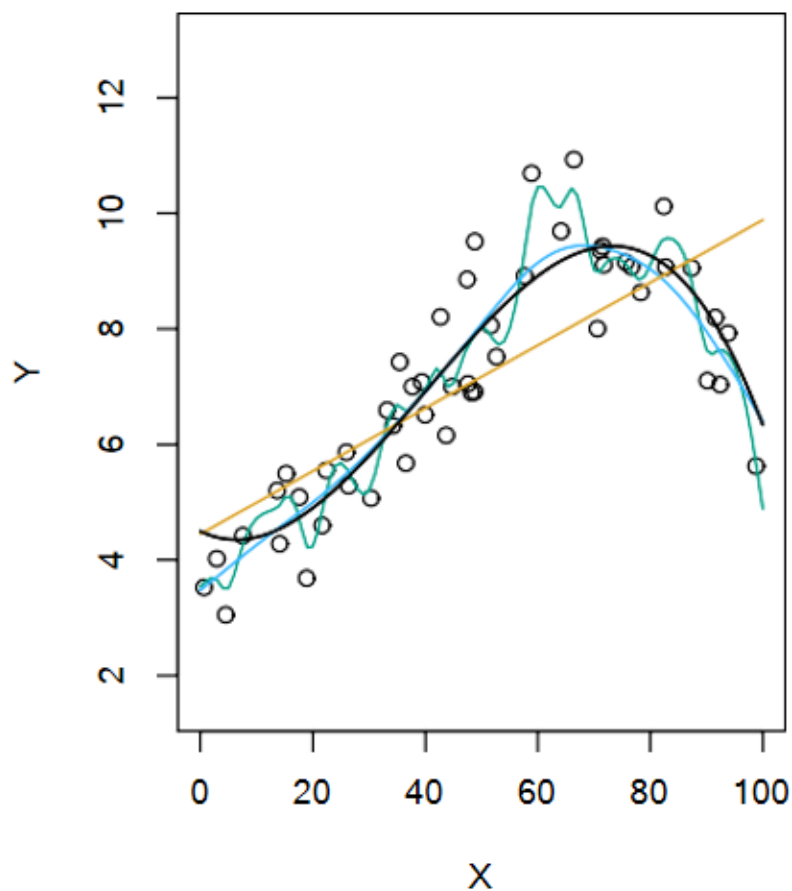
我们可以计算Tr上的训练均方误差:

$$\text{MSE}_{\text{Tr}} = \text{Ave}_{i \in \text{Tr}} [y_i - \hat{f}(x_i)]^2$$

这可能会偏向于更过拟合的模型。

使用新的测试数据 $\text{Te} = \{x_i, y_i\}_1^M$ 计算:

$$\text{MSE}_{\text{Te}} = \text{Ave}_{i \in \text{Te}} [y_i - \hat{f}(x_i)]^2$$

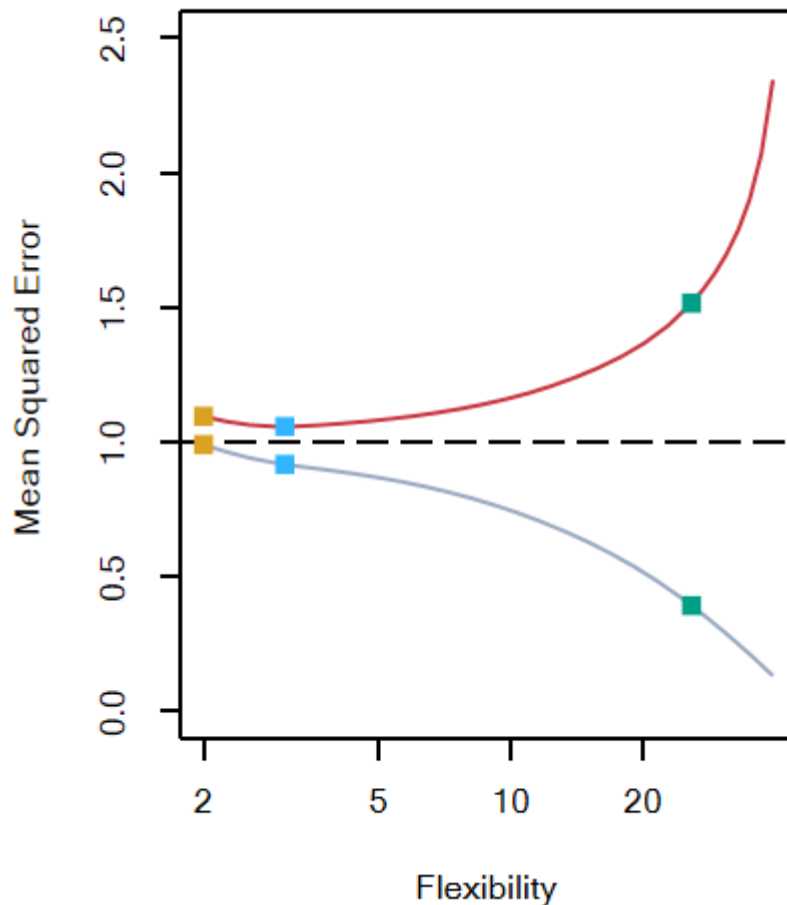
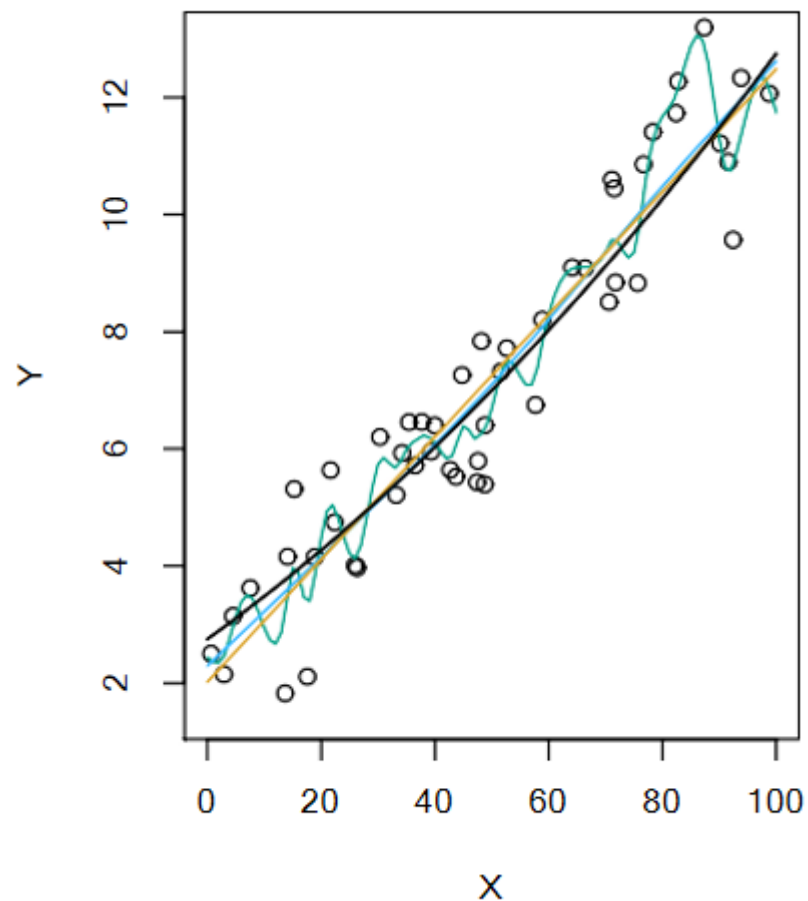


三种 f 的估计:

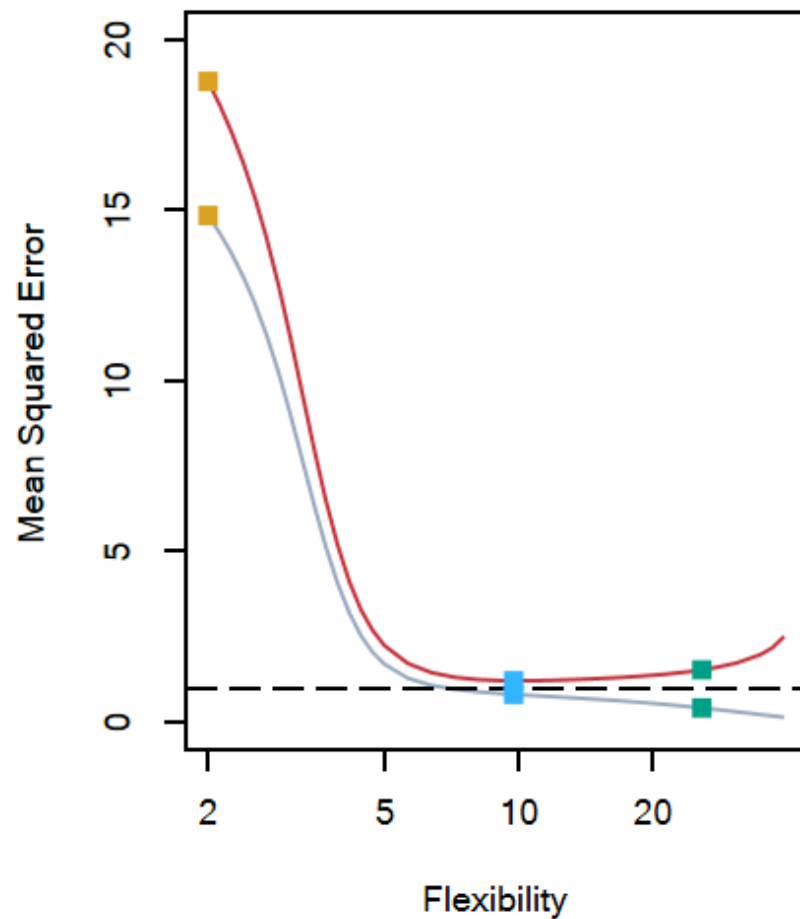
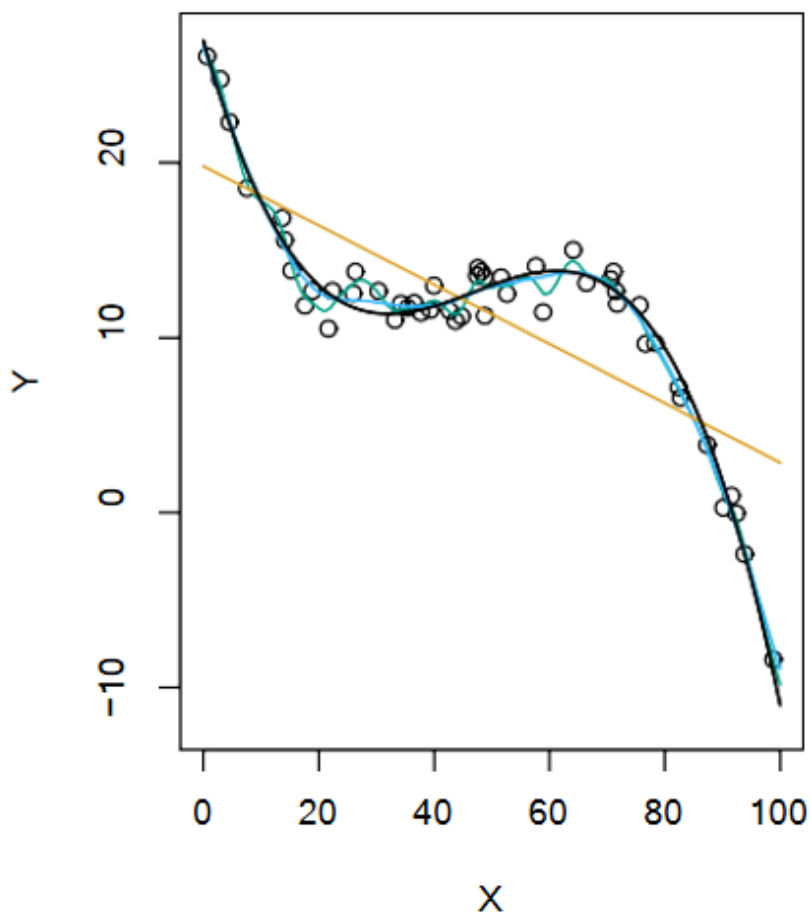
- 线性回归
(橙色曲线)
- 两条光滑样条拟合
(绿色和蓝色曲线)

左: 真实函数 f 模拟产生的数据, 黑色曲线表示。

右: 训练均方误差 (灰色曲线), 测试均方误差 (红色曲线), 所有方法都已使测试均方误差尽可能最小。



这里真实的 f 平滑，
所以更平滑的拟合和线性模型拟合的较好。



这里的真实的 f 是摇摆的和低噪声，
所以更灵活的拟合做得更好。

假设我们对某些训练数据 Tr 拟合一个模型 $\hat{f}(x)$ ，并让 (x_0, y_0) 是从总体中提取的一个测试观察值。如果真正的模型是：

$$Y = f(X) + \epsilon \text{ (with } f(x) = E(\bar{Y} | X = x))$$

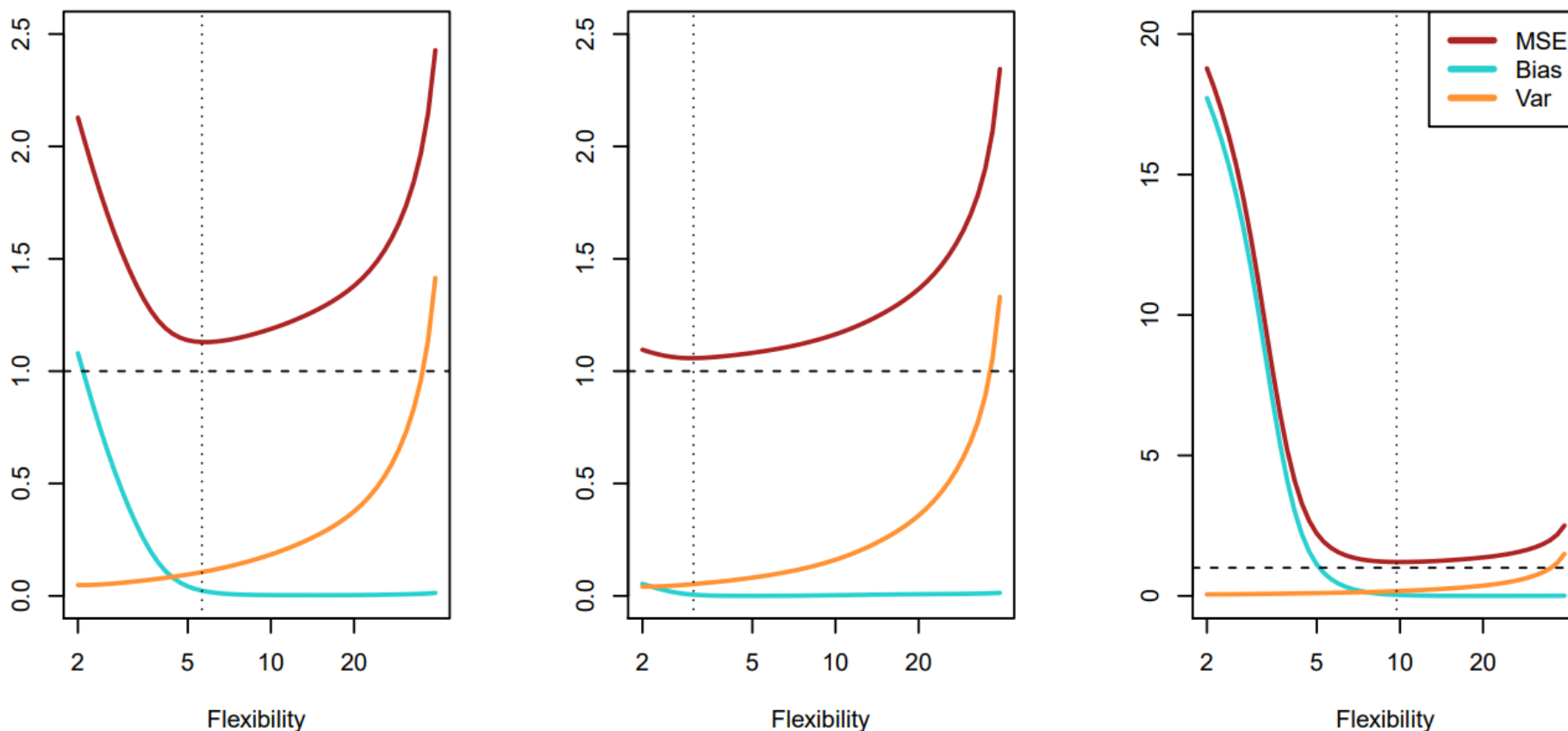
期望测试均方误差：

$$E \left(y_0 - \hat{f}(x_0) \right)^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

即偏差（bias）为：

$$\text{Bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - f(x_0)$$

通常，随着 $\hat{f}(x)$ 灵活性的增加，其方差增加，偏差减少。因此，基于平均测试误差来选择灵活性就相当于一种偏差-方差权衡。

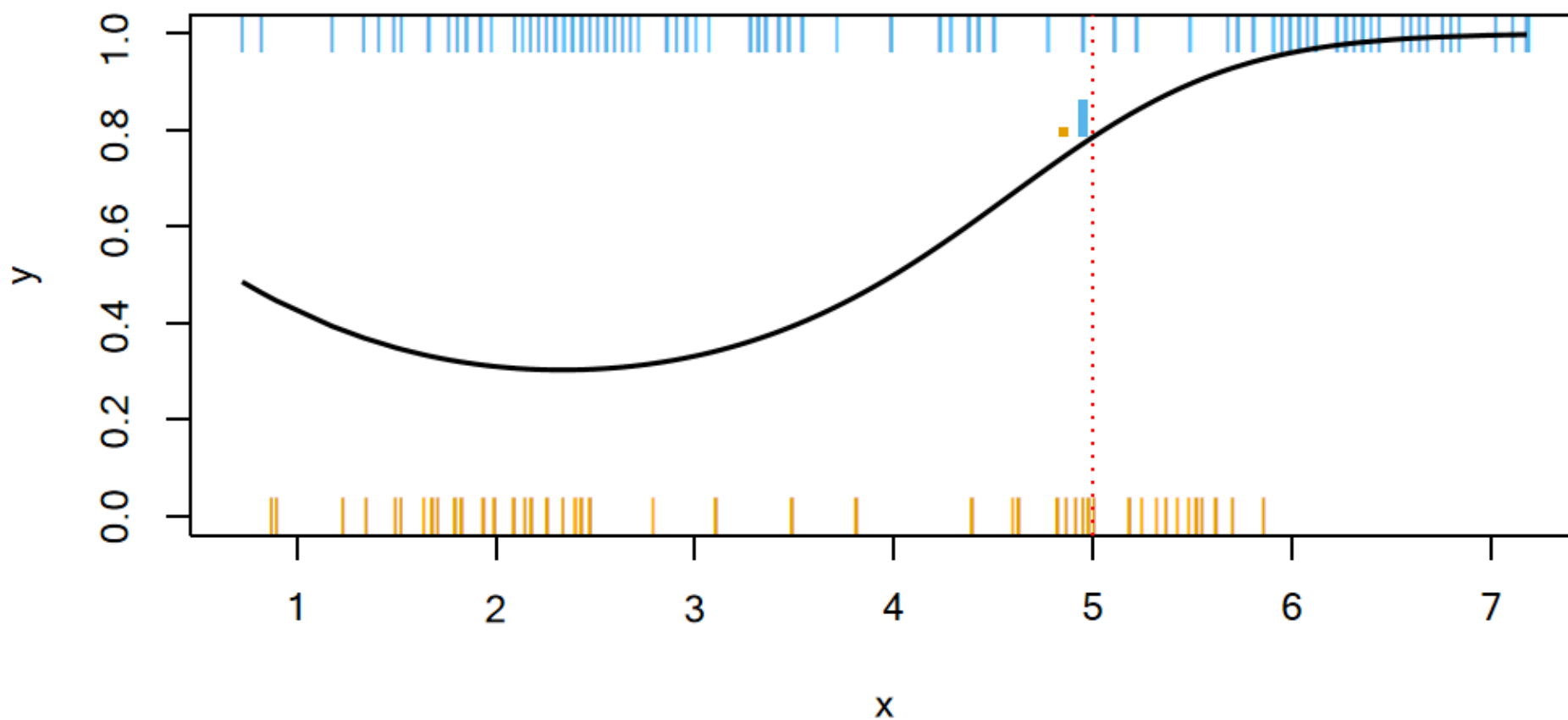


分别表示上图三个数据集的平方偏差（蓝色曲线）、方差（橙色曲线）、不可约误差（虚线）、测试均方误差（红色曲线）。垂直的点线表示最小测试均方误差所对应的光滑度。

分类模型中的响应变量 Y 是定性的，例如，电子邮件是 $C = (\text{spam}, \text{ham})$ (ham =好的电子邮件)中之一，数字类别是 $C = \{0, 1, \dots, 9\}$ 。目标是：

- 构建一个分类器 $C(X)$ ，将 C 的类标签分配给未来未标记的观察 X 。
- 评估每种分类模型的不确定性。
- 了解不同预测因素之间的作用 $X = (X_1, X_2, \dots, X_p)$ 。

有理想的 $C(X)$ 吗?



贝叶斯分类器

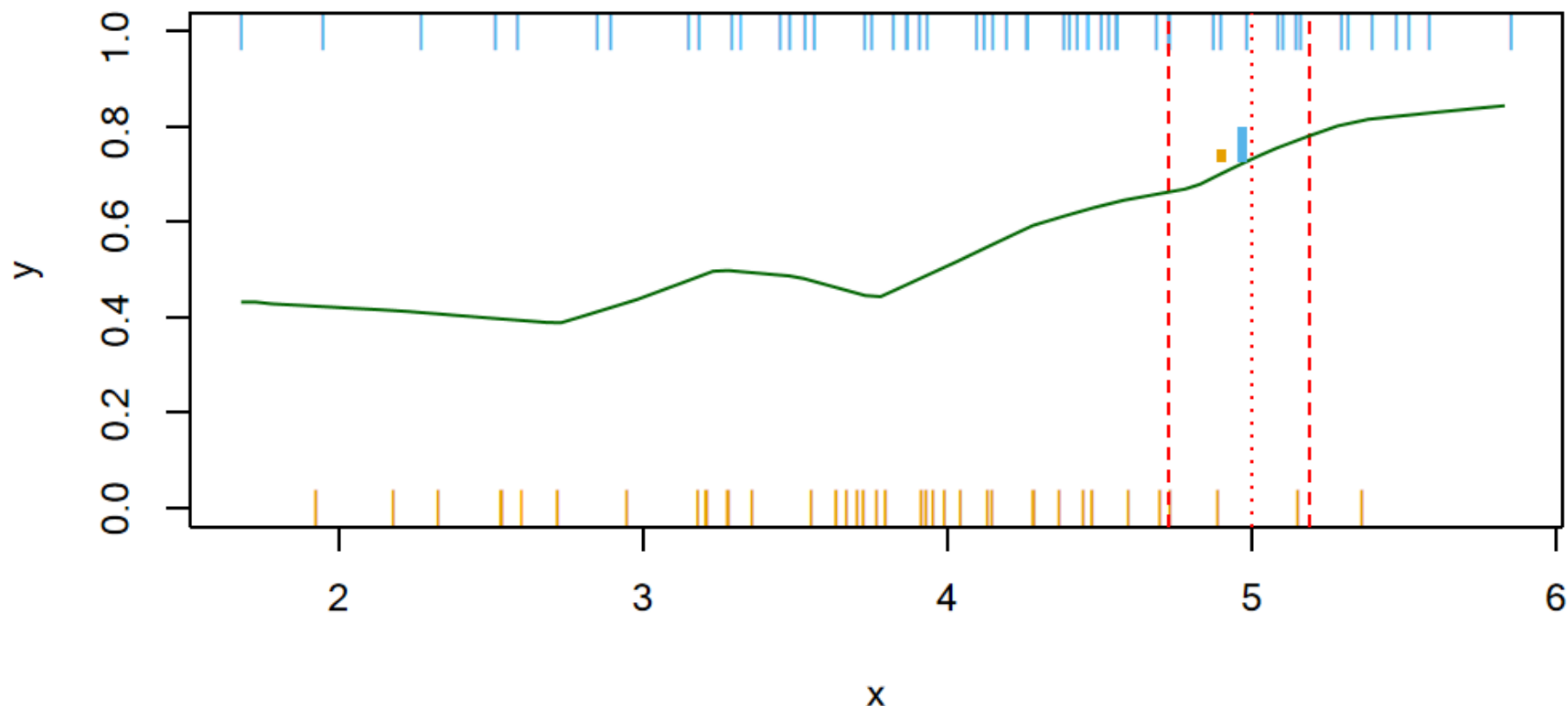
假设C中的K个元素编号为1,2, ..., k .让

$$p_k(x) = \Pr(Y = k|X = x), \quad k = 1, 2, \dots, K$$

这些是在x点的条件概率。例如，见 $x = 5$ 处的条形图。那么x处的**贝叶斯分类器**为：

$$C(x) = j \text{ if } p_j(x) = \max\{p_1(x), p_2(x), \dots, p_K(x)\}$$

k最近邻



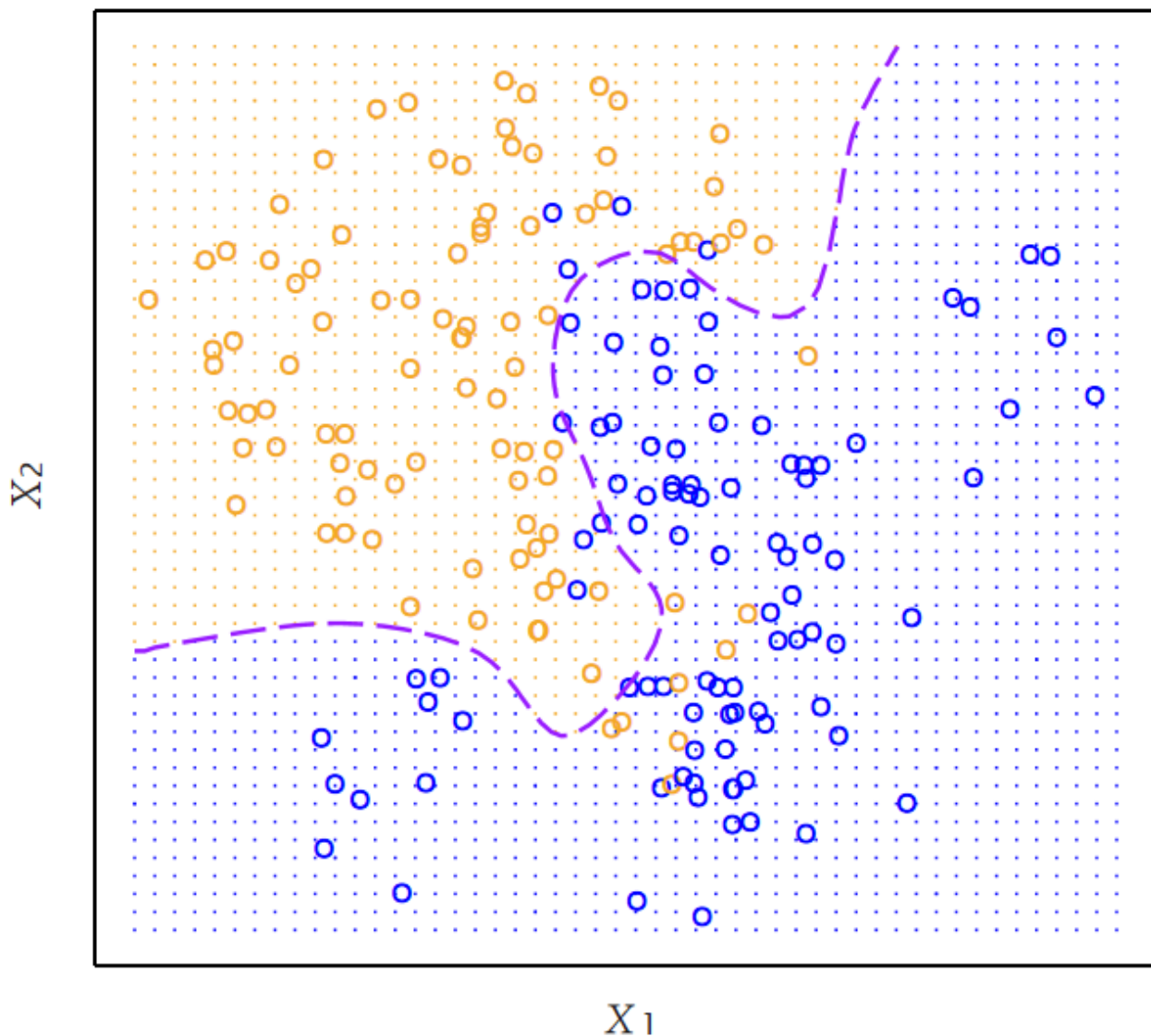
最近邻平均法可与之前一样使用，也会随着维度的增长而分解。但对 $\hat{C}(x)$ 的影响小于对 $\hat{p}_k(x)$, $k = 1, \dots, K$ 。

$$C(x) = j \text{ if } p_j(x) = \max\{p_1(x), p_2(x), \dots, p_K(x)\}$$

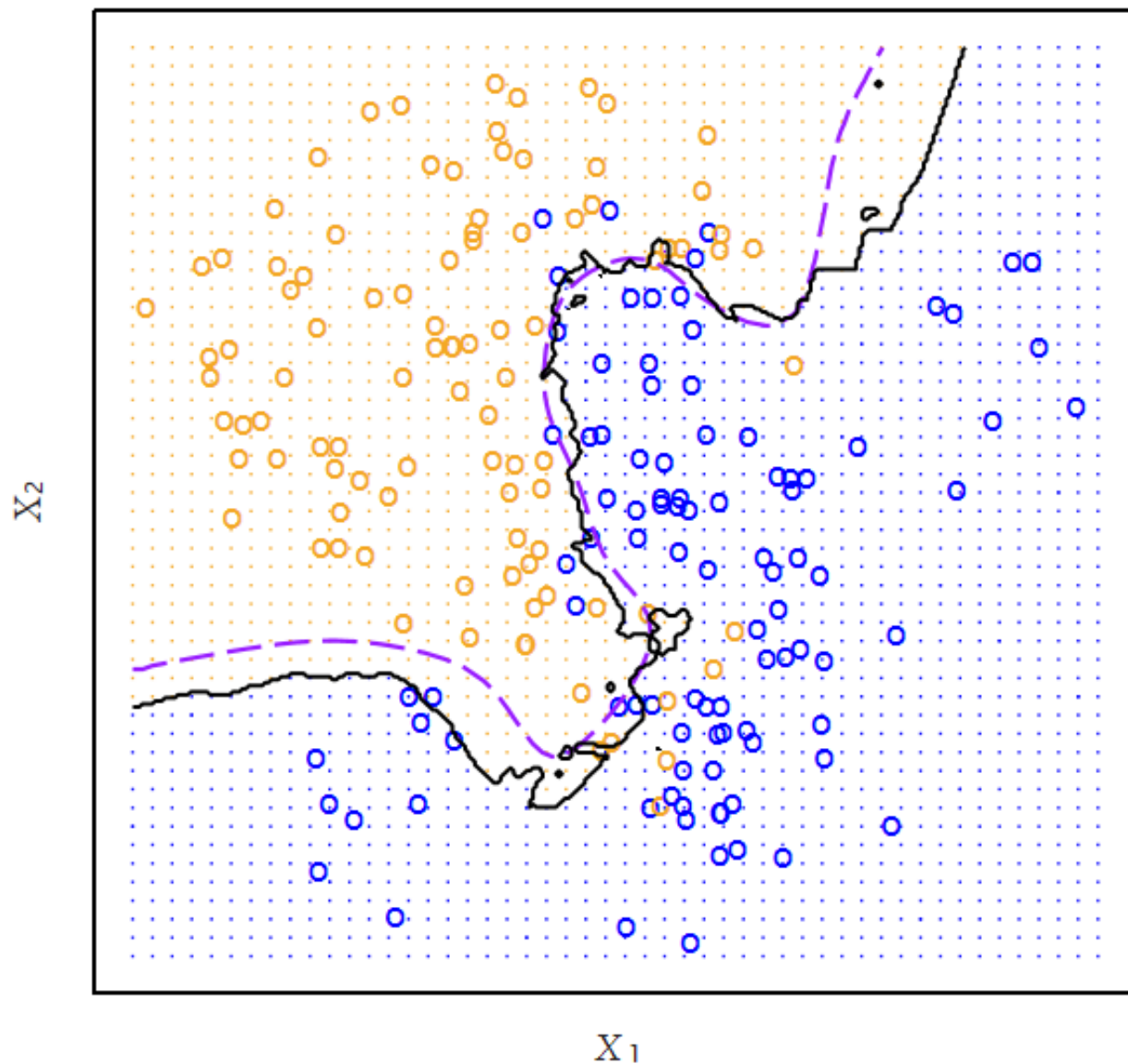
- 通常我们用误分类错误率来衡量 $\hat{C}(x)$ 的性能:

$$\text{Err}_{\text{T}_e} = \text{Ave}_{i \in \text{T}_e} I[y_i \neq \hat{C}(x_i)]$$

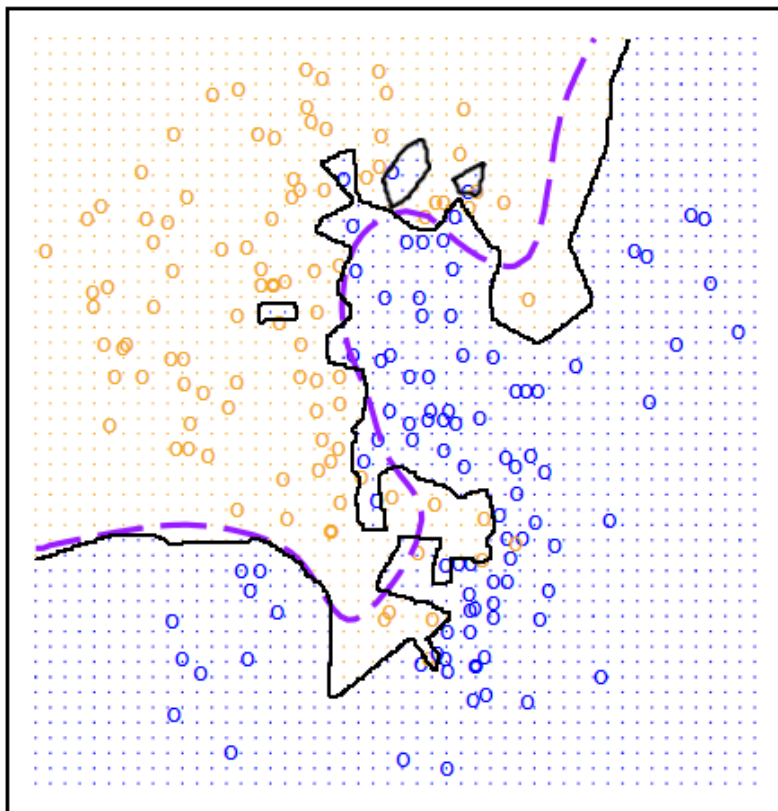
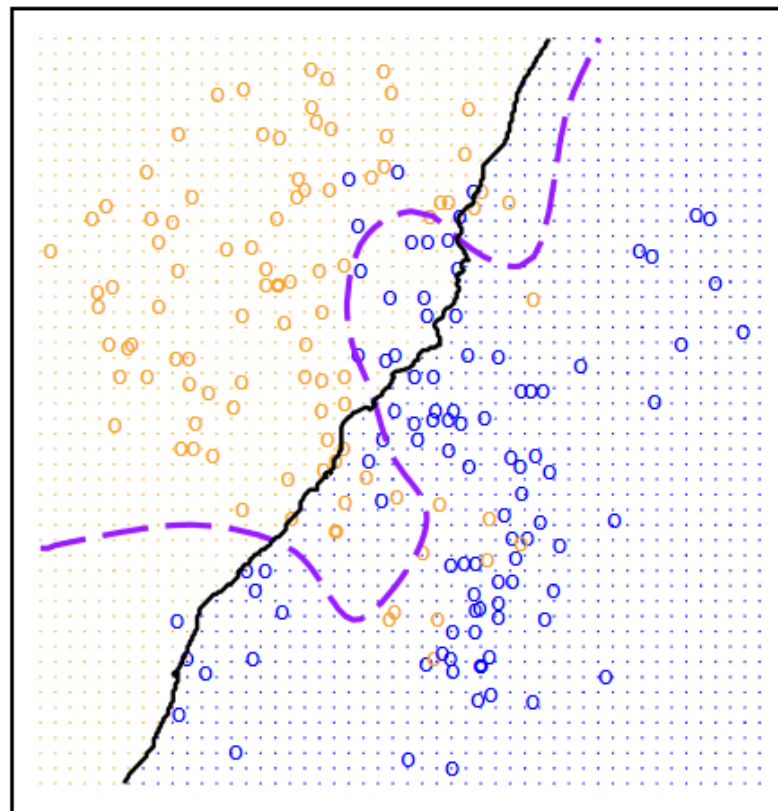
- 贝叶斯分类器(使用真实的 $p_k(x)$)的误差最小(在总体中)。
- 支持向量机为 $C(x)$ 构建结构化模型。
- 我们还将构建结构化模型来表示 $p_k(x)$ 。例如逻辑回归，广义可加性模型。



- 从二分类数据中的每一类中抽取100个观测值组成一个模拟数据集。
- 观测点由蓝色和橙色表示。紫色的虚线表示贝叶斯决策边界。

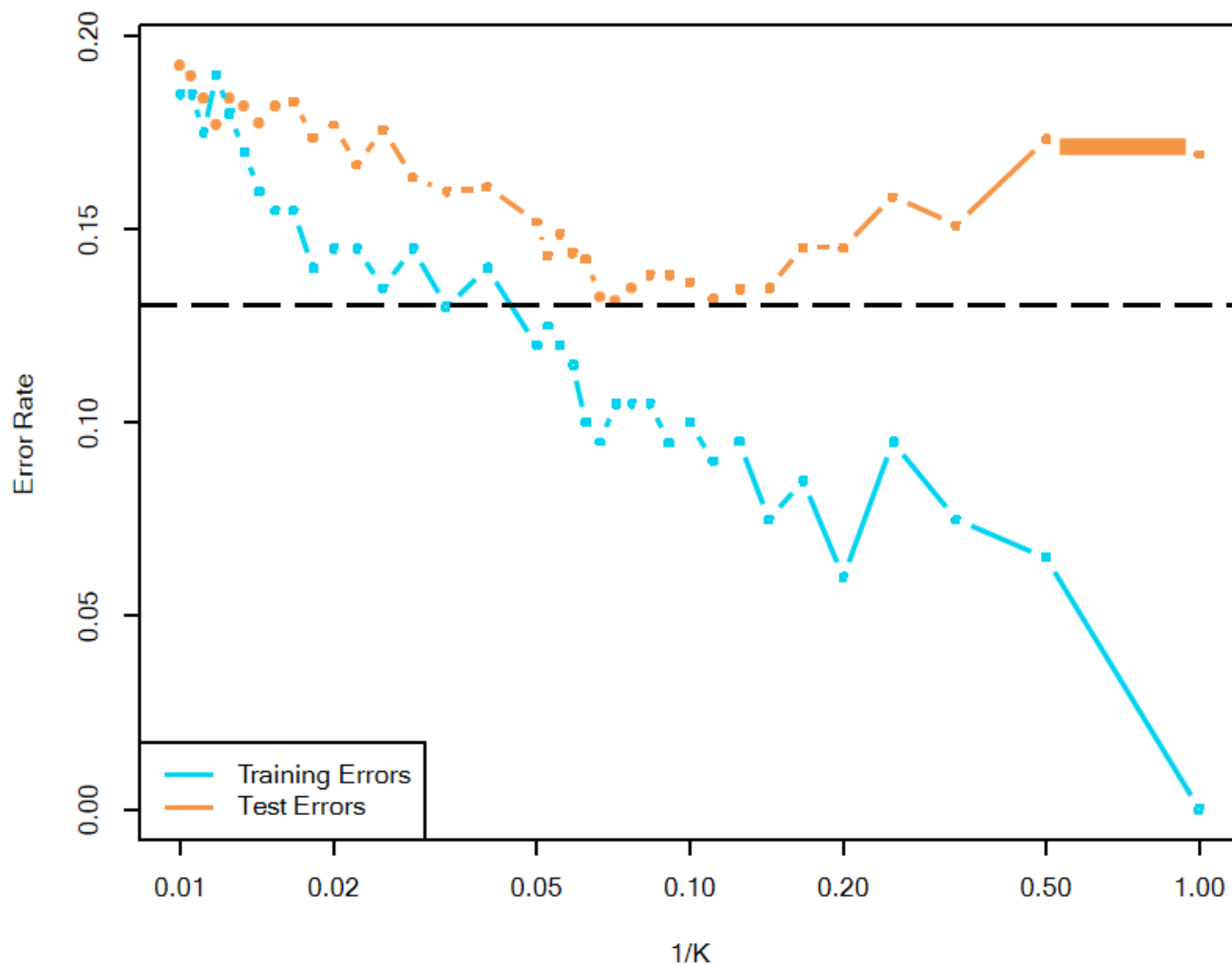
KNN: $K=10$ 

黑色曲线表示KNN方法用于上图数据生成的决策边界，这里 $k=10$.

KNN: $K=1$ KNN: $K=100$ 

在 $K=2$ 和 $K=100$ 两种设置下KNN决策边界的比较：

- 当 $K=1$ 时，决策边界相当不规则；
- 当 $K=100$ 时，模型光滑度下降。紫色虚线为贝叶斯决策边界。



用上述数据所生成的KNN分类器训练错误率（蓝色，200观测点），和测试错误率（橙色，5000观测点）对模型光滑度的变化曲线。