

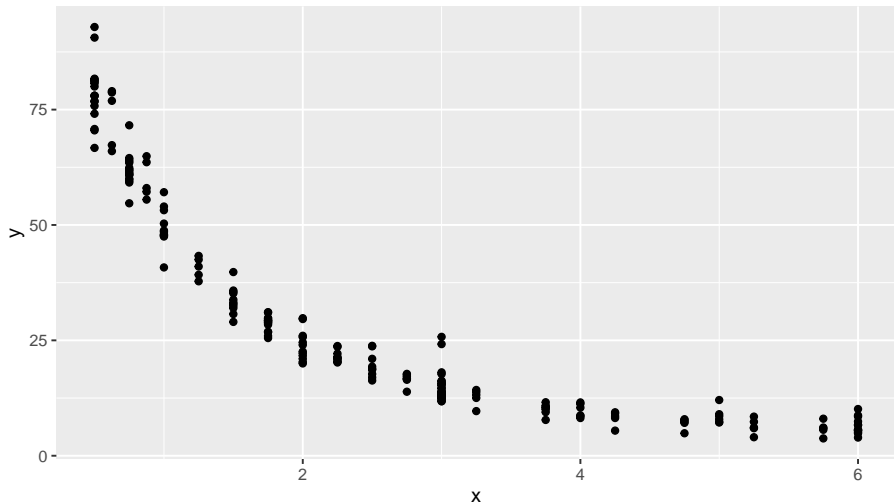
Regression Analysis

Chapter 11: Nonlinear Regression

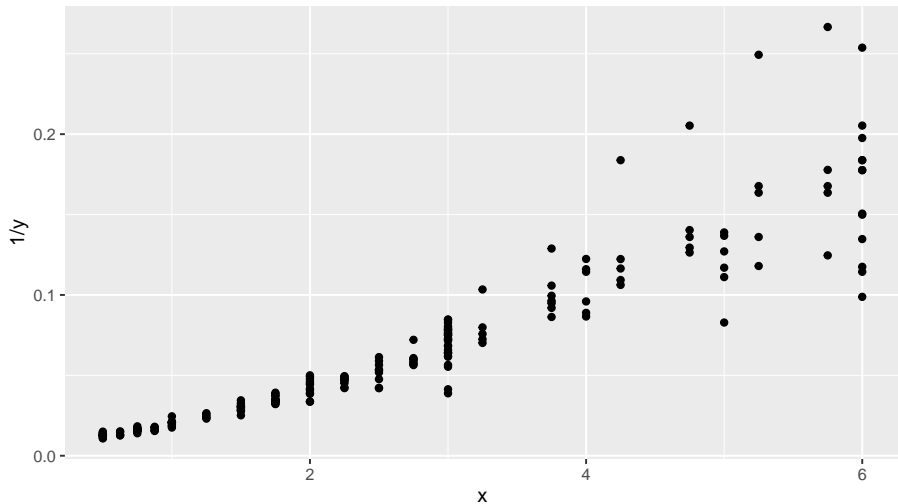
Shaobo Jin

Department of Mathematics

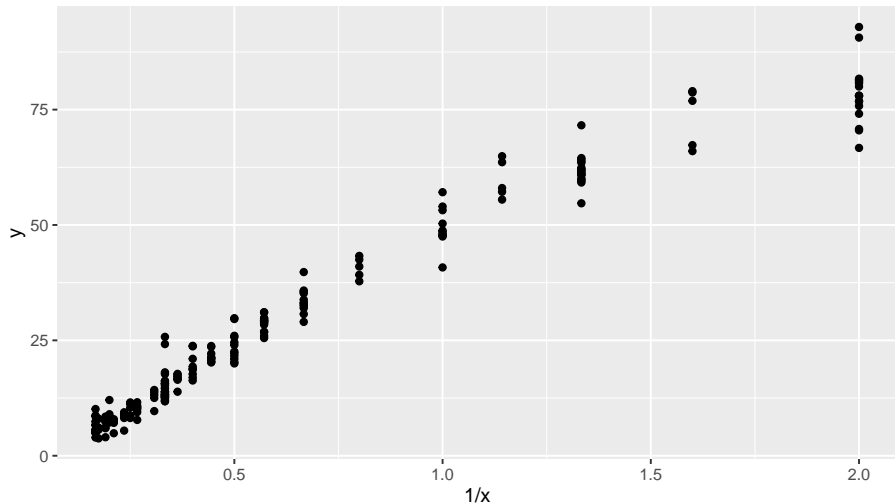
An Example: A Nonlinear Pattern



An Example: Transform y



An Example: Transform x



Nonlinear Mean Functions

In general, we consider

$$E(Y \mid \mathbf{X} = \mathbf{x}) = m(\mathbf{x}, \boldsymbol{\theta}),$$

where $\boldsymbol{\theta}$ is a vector of parameters. The function form $m(\mathbf{x}, \boldsymbol{\theta})$ can be

- ① completely specified (e.g., parametric)
- ② partially data-driven (e.g., semi-parametric)
- ③ data-driven (e.g., nonparametric)

Parametric Case

Suppose that we fully specify the function form of $m(\mathbf{x}, \boldsymbol{\theta})$, up to the values of $\boldsymbol{\theta}$.

- For example, for the data in the above example, we assume

$$E(Y | X = x) = \frac{\exp(-ax)}{b + cx},$$

where $\boldsymbol{\theta} = [a \ b \ c]^T$ is the vector of parameters.

- The model is neither linear in the parameters nor in the regressors.

Estimation

Estimation of $\boldsymbol{\theta}$ can be done by maximum likelihood or least squares.

- 1 In maximum likelihood, we assume that $Y \mid \mathbf{x}$ follows some density $f(y \mid \mathbf{x}; \boldsymbol{\theta})$ and maximize the likelihood function

$$\prod_{i=1}^n f(y_i \mid \mathbf{x}_i; \boldsymbol{\theta}).$$

- 2 In least squares, we minimize the residual sum of squares

$$\text{RSS} = \sum_{i=1}^n w_i [y_i - m(\mathbf{x}_i, \boldsymbol{\theta})]^2,$$

for some weights $\{w_i\}$ (can be used if e.g., $\text{Var}(y_i \mid \mathbf{x}_i) = \sigma^2/w_i$).

However, the closed form solution is generally not available.

General Problem

- Consider a general problem that, for a scalar-valued function $h(\boldsymbol{\beta})$, we need to find the solution of

$$\mathbf{0} = \frac{\partial h(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}.$$

- The solution is approximately the solution of

$$\mathbf{0} = \frac{\partial h(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \approx \frac{\partial h(\boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}} + \frac{\partial^2 h(\boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(t)})$$

for some known $\boldsymbol{\theta}^{(t)}$, which yields

$$\boldsymbol{\theta} = \boldsymbol{\theta}^{(t)} - \left[\frac{\partial^2 h(\boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right]^{-1} \frac{\partial h(\boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}},$$

if the Hessian matrix is invertible.

Newton-Raphson Method or Newton's Method

We can name a first guess of $\boldsymbol{\theta}$, $\boldsymbol{\theta}^{(0)}$, and update the parameter estimates using

$$\begin{aligned}\boldsymbol{\theta}^{(1)} &= \boldsymbol{\theta}^{(0)} - \left[\frac{\partial^2 h(\boldsymbol{\theta}^{(0)})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right]^{-1} \frac{\partial h(\boldsymbol{\theta}^{(0)})}{\partial \boldsymbol{\theta}}, \\ \boldsymbol{\theta}^{(2)} &= \boldsymbol{\theta}^{(1)} - \left[\frac{\partial^2 h(\boldsymbol{\theta}^{(1)})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right]^{-1} \frac{\partial h(\boldsymbol{\theta}^{(1)})}{\partial \boldsymbol{\theta}}, \\ &\vdots\end{aligned}$$

until $\frac{\partial h(\boldsymbol{\theta}^{(t+1)})}{\partial \boldsymbol{\theta}}$ is sufficiently close to $\mathbf{0}$ or $\boldsymbol{\theta}^{(t+1)}$ and $\boldsymbol{\theta}^{(t)}$ are sufficiently close.

Newton-Gauss Method for Nonlinear Least Squares

Suppose that

$$m(\mathbf{x}_i, \boldsymbol{\theta}) \approx m(\mathbf{x}_i, \boldsymbol{\theta}^{(t)}) + \left[\frac{\partial m(\mathbf{x}_i, \boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}} \right]^T (\boldsymbol{\theta} - \boldsymbol{\theta}^{(t)}).$$

Then the residual sum of squares satisfies

$$\text{RSS} \approx \sum_{i=1}^n w_i \left\{ y_i - m(\mathbf{x}_i, \boldsymbol{\theta}^{(t)}) - \left[\frac{\partial m(\mathbf{x}_i, \boldsymbol{\theta}^{(t)})}{\partial \boldsymbol{\theta}} \right]^T (\boldsymbol{\theta} - \boldsymbol{\theta}^{(t)}) \right\}^2,$$

where $y_i - m(\mathbf{x}_i, \boldsymbol{\theta}^{(t)})$ is our new response and $\boldsymbol{\theta} - \boldsymbol{\theta}^{(t)}$ is our new parameter vector. Its minimizer has a closed form expression. We iterate until no more improvements can be made.

Inference

Under some regularity conditions, $\hat{\boldsymbol{\theta}} \mid \mathbf{X}$ is approximately

$$N(\boldsymbol{\theta}, \sigma^2 [\mathbf{U}^T(\boldsymbol{\theta}) \mathbf{W} \mathbf{U}(\boldsymbol{\theta})]),$$

for large enough n .

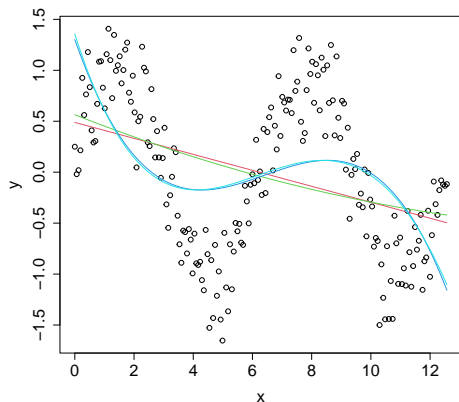
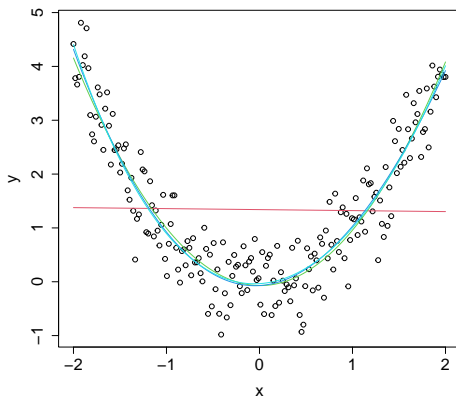
An estimator of σ^2 is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n w_i [y_i - m(\mathbf{x}_i, \hat{\boldsymbol{\theta}})]^2}{n - p},$$

where p is the number of parameters in the mean function.

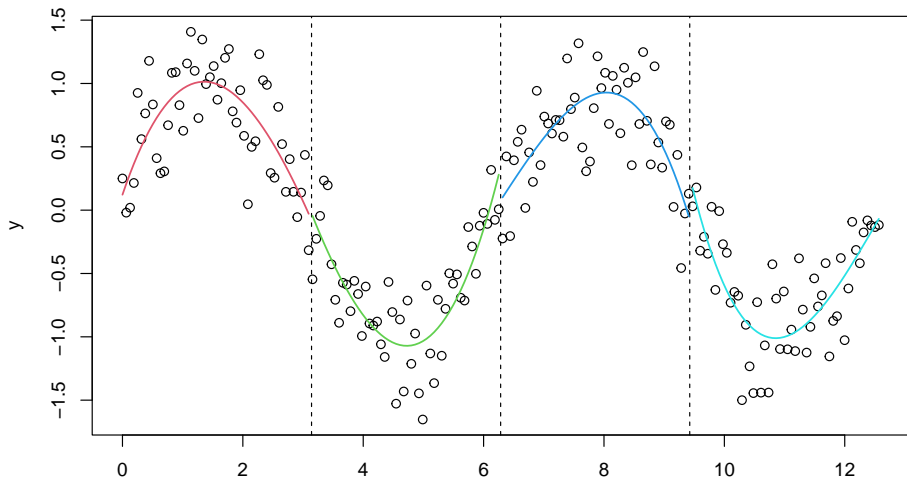
Specify Mean Function

It is not always easy to specify the closed form expression of the mean function $m(\mathbf{x}, \boldsymbol{\theta})$.



Alternative: Piecewise Polynomial

We partition the data into several parts and fit polynomials to each part separately.



Piecewise polynomial

A piecewise polynomial is obtained by

- 1 partitioning the range of x into contiguous intervals using the **knots**,
- 2 Between every two consecutive knots, fitting a polynomial model (in x) to the data points in the interval.

In practice, it is common to use the **cubic polynomials** (with degree 3 and order 4).

Example: Piecewise Linear

Consider two knots ξ_1 and ξ_2 . We fit three linear models

$$\begin{aligned}\text{for } x < \xi_1 : \quad & E_1(Y | X = x) = \beta_0^{(1)} + \beta_1^{(1)}x, \\ \text{for } \xi_1 \leq x < \xi_2 : \quad & E_2(Y | X = x) = \beta_0^{(2)} + \beta_1^{(2)}x, \\ \text{for } x \geq \xi_2 : \quad & E_3(Y | X = x) = \beta_0^{(3)} + \beta_1^{(3)}x.\end{aligned}$$

It is the same as

$$\begin{aligned}E(Y | X = x) = & 1(x < \xi_1) \left(\beta_0^{(1)} + \beta_1^{(1)}x \right) \\ & + 1(\xi_1 \leq x < \xi_2) \left(\beta_0^{(2)} + \beta_1^{(2)}x \right) \\ & + 1(x \geq \xi_2) \left(\beta_0^{(3)} + \beta_1^{(3)}x \right).\end{aligned}$$

However, the model fitted by piecewise polynomials are discontinuous at knots. We want our fitted model to be continuous.

Example: Piecewise Linear

In order to achieve continuity, we need to impose restrictions:

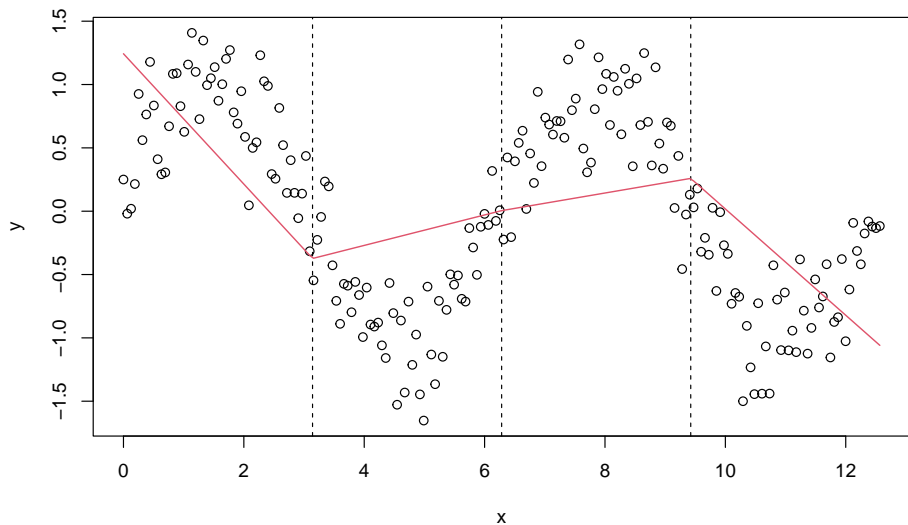
$$\begin{aligned}\beta_0^{(1)} + \beta_1^{(1)}\xi_1 &= \beta_0^{(2)} + \beta_1^{(2)}\xi_1, \\ \beta_0^{(2)} + \beta_1^{(2)}\xi_2 &= \beta_0^{(3)} + \beta_1^{(3)}\xi_2.\end{aligned}$$

Consequently, $\beta_0^{(1)}$ and $\beta_0^{(2)}$ are not free parameters and our model becomes

$$\begin{aligned}\mathbb{E}(Y \mid X = x) &= \left(\beta_0^{(3)} + \beta_1^{(3)}\xi_2 - \beta_1^{(2)}\xi_2 + \beta_1^{(2)}\xi_1 - \beta_1^{(1)}\xi_1 \right) \\ &\quad + \beta_1^{(1)}x + \left(\beta_1^{(2)} - \beta_1^{(1)} \right) \max(0, x - \xi_1) \\ &\quad + \left(\beta_1^{(3)} - \beta_1^{(2)} \right) \max(0, x - \xi_2).\end{aligned}$$

It is the same as regress Y on the intercept, x , $\max(0, x - \xi_1)$ and $\max(0, x - \xi_2)$.

Example: Piecewise Linear



Example: Piecewise Cubic Polynomial

Consider two knots ξ_1 and ξ_2 . We fit three cubic polynomials

$$\text{for } x < \xi_1 : \quad E_1(Y | X = x) = \beta_0^{(1)} + \beta_1^{(1)}x + \beta_2^{(1)}x^2 + \beta_3^{(1)}x^3,$$

$$\text{for } \xi_1 \leq x < \xi_2 : \quad E_2(Y | X = x) = \beta_0^{(2)} + \beta_1^{(2)}x + \beta_2^{(2)}x^2 + \beta_3^{(2)}x^3,$$

$$\text{for } x \geq \xi_2 : \quad E_3(Y | X = x) = \beta_0^{(3)} + \beta_1^{(3)}x + \beta_2^{(3)}x^2 + \beta_3^{(3)}x^3.$$

It is the same as

$$\begin{aligned} E(Y | X = x) = & 1(x < \xi_1) \left(\beta_0^{(1)} + \beta_1^{(1)}x + \beta_2^{(1)}x^2 + \beta_3^{(1)}x^3 \right) \\ & + 1(\xi_1 \leq x < \xi_2) \left(\beta_0^{(2)} + \beta_1^{(2)}x + \beta_2^{(2)}x^2 + \beta_3^{(2)}x^3 \right) \\ & + 1(x \geq \xi_2) \left(\beta_0^{(3)} + \beta_1^{(3)}x + \beta_2^{(3)}x^2 + \beta_3^{(3)}x^3 \right). \end{aligned}$$

However, we still cannot guarantee continuity. We want our fitted model to be continuous and **smooth** (sufficiently many continuous derivatives).

Cubic Spline

In order to produce a continuous and smooth fitted curve, we will impose the following constraints.

- ① The fitted curve must be continuous everywhere, including the knots.
- ② The fitted curve has continuous first and second order derivatives.

If piecewise cubic polynomials are used, then we have a **cubic spline**.

Cubic Spline: Restrictions

In order to achieve continuity, we need

$$\begin{aligned}E_1(Y | X = \xi_1) &= E_2(Y | X = \xi_1) \\E_2(Y | X = \xi_2) &= E_3(Y | X = \xi_2).\end{aligned}$$

In order to achieve smoothness, we need,

$$\begin{aligned}\frac{dE_1(Y | X = x)}{dx} \Big|_{x=\xi_1} &= \frac{dE_2(Y | X = x)}{dx} \Big|_{x=\xi_1} \\ \frac{dE_2(Y | X = x)}{dx} \Big|_{x=\xi_2} &= \frac{dE_3(Y | X = x)}{dx} \Big|_{x=\xi_2} \\ \frac{d^2E_1(Y | X = x)}{dx^2} \Big|_{x=\xi_1} &= \frac{d^2E_2(Y | X = x)}{dx^2} \Big|_{x=\xi_1} \\ \frac{d^2E_2(Y | X = x)}{dx^2} \Big|_{x=\xi_2} &= \frac{d^2E_3(Y | X = x)}{dx^2} \Big|_{x=\xi_2}.\end{aligned}$$

Example: Cubic Spline

Consider two knots ξ_1 and ξ_2 . With the continuity and smoothness requirements, we get

$$\begin{aligned} E(Y \mid X = x) = & \beta_0^{(3)} + \left(\beta_1^{(2)} - \beta_1^{(1)}\right) \xi_1^3 + \left(\beta_1^{(3)} - \beta_1^{(2)}\right) \xi_2^3 \\ & + \beta_1^{(1)} x + \beta_2^{(1)} x^2 + \beta_3^{(1)} x^3 \\ & + \left(\beta_3^{(2)} - \beta_3^{(1)}\right) [\max(0, x - \xi_1)]^3 \\ & + \left(\beta_3^{(3)} - \beta_3^{(2)}\right) [\max(0, x - \xi_2)]^3. \end{aligned}$$

It is the same as regress Y on the intercept, x , x^2 , x^3 , $[\max(0, x - \xi_1)]^3$, and $[\max(0, x - \xi_2)]^3$.

Cubic Spline

Suppose that we have K knots (excluding the lower and upper limits of the range). Then, there are

$$4(K + 1) - K - K - K = K + 4$$

free parameters to be estimated in cubic spline. That is, a cubic spline with K knots has $K + 4$ **degrees of freedom**. That is, the cubic spline is equivalent to

$$E(Y | X = x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \sum_{k=1}^K \beta_{k+3} [\max(0, x - \xi_k)]^3.$$

In general, **spline** is a function defined by piecewise polynomials with continuity and smoothness conditions.

Still Not Necessarily Enough

- The fit of a cubic spline is often poor for very small or very large x values, due to the lack of information and large variation.
- We need to impose additional boundary constraints, i.e. the curve is linear in the region where X is smaller than or larger than the observed values. Then we have a **natural spline**.
- If we impose the boundary constraints to a cubic spline, we have a **natural cubic spline**.

More General Approach: Basis Expansion

From the above examples, we have

$$E(Y | X = x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \sum_{k=1}^K \beta_{k+3} [\max(0, x - \xi_k)]^3.$$

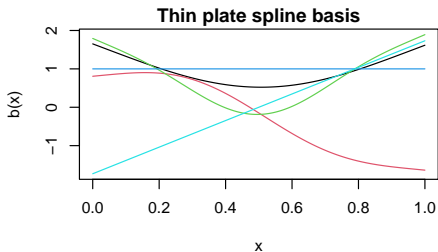
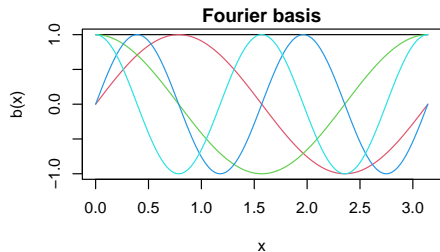
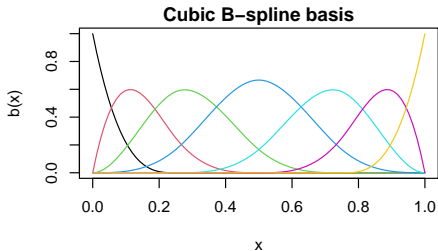
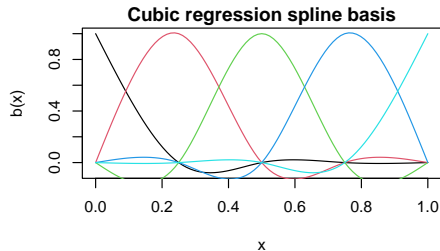
It is equivalent to using $1, x, x^2, x^3$, and $[\max(0, x - \xi_k)]^3$ as **basis functions** and performing a “global” regression.

- We choose a series of functions $\{b_k(x)\}$ and use global data to fit

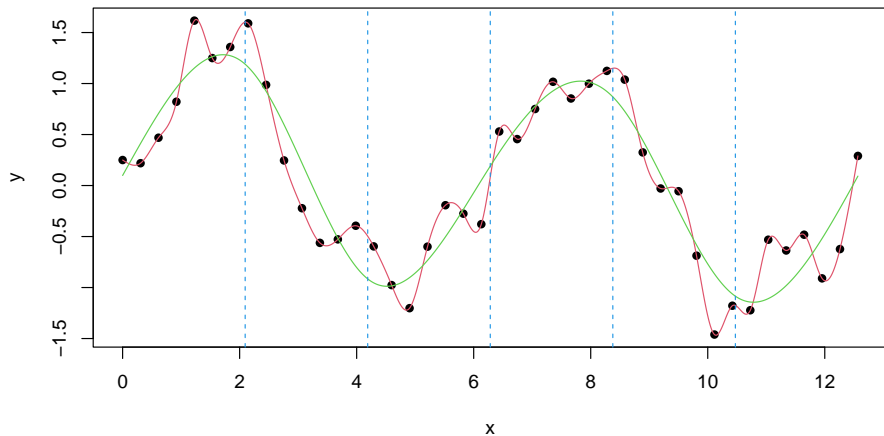
$$E(Y | X = x) = \sum_k \beta_k b_k(x).$$

- $b_k(x)$'s are the **basis functions** and $\sum_k \beta_k b_k(x)$ is the **basis expansion**.
- The basis functions are fixed ahead of time (known forms) and the shape of $E(Y | X = x)$ is controlled by varying β_k 's.

Choice of Basis Functions



Overfitting



Wiggleness

- In practice, we often want to find a function $m(x)$ that minimizes

$$\sum_{i=1}^n [y_i - m(x_i)]^2 + \lambda \int [m''(t)]^2 dt.$$

to avoid overfitting, where the smoothness is controlled by

$$\int [m''(t)]^2 dt.$$

- In general, we want to minimize

loss function + penalty for wiggleness.

- Penalties other than the derivatives are also possible, such as

$$\sum (\beta_j - \beta_{j-1})^2.$$

Ridge Regression Perspective

Suppose that

$$m(x) = \sum_k \beta_k b_k(x) = \mathbf{b}^T(x) \boldsymbol{\beta}.$$

Then

$$\begin{aligned} & \sum_{i=1}^n [y_i - m(x_i)]^2 + \lambda \int [m''(t)]^2 dt \\ &= \sum_{i=1}^n [y_i - \mathbf{b}^T(x) \boldsymbol{\beta}]^2 + \lambda \boldsymbol{\beta}^T \left[\int \frac{d^2 \mathbf{b}(t)}{dt^2} \frac{d^2 \mathbf{b}^T(t)}{dt^2} dt \right] \boldsymbol{\beta}, \end{aligned}$$

which is simply a [ridge regression](#) with regression coefficients $\boldsymbol{\beta}$.

Additive Model

- To account for non-linearity, we can assume

$$E(Y_i | \mathbf{x}_i) = m(x_{1i}, x_{2i}, \dots, x_{pi}),$$

where the function form $m()$ is estimated from the data.

- However, this formulation suffers from [curse of dimensionality](#).
- In practice, we often consider the [generalized additive model \(GAM\)](#), such as

$$m(x_{1i}, x_{2i}, \dots, x_{pi}) = f_1(x_{1i}) + f_2(x_{2i}) + f_2(x_{2i}) + \dots + f_p(x_{pi}),$$

$$m(x_{1i}, x_{2i}, \dots, x_{pi}) = f_1(x_{1i}) + f_{2,3}(x_{2i}, x_{3i}) + \dots + f_p(x_{pi}),$$

where all $f()$ have unknown forms and are estimated.

- Roughly speaking, GAM uses basis expansions to approximate unknown functions forms, and uses some penalty terms to control the wiggleness.

Kernel Regression

Suppose that we want to model $E(Y | X = x)$.

- The observed $\{x_i\}$ that are close to x should carry more information about Y than $\{x_i\}$ that are far away.
- More informative $\{x_i\}$ should be given higher weights.

The **kernel regression** estimator is

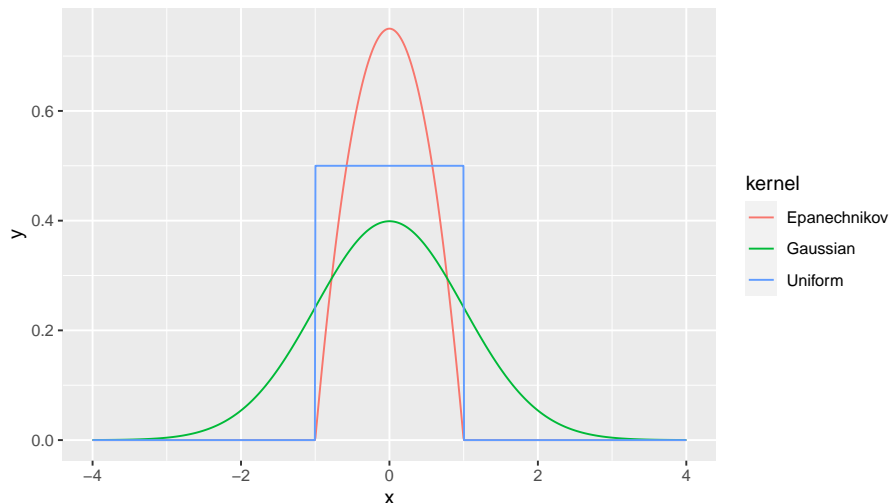
$$\hat{m}(x) = \frac{\sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) y_i}{\sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)},$$

where $K(\cdot)$ is a **kernel**, and h is a **bandwidth**. This estimator is called the **Nadaraya-Watson estimator**.

- A large h typically means a large bias and a small variance.
- A small h typically means a small bias and a high variance.

Kernel Function for Kernel Regression

A non-negative function $K(x)$ is a **kernel function** in kernel regression if $\int K(x) dx = 1$ and $K(x) = K(-x)$.



Local Regression

The **local regression** also weigh the “close points” more than the “remote points”, but the regression function in a local neighborhood of x is approximated by a polynomial with estimated coefficients.

- Say, within a neighborhood of x_0 ,

$$m(x) \approx m(x_0) + m'(x_0)(x - x_0) + \frac{1}{2}m''(x_0)(x - x_0)^2.$$

- For fixed x , the local regression estimator minimizes

$$\sum_{i=1}^n K\left(\frac{x_i - x}{h}\right) \left[y_i - \beta_0 - \beta_1(x_i - x) - \beta_2(x_i - x)^2 \right]^2.$$

- The estimated model is $\hat{m}(x) = \hat{\beta}_0$.

Limitation

A limitation of classic parametric regression is that we need to specify the form of the covariates \mathbf{x} in $\mathbf{x}^T \boldsymbol{\beta}$ ourselves.

The gradient of the log-likelihood function is

$$\frac{\partial \text{RSS}}{\partial \beta_k} = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta}) x_{ik}.$$

If the [gradient descend](#) algorithm is used to update $\boldsymbol{\beta}$, the initial guess $\boldsymbol{\beta}^{(0)} = \mathbf{0}$ of $\boldsymbol{\beta}$ is improved by

$$\begin{aligned}\boldsymbol{\beta}^{(1)} &= \boldsymbol{\beta}^{(0)} - \sum_{i=1}^n y_i \mathbf{x}_i = - \sum_{i=1}^n y_i \mathbf{x}_i, \\ \boldsymbol{\beta}^{(2)} &= \boldsymbol{\beta}^{(1)} - \sum_{i=1}^n \left(y_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(1)} \right) \mathbf{x}_i,\end{aligned}$$

until no updates can be made.

Euclidean Inner Product

- Hence, we can expect the final estimate of β to be a linear combination of \mathbf{x}_i , say

$$\hat{\beta} = \sum_{i=1}^n d_i \mathbf{x}_i,$$

for some coefficients $\{d_i\}$.

- For a new \mathbf{x}_0 , the predicted value of Y is

$$\mathbf{x}_0^T \hat{\beta} = \sum_{i=1}^n d_i \mathbf{x}_0^T \mathbf{x}_i = \sum_{i=1}^n d_i \langle \mathbf{x}_i, \mathbf{x}_0 \rangle,$$

where $\langle \mathbf{x}_i, \mathbf{x}_0 \rangle$ denotes the Euclidean inner product.

New Features

If we are not satisfied with the results using \mathbf{x} , we can create a vector of new features $\boldsymbol{\delta}(\mathbf{x}_i)$ and fit a new model

$$\mathbb{E}(Y_i \mid \mathbf{x}_i) = \boldsymbol{\gamma}^T \boldsymbol{\delta}(\mathbf{x}_i).$$

For a new \mathbf{x}_0 , the predicted value of Y is

$$\begin{aligned} \boldsymbol{\delta}^T(\mathbf{x}_0) \hat{\boldsymbol{\gamma}} &= \sum_{i=1}^n d_i^* \boldsymbol{\delta}^T(\mathbf{x}_0) \boldsymbol{\delta}(\mathbf{x}_i) \\ &= \sum_{i=1}^n d_i^* \langle \boldsymbol{\delta}(\mathbf{x}_0), \boldsymbol{\delta}(\mathbf{x}_i) \rangle, \end{aligned}$$

for some coefficients $\{d_i^*\}$.

Kernel Function and Kernel Matrix

A function $\kappa(\mathbf{x}, \mathbf{z})$ is a **kernel function** if

- ① it is symmetric, $\kappa(\mathbf{x}, \mathbf{z}) = \kappa(\mathbf{z}, \mathbf{x})$,
- ② the **kernel matrix** \mathbf{K} with (i, j) th entry $\kappa(\mathbf{x}_i, \mathbf{x}_j)$ is positive semi-definite for all $\mathbf{x}_1, \dots, \mathbf{x}_n$.

If $\kappa(\mathbf{x}, \mathbf{z})$ is a kernel function, then you must be able to find a function $\delta(\cdot)$ such that

$$\kappa(\mathbf{x}, \mathbf{z}) = \delta^T(\mathbf{x}) \delta(\mathbf{z}).$$

As a kernel function, we will have an **eigen-decomposition**

$$\kappa(\mathbf{x}, \mathbf{z}) = \sum_{m=1}^{\infty} \rho_m e_m(\mathbf{x}) e_m(\mathbf{z}),$$

for some eigenvalues ρ_k and eigenfunctions $e_m(\mathbf{x})$.

Kernel Trick

Kernel trick is a commonly used trick to create new features from your original observed features.

- If the new features are created using $\delta(\mathbf{x})$, then the prediction is

$$\delta^T(\mathbf{x}_0) \hat{\gamma} = \sum_{i=1}^n d_i^* \langle \delta(\mathbf{x}_0), \delta(\mathbf{x}_i) \rangle.$$

- If κ is the corresponding kernel function, the prediction is equivalent to

$$\delta^T(\mathbf{x}_0) \hat{\gamma} = \sum_{i=1}^n d_i^* \kappa(\mathbf{x}_0, \mathbf{x}_i).$$

- This means that we only need to choose the kernel function κ , if we need to create new features from the raw covariates \mathbf{x} .

Benefits of Kernel Trick

The eigen-decomposition implies that

$$\begin{aligned}
 \sum_{j=1}^n d_j^* \kappa(\mathbf{x}_0, \mathbf{x}_j) &= \sum_{j=1}^n d_j^* \sum_{m=1}^{\infty} \rho_m e_m(\mathbf{x}_0) e_m(\mathbf{x}_j) \\
 &= \sum_{m=1}^{\infty} \underbrace{\left(\sum_{j=1}^n d_j^* \sqrt{\rho_m} e_m(\mathbf{x}_j) \right)}_{\text{our new coefficients } d_j^*} \underbrace{\sqrt{\rho_m} e_m(\mathbf{x}_0)}_{\text{our new feature } \phi_m(\mathbf{x})}
 \end{aligned}$$

That is, we are using the vector of new features $\boldsymbol{\delta}(\mathbf{x})$, possibly of infinite dimension, to model ϕ_j .

Reproducing Kernel Hilbert Space (RKHS) Regression

- The kernel trick means that we first choose a kernel function, compute the kernel matrix \mathbf{K} using the data matrix \mathbf{X} , and use \mathbf{K} as the data matrix to fit our logistic regression model.
- However, when the data matrix \mathbf{X} is of dimension $n \times p$, then the kernel matrix \mathbf{K} is $n \times n$ and we have n regression coefficients to estimate. We often cannot obtain a unique estimator.
- Hence, we often use **penalized least squares**: the coefficients can be estimated by minimizing

$$(\mathbf{y} - \mathbf{K}\boldsymbol{\gamma})^T (\mathbf{y} - \mathbf{K}\boldsymbol{\gamma}) + \lambda \boldsymbol{\gamma}^T \mathbf{K} \boldsymbol{\gamma},$$

where λ is the **tuning parameter**. It is often chosen by cross validation.

Multivariate Gaussian Distribution

The **univariate normal distribution** with mean μ and variance σ^2 has the probability density function

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}, \quad -\infty < x < \infty.$$

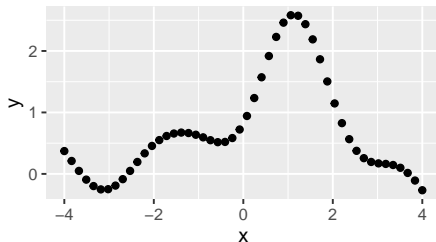
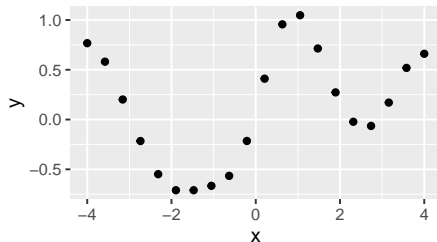
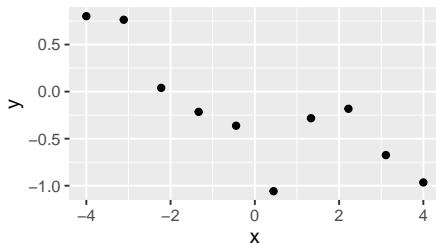
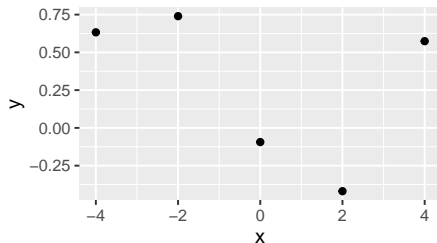
We often denote it by $X \sim N(\mu, \sigma^2)$. If $Z \sim N(0, 1)$, then $X = \sigma Z + \mu \sim N(\mu, \sigma^2)$.

Let $\mathbf{Z} = [Z_1 \ Z_2 \ \cdots \ Z_p]^T$ be a random vector, each $Z_j \sim N(0, 1)$, and Z_j is independent of Z_k for any $j \neq k$. Let \mathbf{A} be a constant matrix and $\boldsymbol{\mu}$ be a constant vector. Then,

$$\mathbf{X} = \mathbf{AZ} + \boldsymbol{\mu}$$

follows a p -dimensional **multivariate normal distribution**. It is denoted by $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Gaussian Random Numbers



Conditional Distribution

Let $\begin{bmatrix} Y_0 \\ \mathbf{y} \end{bmatrix}$ be distributed as $N_n \left(\begin{bmatrix} \mu_0 \\ \boldsymbol{\mu} \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{11} & \mathbf{K}_{12} \\ \mathbf{K}_{21} & \mathbf{K}_{22} \end{bmatrix} \right)$. Then the conditional distribution of Y_0 given that \mathbf{y} , is

$$Y_0 \mid \mathbf{y} \sim N \left\{ \mu_0 + \mathbf{K}_{12} \mathbf{K}_{22}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \mathbf{K}_{11} - \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{K}_{21} \right\},$$

provided that \mathbf{K}_{22} is invertible.

Gaussian Process Regression

Suppose that we have observed (\mathbf{y}, \mathbf{X}) and that we want to predict the value of Y_0 at \mathbf{x}_0 . For simplicity, we assume that $\boldsymbol{\mu} = \mathbf{0}$ and $\boldsymbol{\mu}_0 = \mathbf{0}$.

Then, $Y_0 \mid \mathbf{y}$ is Gaussian with

$$\mathbb{E}[Y_0 \mid \mathbf{y}] = \tau^2 \mathbf{K}_{12} (\tau^2 \mathbf{K}_{22} + \sigma^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\text{Var}[Y_0 \mid \mathbf{y}] = \tau^2 \mathbf{K}_{11} - (\tau^2 \mathbf{K}_{12}) (\tau^2 \mathbf{K}_{22} + \sigma^2 \mathbf{I})^{-1} (\tau^2 \mathbf{K}_{21}),$$

where τ^2 is a tuning parameter, and σ^2 is the error variance (also a tuning parameter) in $Y = \text{GP} + \text{error}$.

- ① The predicted value is $\hat{Y}_0 = \tau^2 \mathbf{K}_{12} (\tau^2 \mathbf{K}_{22} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$.
- ② The confidence interval is constructed by

$$\hat{Y}_0 \pm 1.96 \sqrt{K_{11} - \mathbf{K}_{12} (\mathbf{K}_{22} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{21}}.$$