

# Chapter 1. Nonparametric Curve Estimation

## Part 2

January 20, 2007

### 1 Estimation of density function

Suppose  $X_1, \dots, X_n$  are IID and have a common density function  $f(x)$ . We can draw the histogram to show the density function.

**Example 1.1** Suppose we draw 500 observations from  $N(0, 1)$ . then its histogram is as follows (left panel). For the data in [data](#) the histogram is the right panel

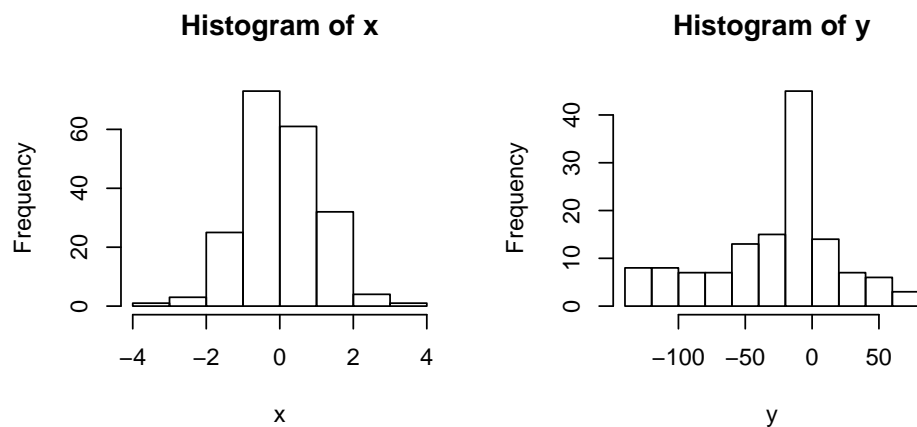


Figure 1: calculation for example 1.1: blue curves are true density function and the red lines are the estimated density functions. [\(code\)](#)

Suppose  $K()$  is a kernel function (a symmetric density function). We estimate the density by

$$\hat{f}(x) = (nh)^{-1} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

Then it is easy to see that  $\hat{f}(x)$  is density function because

$$\int \hat{f}(x) dx = 1.$$

It is more convenient to write it as

$$\hat{f}(x) = (n)^{-1} \sum_{i=1}^n K_h(x - X_i)$$

where

$$K_h(x - X_i) = \frac{1}{h} K\left(\frac{x - X_i}{h}\right).$$

**Lemma 1.2** Suppose  $X_1, \dots, X_n$  are IID and have a common density function  $f(x)$ .

$$\hat{f}(x) \rightarrow f(x), \quad \text{as } n \rightarrow \infty \text{ and } h \rightarrow 0.$$

In other words,  $\hat{f}(x) \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n K_h(X_i - x)$  is a consistent estimator of  $f(x)$ .

In practice, the bandwidth is chosen by the rule-of-thumb bandwidth selection as

$$h = 1.06 s_x n^{-1/5}.$$

where

$$s_x = \sqrt{n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2}$$

**Example 1.3 (Simulation)** samples  $X_1, \dots, X_n$  with size  $n$  are drawn from population  $X \sim N(0, 1)$ . We can estimate the density function of  $X$  as

$$\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)$$

where  $h$  is the bandwidth and  $K(x)$  is a symmetric density function.  $h$  is selected by the rule-of-thumb. The estimated density is shown in Fig. 2

**Example 1.4** The motorcycle data set:  $n = 133$ , The estimated pdf of  $X$

$$\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)$$

where  $h$  is the bandwidth and  $K(x)$  is a symmetric density function.  $h$  is selected by the rule-of-thumb

The estimated pdf of  $Y$

$$\hat{f}_n(y) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{Y_i - y}{h}\right)$$

where  $h$  is the bandwidth and  $K(x)$  is a symmetric density function.  $h$  is selected by the rule-of-thumb. The estimated density is shown in Fig. 3

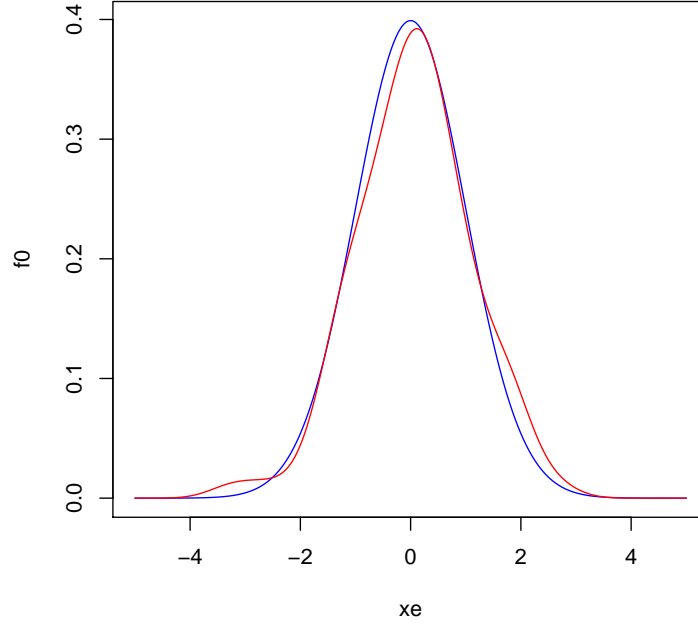


Figure 2: calculation for Example 1.3: blue curves are true density function and the red lines are the estimated density functions. [\(code\)](#)

## 2 Statistical properties of the kernel smoothing

We need to check the errors of the estimators.

### 2.1 NW-estimator

Consider model

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

where  $\varepsilon$  is independent of  $X$ . For easy of exposition, we assume that  $m(x)$  has third order derivative. Again, suppose we have the random observations

obs. 1	$X_1$	$Y_1$
obs. 2	$X_2$	$Y_2$
		$\dots$
obs. n	$X_n$	$Y_n$

or we can write the model as

$$Y_i = m(X_i) + \varepsilon_i$$

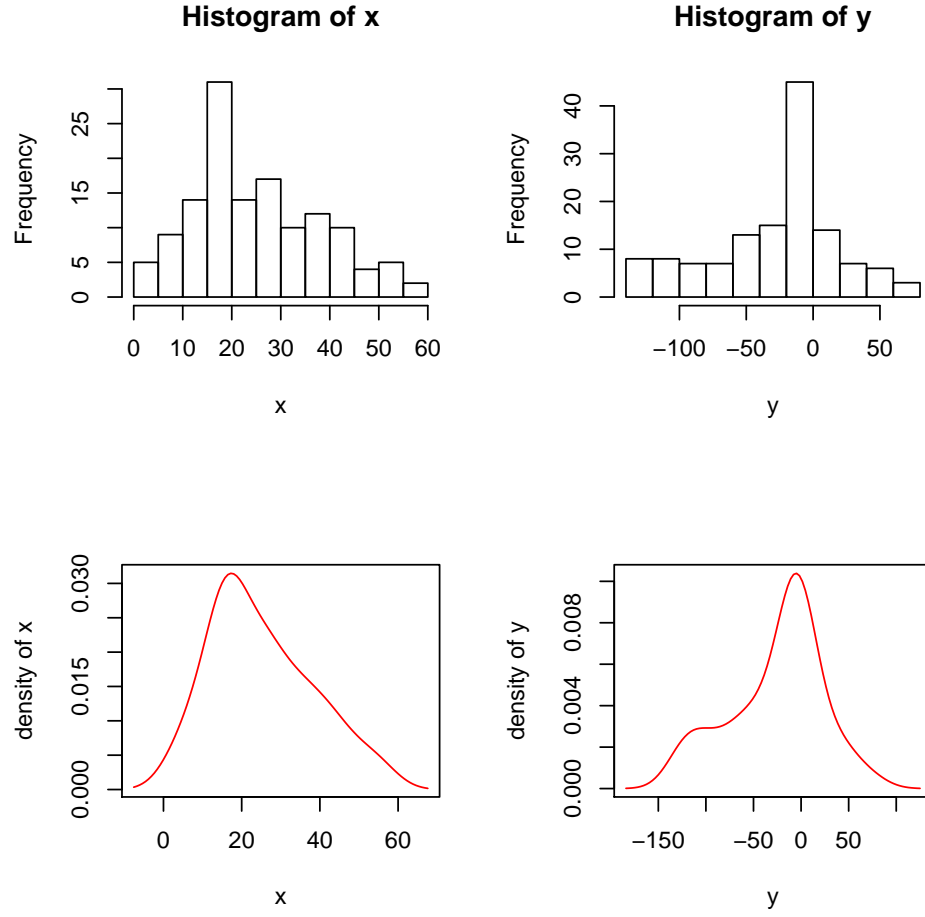


Figure 3: Calculation of Example 1.4. the estimated density functions. [\(code\)](#)

### A. Taylor Expansion of a smoothing function

For any point  $x$ , if  $X_i$  is close to  $x$  we have the Taylor expansion

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

The model can be written as

$$Y_i \approx m(x) + m'(x)(X_i - x) + \frac{1}{2}m''(x)(X_i - x)^2 + \varepsilon_i \quad (2.1)$$

### B. expansion of NW-estimator

Recall the NW-estimator is

$$\hat{m}(x) = \sum_{i=1}^n K_h(X_i - x)Y_i / \sum_{i=1}^n K_h(X_i - x).$$

By (2.1), we have

$$\hat{m}(x) \approx m(x) + \frac{\sum_{i=1}^n K_h(X_i - x) [m'(x)(X_i - x) + \frac{1}{2}m''(x)(X_i - x)^2 + \varepsilon_i]}{\sum_{i=1}^n K_h(X_i - x)}$$

We have

$$bias(\hat{m}(x)) \approx \frac{\sum_{i=1}^n K_h(X_i - x) [m'(x)(X_i - x) + \frac{1}{2}m''(x)(X_i - x)^2]}{\sum_{i=1}^n K_h(X_i - x)}$$

and

$$Var(\hat{m}(x)) = Var\left(\frac{\sum_{i=1}^n K_h(X_i - x)\varepsilon_i}{\sum_{i=1}^n K_h(X_i - x)}\right)$$

**Proposition 2.1** *Suppose that  $h \rightarrow 0$  and  $nh \rightarrow \infty$  as  $n \rightarrow \infty$ . Then at every point of continuity of  $m(x)$ ,  $f(x)$  with  $f(x) > 0$ , we have*

$$E|\hat{m}(x) - m(x)|^2 = bias^2 + variance \rightarrow 0$$

as  $n \rightarrow \infty$  and  $h \rightarrow 0$

**Theorem 2.2** *Suppose that  $h \rightarrow 0$  and  $nh \rightarrow \infty$  as  $n \rightarrow \infty$ .  $EY^2 < \infty$ . Then at every point of continuity of  $m(x)$ ,  $f(x)$  with  $f(x) > 0$ , we have*

$$biase \stackrel{def}{=} E\hat{m}(x) - m(x) = c_2\{\frac{1}{2}m''(x) + f^{-1}(x)m'(x)f'(x)\}h^2$$

and

$$variance \stackrel{def}{=} var(\hat{m}(x) - m(x)) \rightarrow \frac{d_0\sigma^2}{nhf(x)}.$$

where  $d_0 = \int K(v)^2 dv$ .

Therefore,  $\hat{m}(x)$  is a biased estimator of  $m(x)$ . The choice of  $h$  should minimize

$$\begin{aligned} E|\hat{m}(x) - m(x)|^2 &= biase^2 + variance \\ &= c_2^2\{\frac{1}{2}m''(x) + f^{-1}(x)m'(x)f'(x)\}^2h^4 + \frac{d_0\sigma^2}{nhf(x)} \end{aligned}$$

The optimal bandwidth is than

$$h_{opt} = \left\{ \frac{d_0\sigma^2}{4f(x)c_2^2\{\frac{1}{2}m''(x) + f^{-1}(x)m'(x)f'(x)\}^2} \right\}^{1/5} n^{-1/5}.$$

**Example 2.3** *n samples are drawn from*

$$Y = \sin(2\pi X) + 0.2\varepsilon.$$

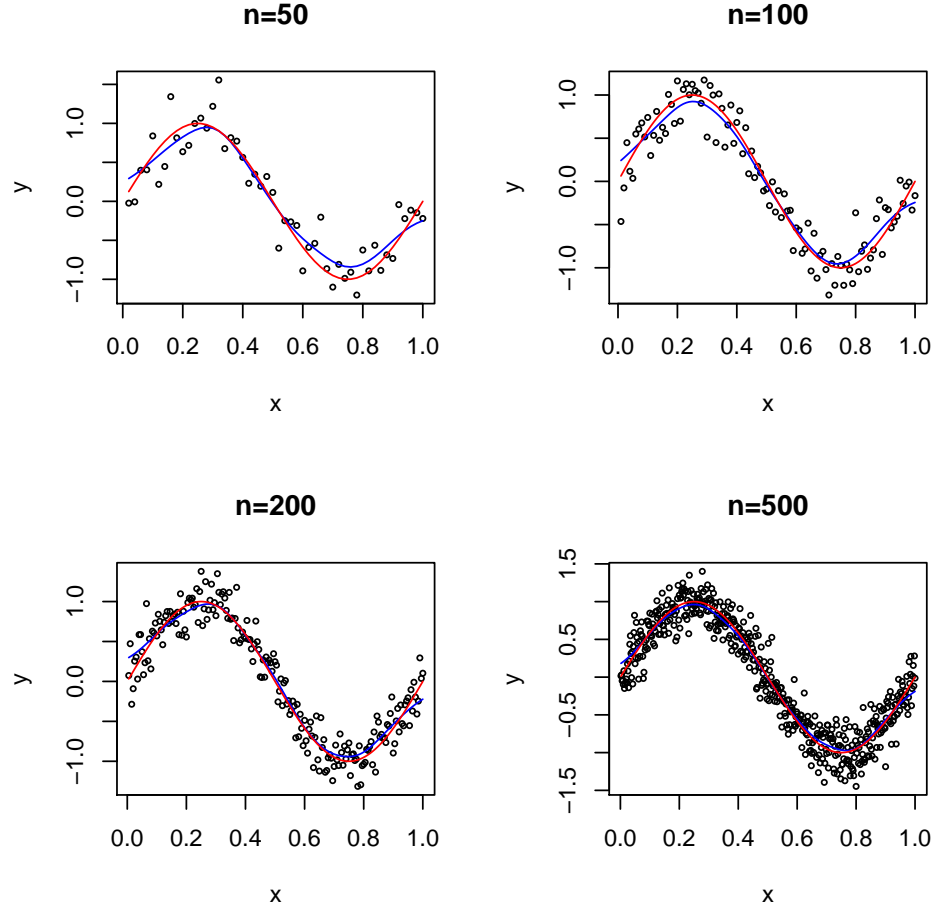


Figure 4: Calculation for Example 2.3 blue curves are true regression function and the red lines are the estimated functions

where  $X \in \text{Uniform}(0, 1)$  and  $\varepsilon \sim N(0, 1)$ .

Estimated function is shown in Fig 4

We can see from the figure, as  $n$  increase, the estimated curve tends to the true curve function. Note that

$$[\sin(2\pi x)]'' = -(2\pi)^2 \sin(2\pi x)$$

which has big absolute value at  $x = 1/4$  and  $3/4$ . Therefore, the estimated curve has big bias at those points.

**Example 2.4**  $n$  samples are drawn from

$$Y = \exp(-20X^2) + 0.2\varepsilon.$$

where  $X \in \text{Uniform}(-1, 1)$  and  $\varepsilon \sim N(0, 1)$ . Estimated function is shown in Fig 5

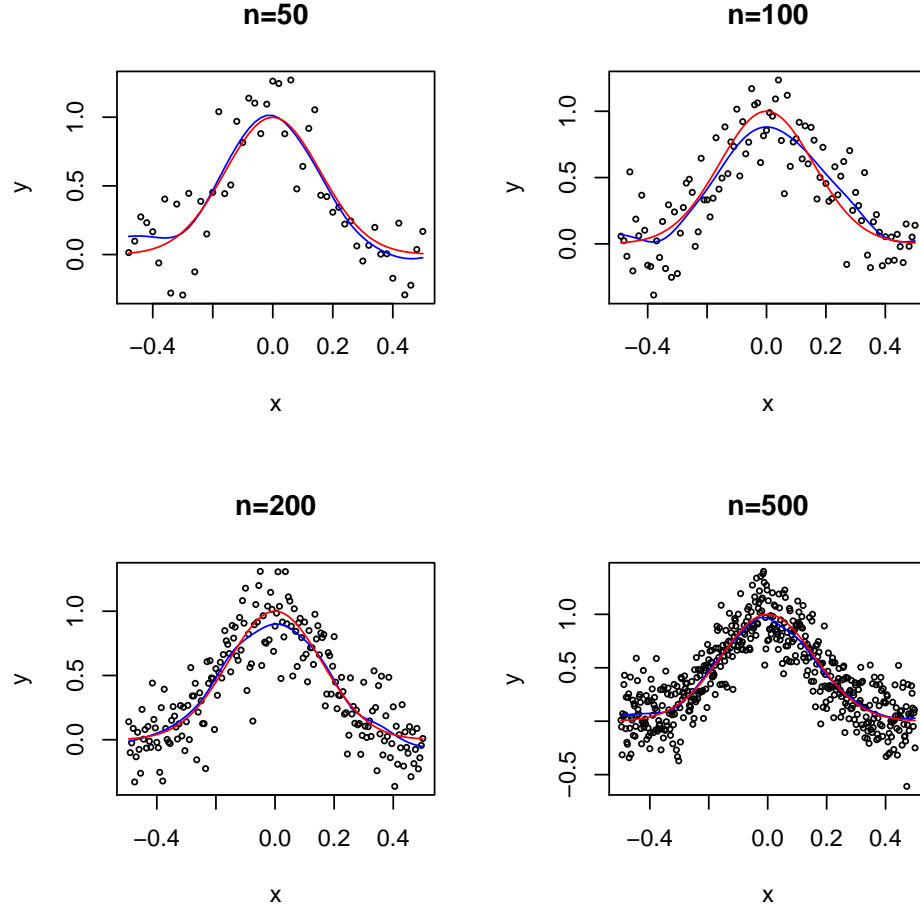


Figure 5: Calculation for Example 2.4: blue curves are true regression function and the red lines are the estimated functions [\(code\)](#)

**Lemma 2.5** *If  $X_1, \dots, X_n$  are IID and have common density function  $f(x)$ . If  $\varepsilon$  is independent of  $X$  and  $\text{Var}(\varepsilon) = \sigma^2$ , then*

$$\sqrt{nh} * n^{-1} \sum_{i=1}^n K_h(X_i - x) \varepsilon_i \xrightarrow{D} N(0, d_0 f(x) \sigma^2)$$

where  $d_0 = \int K(v)^2 dv$ .

**Theorem 2.6**

$$\sqrt{nh} \{ \hat{m}(x) - m(x) - \frac{1}{2} c_2 [m''(x) + 2f^{-1}(x)m'(x)f'(x)]h^2 \} \xrightarrow{D} N(0, \frac{d_0 \sigma^2}{f(x)})$$

If  $\sqrt{nh}h^2 \rightarrow 0$ , then

$$\sqrt{nh} \{ \hat{m}(x) - m(x) \} \xrightarrow{D} N(0, \frac{d_0 \sigma^2}{f(x)})$$

The point-wise confidence band is

$$P\left\{\hat{m}(x) - 1.96\sqrt{\frac{d_0\sigma^2}{nhf(x)}} \leq m(x) \leq \hat{m}(x) + 1.96\sqrt{\frac{d_0\sigma^2}{nhf(x)}}\right\} \approx 0.95$$

In using the above theorem to draw the confidence band, we need to calculate a number of values

1. Gaussian Kernel:

$$c_0 = 1, c_1 = 0, c_2 = 1, d_0 = 0.2821$$

For Epanechnikov kernel

$$c_0 = 1, c_1 = 0, c_2 = 0.2, d_0 = 0.6$$

2.  $\sigma^2$

$$\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (Y_i - \hat{m}(X_i))^2$$

3. the density function of  $X$

$$\hat{f}_n(x) = n^{-1} \sum_{i=1}^n K_h(X_i - x)$$

**Example 2.7 (Simulation)** consider model  $Y = \cos(\pi X) + 0.2\varepsilon$ , where  $X \sim \text{uniform}(0, 1)$  and  $\varepsilon \sim N(0, 1)$ . with sample size  $n$ . Consider  $h = 0.1$ . The estimated regression curve

$$\hat{m}(x) = \frac{\sum_{i=1}^n K_h(X_i - x)Y_i}{\sum_{i=1}^n K_h(X_i - x)}$$

The estimated pdf of  $X$ .

$$\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)$$

where  $h$  is the bandwidth and  $K(x)$  is a symmetric density function. The estimator of  $\sigma^2$  is

$$\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (Y_i - \hat{m}(X_i))^2.$$

The estimated function is shown in Fig 6.

**Example 2.8** The motorcycle data set:  $n = 133$ , consider  $h = 1.5$ . The estimated regression curve

$$\hat{m}(x) = \frac{\sum_{i=1}^n K_h(X_i - x)Y_i}{\sum_{i=1}^n K_h(X_i - x)}$$



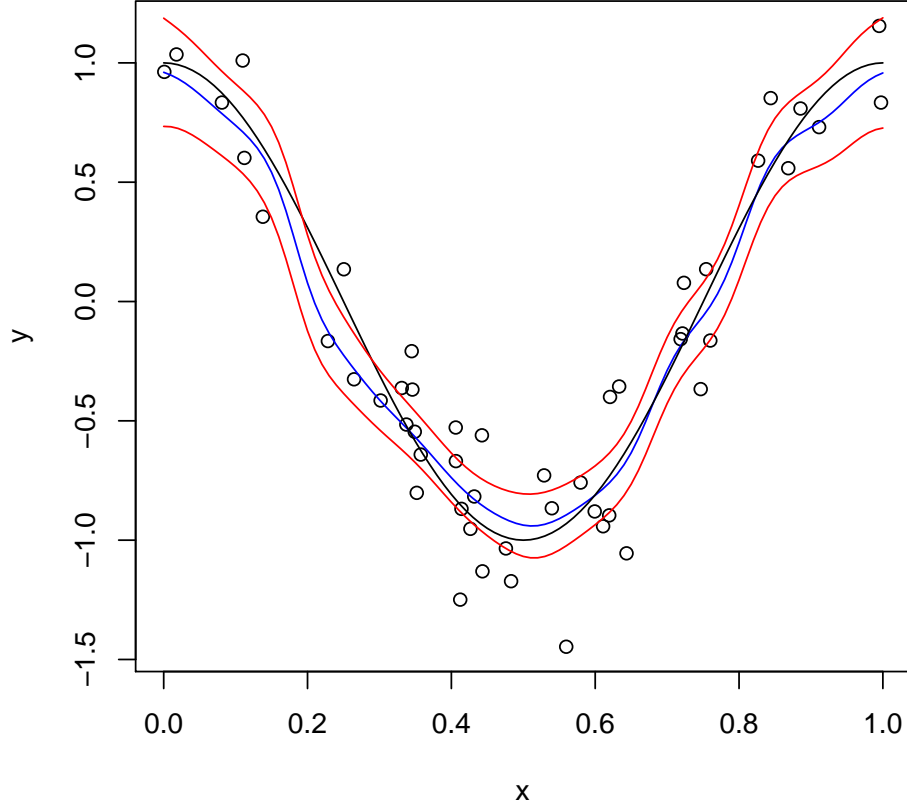


Figure 6: calculation for Example 2.7 The black line is the true function. The blue line in the central is the estimated regression function, the upper and lower red lines are the 95% point-wise confidence band. [\(code\)](#)

The estimated pdf of  $X$ .

$$\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)$$

where  $h$  is the bandwidth and  $K(x)$  is a symmetric density function. The estimator of  $\sigma^2$  is

$$\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (Y_i - \hat{m}(X_i))^2.$$

The estimated function is shown in Fig 7

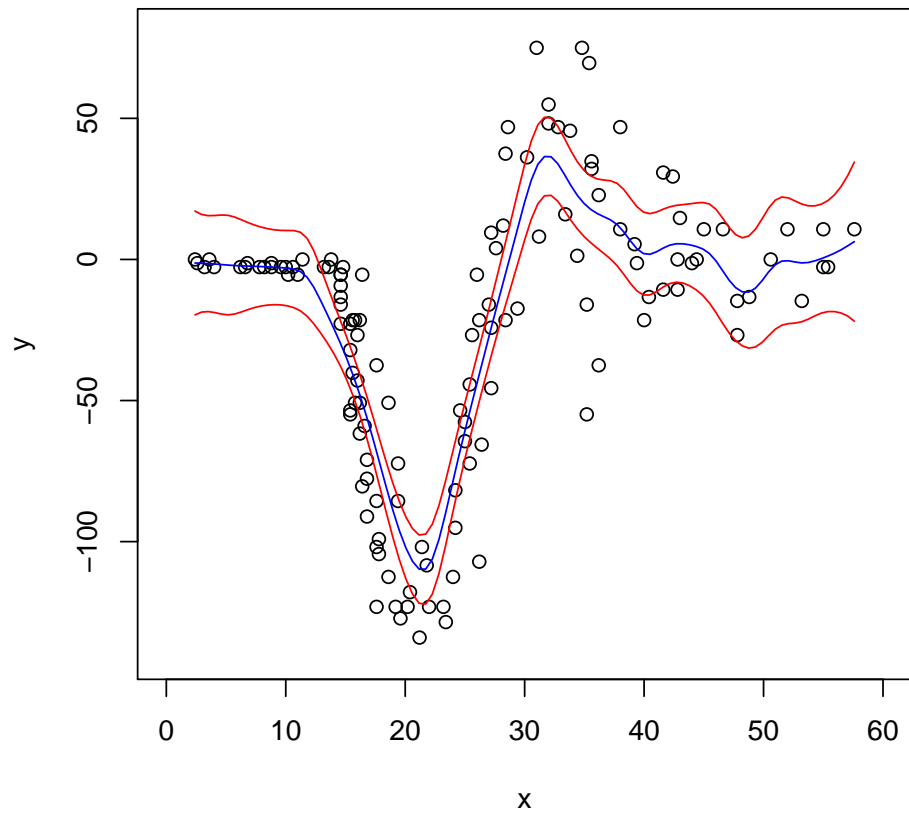


Figure 7: Calculation for Example 2.8. The line in the central is the estimated regression function, the upper and lower lines are the 95% point-wise confidence band. [\(code\)](#)