

NONLINEAR REGRESSION MODELING USING REGULARIZED LOCAL LIKELIHOOD METHOD

YOSHISUKE NONAKA* AND SADANORI KONISHI

*Graduate School of Mathematics, Kyushu University, 6-10-1 Hakozaki, Higashi-Ku,
Fukuoka 812-8581, Japan, e-mail: konishi@math.kyushu-u.ac.jp*

(Received December 16, 2003; revised August 6, 2004)

Abstract. We introduce a nonlinear regression modeling strategy, using a regularized local likelihood method. The local likelihood method is effective for analyzing data with complex structure. It might be, however, pointed out that the stability of the local likelihood estimator is not necessarily guaranteed in the case that the structure of system is quite complex. In order to overcome this difficulty, we propose a regularized local likelihood method with a polynomial function which unites local likelihood and regularization. A crucial issue in constructing nonlinear regression models is the choice of a smoothing parameter, the degree of polynomial and a regularization parameter. In order to evaluate models estimated by the regularized local likelihood method, we derive a model selection criterion from an information-theoretic point of view. Real data analysis and Monte Carlo experiments are conducted to examine the performance of our modeling strategy.

Key words and phrases: Information criteria, local maximum likelihood estimates, model selection, generalized linear models, regularization.

1. Introduction

Local likelihood estimation has received considerable attention as a useful technique for analyzing data with complex structure (Tibshirani and Hastie (1987), Hjort and Jones (1996), Copas and Eguchi (1998), Eguchi and Copas (1998), Loader (1999), Eguchi and Kim (2001), Eguchi *et al.* (2003), and so on).

A local likelihood function is constructed based on first considering a parametric model for the unknown true model. It is defined as a locally weighted log-likelihood with weights determined by a kernel function and a bandwidth. When a large bandwidth is chosen, the estimator is close to the maximum likelihood estimator and tends to have a large bias. On the other hand, when a small bandwidth is chosen, the estimator depends much on the data points and tends to have a large variance. In the local modeling, bandwidth plays an important role for controlling the trade-off between bias and variance of the estimator (Wand and Jones (1995), Simonoff (1996)).

Issues still remain in constructing nonlinear regression models based on the local likelihood from a finite and noisy set of data. First, the stability of local likelihood estimators is not guaranteed in the case that the structure of the system is quite complex.

*Now at Biostatistics Center, Kurume University, 67 Asahi-Machi, Kurume, Fukuoka 830-0011, Japan, e-mail: nonaka_yoshisuke@kurume-u.ac.jp.

In order to overcome this issue, we introduce a regularized local likelihood function with regularization parameter that controls the local likelihood function and the complexity of a nonlinear regression model. Second, the local likelihood procedure requires the choice of a bandwidth, the degree of the polynomial and also a regularization parameter. In order to choose these adjusted parameters, we derive model selection and evaluation criteria. These criteria are derived under model misspecification both for distributional and structural assumptions, which is usually the case in practice. Our modeling strategy can be easily applied to analyze multi-dimensional continuous data, and clear improvements are obtained for the use of the regularization parameter in the regularized local likelihood functions.

This article is organized as follows. In Section 2 we describe the proposed regularized local likelihood method in the context of generalized linear models and present an information criterion for evaluating the estimated models. In Section 3 we describe the multivariate regularized local likelihood method. Section 4 includes some applications to real data sets and numerical results, in which we use the regularized local likelihood method with Gaussian and logistic models. Some concluding remarks are given in Section 5.

2. Regularized local likelihood method

Local likelihood method has both parametric and nonparametric theoretical aspects. Eguchi and Kim (2001) bridged a gap between these theories for the local likelihood method. In this paper we call regression models based on the local likelihood method “nonlinear” regression models, since a modeling strategy is used for a curve fitting and a surface fitting. We consider a nonlinear modeling strategy based on the regularized local likelihood method in the context of generalized linear models (Nelder and Wedderburn (1972), McCullagh and Nelder (1989), Green and Silverman (1994), Fan *et al.* (1995)). We first introduce the regularized local likelihood method in the case of univariate explanatory variables.

2.1 Model

Suppose we have n independent observations $\{(y_i, x_i); i = 1, 2, \dots, n\}$, where y_i are random response variables and x_i are univariate explanatory variables. It is assumed that the responses y_i are generated from an unknown true distribution $G(y | x)$ with density $g(y | x)$. Generally, a regression model consists of a random component which specifies the distribution of the response Y and systematic component which presents the structure of the conditional expectation $m(x) = E[Y | x]$.

It is assumed that Y has a distribution in the exponential family, taking the form

$$(2.1) \quad f(y | x) = \exp \left\{ \frac{\theta(x)y - b(\theta(x))}{\psi(x)} + c(y, \psi(x)) \right\},$$

where $b(\cdot)$ and $c(\cdot, \cdot)$ are known functions. Here, function $\theta(\cdot)$ is called the canonical parameter and function $\psi(\cdot)$ is called the dispersion parameter. We assume that the log-likelihood $\sum_{i=1}^n \log f(y_i | x_i)$ satisfies the Bartlett (1954) identities. The mean and variance of Y can be derived easily from the well known relations

$$(2.2) \quad m(x) = E[Y | x] = b'(\theta(x)), \quad \text{Var}[Y | x] = \psi(x)b''(\theta(x)),$$

where $b'(\cdot)$ and $b''(\cdot)$ are differential of first and second order respectively. In usual generalized linear models, the unknown regression function $m(x)$ is modeled linearly via the known link function $l(\cdot)$:

$$(2.3) \quad l(m(x)) = \beta_0 + \beta_1 x.$$

The function l links the regression function to a linear space of the covariates. If $l = (b')^{-1}$, then l is called the canonical link function since in the case $l(m(x))$ is the canonical parameter in the exponential family (2.1). Expressions (2.2) and (2.3) characterize the generalized linear models.

In some situations, the use of the linear relationship (2.3) has a problem. Trial and error are required in order to search for a reasonable parametric link function. So we use the technique with local polynomial. The aim of the local modeling approach is searching for a model that describes the data well.

Here, we focus on estimating $\theta(x)$ since estimating $m(x)$ is equivalent to estimating $\theta(x) = (b')^{-1}(m(x))$. Assume that the function $\theta(x)$ is at least p -times differentiable at a point x_0 . Then $\theta(x)$ can be approximated locally by a polynomial of degree p for x in a neighbourhood of x_0

$$(2.4) \quad \begin{aligned} \theta(x) &\approx \theta(x_0) + \theta^{(1)}(x_0)(x - x_0) + \dots + \theta^{(p)}(x_0)(x - x_0)^p/p! \\ &= \beta(x_0)^T \mathbf{x}(x; x_0), \end{aligned}$$

where $\beta_j(x_0) = \theta^{(j)}(x_0)/j!$ ($j = 0, 1, \dots, p$), $\beta(x_0) = (\beta_0(x_0), \beta_1(x_0), \dots, \beta_p(x_0))^T$ and $\mathbf{x}(x; x_0) = (1, x - x_0, \dots, (x - x_0)^p)^T$. Then the data $\{y_1, y_2, \dots, y_n\}$ are summarized by a model from a class of probability densities

$$(2.5) \quad \begin{aligned} &f(y_i | x_i; \beta(x_0), \psi(x_0)) \\ &= \exp \left\{ \frac{\beta(x_0)^T \mathbf{x}(x_i; x_0) y_i - b(\beta(x_0)^T \mathbf{x}(x_i; x_0))}{\psi(x_0)} + c(y_i, \psi(x_0)) \right\}, \end{aligned}$$

where $\psi(x_0)$ is the dispersion parameter at x_0 .

2.2 Estimation

The unknown parameters $\beta(x_0)$ and $\psi(x_0)$ in (2.5) may be estimated by maximizing the local log-likelihood function. However, when fitting a nonlinear model to data with complex structure, the local likelihood method does not yield satisfactory results. Instead of maximizing this function, we propose choosing the parameters $\beta(x_0)$ and $\psi(x_0)$ to maximize the regularized local log-likelihood function:

$$(2.6) \quad \begin{aligned} &RL(\beta(x_0), \psi(x_0); p, h, \lambda) \\ &= \sum_{i=1}^n w_h(x_i; x_0) \log f(y_i | x_i; \beta(x_0), \psi(x_0)) - n\lambda \beta(x_0)^T K \beta(x_0)/2 \\ &= \sum_{i=1}^n w_h(x_i; x_0) \left\{ \frac{\beta(x_0)^T \mathbf{x}(x_i; x_0) y_i - b(\beta(x_0)^T \mathbf{x}(x_i; x_0))}{\psi(x_0)} + c(y_i, \psi(x_0)) \right\} \\ &\quad - n\lambda \beta(x_0)^T K \beta(x_0)/2. \end{aligned}$$

The first term in the right side of the equation (2.6) is a usual local log-likelihood function, where $w_h(x; x_0)$ is a weight function and we use the Gaussian kernel:

$$(2.7) \quad w_h(x; x_0) = (2\pi h^2)^{-1/2} \exp \left\{ -\frac{(x - x_0)^2}{2h^2} \right\}, \quad x \in (-\infty, \infty).$$

h is a bandwidth which controls the smoothness of estimated curve. The second term in the right side of the equation (2.6) is a regularization term. Typical forms for the $(p+1) \times (p+1)$ matrix K are given in the following forms

$$(2.8) \quad P1 : K = \begin{pmatrix} 0 & \mathbf{0}_p^T \\ \mathbf{0}_p & D_k^T D_k \end{pmatrix}, \quad P2 : K = \begin{pmatrix} 0 & \mathbf{0}_p^T \\ \mathbf{0}_p & I_p \end{pmatrix},$$

where I_p is a p -dimensional identity matrix, $\mathbf{0}_p$ is a p -dimensional zero vector and D_k is a $(p-k) \times p$ matrix that represents the difference operator given by

$$D_k = \begin{pmatrix} (-1)^0 {}_n C_0 & \cdots & (-1)^k {}_n C_k & 0 & \cdots & \cdots & 0 \\ 0 & (-1)^0 {}_n C_0 & \cdots & (-1)^k {}_n C_k & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & 0 & (-1)^0 {}_n C_0 & \cdots & (-1)^k {}_n C_k \end{pmatrix}$$

with ${}_n C_k = n! / \{k!(n-k)!\}$. The λ is a regularization parameter which controls the local log-likelihood and the complexity of the estimated model.

The estimator $\hat{\beta}(x_0)$ is obtained as a solution of $\partial RL(\beta(x_0), \psi(x_0); p, h, \lambda) / \partial \beta = \mathbf{0}$. This equation is generally nonlinear in $\beta(x_0)$, so we optimize $\beta(x_0)$ by the iterative algorithm. After the estimator $\hat{\beta}(x_0)$ is obtained, the estimator $\hat{\psi}(x_0)$ is given as a solution of $\partial RL(\hat{\beta}(x_0), \psi(x_0); p, h, \lambda) / \partial \psi = 0$. Replacing the unknown parameter $\beta(x_0)$ and $\psi(x_0)$ by $\hat{\beta}(x_0)$ and $\hat{\psi}(x_0)$ respectively and noting that parameter $\theta(x_i)$ in the exponential family (2.1) is directly given by $\theta(x_i) = \beta(x_i)^T \mathbf{x}(x_i; x_i) = \beta_0(x_i)$, we obtain the estimated model as follows

$$(2.9) \quad f(y_i | x_i, \hat{\beta}(x_i), \hat{\psi}(x_i)) = \exp \left\{ \frac{\hat{\beta}_0(x_i) y_i - b(\hat{\beta}_0(x_i))}{\hat{\psi}(x_i)} + c(y_i, \hat{\psi}(x_i)) \right\}.$$

2.2.1 Gaussian model

We consider the following model:

$$(2.10) \quad y_i = m(x_i) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2(x_i)) \quad (i = 1, \dots, n),$$

where $m(\cdot)$ is an unknown smooth function. Then taking $b(\theta(x_i)) = \theta(x_i)^2/2$, $\psi(x_0) = \sigma^2(x_0)$,

$$c(y_i, \psi(x_0)) = -\frac{y_i^2}{2\psi(x_0)} - \frac{1}{2} \log\{2\pi\psi(x_0)\} = -\frac{1}{2} \left\{ \frac{y_i}{\sigma(x_0)} \right\}^2 - \frac{1}{2} \log\{2\pi\sigma^2(x_0)\}$$

and $l(m(x_i)) = m(x_i) = b'(\theta(x_i)) = \theta(x_i) = \beta(x_0)^T \mathbf{x}(x_i; x_0)$ in the exponential family of densities (2.5), we have a nonlinear regression model with Gaussian noise which can be expressed as

$$(2.11) \quad f_N(y_i | x_i; \beta(x_0), \sigma^2(x_0)) = \{2\pi\sigma^2(x_0)\}^{-1/2} \exp \left[-\frac{\{y_i - \beta(x_0)^T \mathbf{x}(x_i; x_0)\}^2}{2\sigma^2(x_0)} \right].$$

A $(p+1)$ -dimensional parameter $\beta(x_0)$ and an error variance $\sigma^2(x_0)$ in equation (2.11) are estimated by the maximization of the regularized local log-likelihood function (2.6).

Then the estimators of $\beta(x_0)$ and $\sigma^2(x_0)$ are explicitly given by

$$(2.12) \quad \begin{aligned} \hat{\beta}(x_0) &= (X^T W X + n\zeta K)^{-1} X^T W \mathbf{y}, \\ \hat{\sigma}^2(x_0) &= \{\mathbf{y} - X\hat{\beta}(x_0)\}^T W \{\mathbf{y} - X\hat{\beta}(x_0)\} / \text{tr}(W), \end{aligned}$$

where $X = (\mathbf{x}(x_1; x_0), \dots, \mathbf{x}(x_n; x_0))^T$, $W = \text{diag}\{w_h(x_i; x_0)\}$, $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$ and $\zeta = \lambda\sigma^2$ is another expression of the regularization parameter. Replacing the unknown parameters $\beta(x_0)$ and $\sigma^2(x_0)$ in (2.11) by their sample estimators $\hat{\beta}(x_i)$ and $\hat{\sigma}^2(x_i)$, we obtain the statistical model

$$(2.13) \quad f_N(y_i | x_i; \hat{\beta}(x_i), \hat{\sigma}^2(x_i)) = \{2\pi\hat{\sigma}^2(x_i)\}^{-1/2} \exp \left[-\frac{\{y_i - \hat{\beta}_0(x_i)\}^2}{2\hat{\sigma}^2(x_i)} \right].$$

The local polynomial ridge regression proposed by Seifert and Gasser (1996, 2000) may be considered as a special case of our regularized local likelihood method. Moreover when $\lambda = 0$, the estimator $\hat{\beta}(x_0)$ in (2.12) is equivalent to the local polynomial estimator (Stone (1977), Cleveland (1979), Fan and Gijbel (1996), Loader (1999)).

Example 2.1. We illustrate the proposed regularized local likelihood modeling by fitting curve to the simulation data. The random samples $\{(y_i, x_i); i = 1, 2, \dots, 100\}$ were generated from the true regression model $y_i = \sin(18\pi x_i) + 5x_i \cos(18\pi x_i - \pi/2) + 4x_i + \varepsilon_i$, $\varepsilon_i \sim N(0, 0.7^2)$, where the design points x_i are uniformly distributed in $[0, 1]$. Figure 1(a) shows the true curve and the scatterplot of the data. We apply the local likelihood method and the regularized local likelihood method to the simulation data, where we fix the degree of polynomial $p = 5$ since we focus on the efficiency of regularization parameter. A broken line in Fig. 1(b) gives the smoothed curve for the bandwidth $h = 0.07$ without the help of regularization ($\lambda = \zeta = 0$). This curve is obviously oversmoothed, but it is impossible to apply the local likelihood method for a smaller h since the matrix $(X^T W X)^{-1}$ in (2.12) is not computable in practice. A solid line in Fig. 1(b) gives the smoothed curve for the bandwidth $h = 0.015$ with the help of

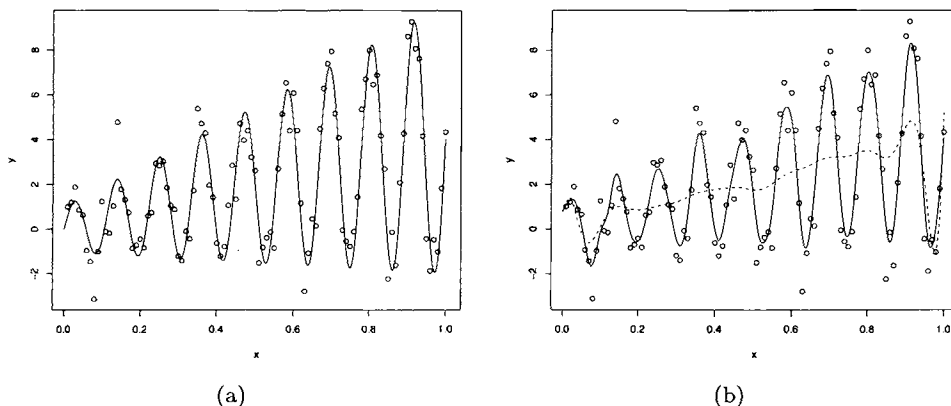


Fig. 1. Comparison of the true curve and the smoothed curves with $p = 5$. (a) shows the true curve $y = \sin(18\pi x) + 5x \cos(18\pi x - \pi/2) + 4x$ and the scatterplot of the data. (b) shows the smoothed curve for the local likelihood methods with $h = 0.07$ (broken line) and the smoothed curve for the regularized local likelihood methods with $h = 0.015$ and $\zeta = 10^{-8}$ (solid line).

regularization, where we use the regularization parameter $\zeta = 10^{-8}$. This curve gives a good representation of the underlying function over the region $[0, 1]$. We observe that by appropriate choice of h and ζ , our nonlinear regression modeling strategy can capture the true structure generating the data for large p .

2.2.2 Logistic model

Let y_1, y_2, \dots, y_n be independent sequences of binary random variables taking values of 0 or 1 with conditional probabilities $\Pr(y_i = 1 \mid x_i) = \pi(x_i)$ and $\Pr(y_i = 0 \mid x_i) = 1 - \pi(x_i)$, where x_i are univariate explanatory variables. Taking $\theta(x_i) = \log[\pi(x_i)/\{1 - \pi(x_i)\}]$, $b(\theta(x_i)) = \log[1 + \exp\{\theta(x_i)\}]$, $\psi(x_0) = 1$, $c(y_i, \psi(x_0)) = 0$ and $l(m(x_i)) = \log[m(x_i)/\{1 - m(x_i)\}] = \theta(x_i) = \beta(x_0)^T \mathbf{x}(x_i; x_0)$ in (2.5), we have a nonlinear logistic regression model as follows

$$(2.14) \quad f_L(y_i \mid x_i; \beta(x_0)) = \pi(x_i; x_0)^{y_i} \{1 - \pi(x_i; x_0)\}^{1-y_i},$$

where

$$(2.15) \quad \pi(x_i; x_0) = \frac{\exp\{\beta(x_0)^T \mathbf{x}(x_i; x_0)\}}{1 + \exp\{\beta(x_0)^T \mathbf{x}(x_i; x_0)\}}.$$

The unknown parameter vector $\beta(x_0)$ is estimated by maximizing the regularized local log-likelihood (2.6), where the matrix K is given by P2 in equation (2.8). Then we can obtain the estimated model as follows

$$(2.16) \quad f_L(y_i \mid x_i, \hat{\beta}(x_i)) = \hat{\pi}(x_i; x_i)^{y_i} \{1 - \hat{\pi}(x_i; x_i)\}^{1-y_i} \quad (i = 1, \dots, n),$$

where $\hat{\pi}(x_i; x_i) = \exp\{\hat{\beta}_0(x_i)\}/[1 + \exp\{\hat{\beta}_0(x_i)\}]$ is the estimated conditional probability.

Nonlinear regression models based on the regularized local likelihood method depend on a bandwidth h , the degree of polynomial p and a regularization parameter λ . In the next subsection we derive a model evaluation criterion for nonlinear regression models estimated by the regularized local likelihood method.

2.3 Model selection

Suppose that we observe a realization of a random variable with distribution as defined in Subsection 2.1. We recall that the statistical model $f(y \mid x, \hat{\beta}(x), \hat{\psi}(x))$ defined in (2.9) is constructed within the generalized linear model framework. We want to assess the closeness of $f(y \mid x, \hat{\beta}(x), \hat{\psi}(x))$ to the true model $g(y \mid x)$ from a predictive point of view.

Konishi and Kitagawa (1996) proposed a model selection criterion as an estimator of the Kullback-Leibler information between the true model and the estimated model. As shown in the Appendix, we can obtain an information criterion to evaluate the estimated model $f(y_i \mid x_i, \hat{\beta}(x_i), \hat{\psi}(x_i))$ in (2.9) with $(p + 2)$ -dimensional parameter estimator $\hat{\theta}(x) = (\hat{\beta}(x)^T, \hat{\psi}(x))^T$. The information criterion based on the regularized local likelihood method with the exponential family is given by

$$(2.17) \quad \text{GIC}_{h,p,\lambda} = -2 \sum_{i=1}^n \left\{ \frac{\hat{\beta}_0(x_i) y_i - b(\hat{\beta}_0(x_i))}{\hat{\psi}(x_i)} + c(y_i, \hat{\psi}(x_i)) \right\} + \frac{2}{n} \sum_{i=1}^n \text{tr}\{\hat{R}(x_i)^{-1} \hat{Q}(x_i)\},$$

where we use the following notations:

$$\begin{aligned}
 (2.18) \quad \hat{Q}(x) &= \frac{1}{n\hat{\psi}(x)} \begin{bmatrix} A^{(2)}X/\hat{\psi}(x) - \lambda K\hat{\beta}(x)\mathbf{1}_n^T \Lambda X & A^{(1)}\mathbf{p} - \hat{\psi}(x)\lambda K\hat{\beta}(x)\mathbf{1}_n^T \mathbf{p} \\ \mathbf{p}^T A^{(1)T} & \hat{\psi}(x)\mathbf{p}^T W \mathbf{p} \end{bmatrix}, \\
 \hat{R}(x) &= \frac{1}{n\hat{\psi}(x)} \begin{bmatrix} X^T W \Gamma X + n\hat{\psi}(x)\lambda K & A^{(1)}\mathbf{1}_n/\hat{\psi}(x) \\ \mathbf{1}_n^T A^{(1)T}/\hat{\psi}(x) & -\hat{\psi}(x)\mathbf{q}^T W \mathbf{1}_n \end{bmatrix}, \\
 A^{(k)} &= X^T W \Lambda^k, \quad \Lambda = \text{diag}\{y_i - b'(\hat{\beta}(x)^T \mathbf{x}(x_i; x))\}, \\
 \Gamma &= \text{diag}\{b''(\hat{\beta}(x)^T \mathbf{x}(x_i; x))\}, \\
 \mathbf{p} &= (p_1, p_2, \dots, p_n)^T, \quad \mathbf{q} = (q_1, q_2, \dots, q_n)^T, \\
 \mathbf{1}_n &= (1, 1, \dots, 1)^T \quad (n \times 1), \\
 p_i &= -\frac{\hat{\beta}(x)^T \mathbf{x}(x_i; x)y_i - b(\hat{\beta}(x)^T \mathbf{x}(x_i; x))}{\hat{\psi}(x)^2} + \frac{\partial c(y_i, \psi(x))}{\partial \psi} \Bigg|_{\hat{\psi}(x)}, \\
 q_i &= \frac{2}{\hat{\psi}(x)^3} \{\hat{\beta}(x)^T \mathbf{x}(x_i; x)y_i - b(\hat{\beta}(x)^T \mathbf{x}(x_i; x))\} + \frac{\partial^2 c(y_i, \psi(x))}{\partial \psi^2} \Bigg|_{\hat{\psi}(x)}.
 \end{aligned}$$

We select a bandwidth h , the degree of polynomial p and a regularization parameter λ which minimize the information criterion (2.17).

2.3.1 *Gaussian model*

Using the equations (2.17) and (2.18), we can obtain the information criterion for the estimated Gaussian model $f_N(y_i | x_i; \hat{\beta}(x_i), \hat{\sigma}^2(x_i))$ in (2.13) as follows;

$$\begin{aligned}
 (2.19) \quad \text{GIC}_{h,p,\lambda} &= \sum_{i=1}^n \left[\log\{2\pi\hat{\sigma}^2(x_i)\} + \frac{\{y_i - \hat{\beta}_0(x_i)\}^2}{\hat{\sigma}^2(x_i)} \right] \\
 &\quad + \frac{2}{n} \sum_{i=1}^n \text{tr}\{\hat{R}(x_i)^{-1}\hat{Q}(x_i)\}, \\
 \hat{Q}(x) &= \frac{1}{n\hat{\sigma}^4(x)} \\
 &\quad \times \begin{bmatrix} A_N^{(2)}X - \lambda\hat{\sigma}^2(x)K\hat{\beta}(x)\mathbf{1}_n^T \Lambda_N X & \frac{1}{2\hat{\sigma}^2(x)}A_N^{(3)}\mathbf{1}_n - \frac{1}{2}A_N^{(1)}\mathbf{1}_n \\ \frac{1}{2\hat{\sigma}^2(x)}\mathbf{1}_n^T A_N^{(3)T} - \frac{1}{2}\mathbf{1}_n^T A_N^{(1)T} & \frac{1}{4\hat{\sigma}^4(x)}\mathbf{1}_n^T \Lambda_N^4 W \mathbf{1}_n - \frac{1}{4}\text{tr}(W) \end{bmatrix}, \\
 \hat{R}(x) &= \frac{1}{n\hat{\sigma}^2(x)} \begin{bmatrix} X^T W X + n\hat{\sigma}^2(x)\lambda K & \frac{1}{\hat{\sigma}^2(x)}A_N^{(1)}\mathbf{1}_n \\ \frac{1}{\hat{\sigma}^2(x)}\mathbf{1}_n^T A_N^{(1)T} & \frac{1}{2\hat{\sigma}^2(x)}\text{tr}(W) \end{bmatrix},
 \end{aligned}$$

where $A_N^{(k)} = X^T W \Lambda_N^k$ and $\Lambda_N = \text{diag}\{y_i - \hat{\beta}(x)^T \mathbf{x}(x_i; x)\}$.

2.3.2 *Logistic model*

In the case of the estimated logistic model $f_L(y_i | x_i, \hat{\beta}(x_i))$ in (2.16), we can obtain the information criterion as follows;

$$\begin{aligned}
 (2.20) \quad \text{GIC}_{h,p,\lambda} &= -2 \sum_{i=1}^n \log f_L(y_i | x_i, \hat{\beta}(x_i)) + \frac{2}{n} \sum_{i=1}^n \text{tr}\{\hat{R}(x_i)^{-1}\hat{Q}(x_i)\}, \\
 \hat{Q}(x) &= \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{\pi}(x_i; x)\}
 \end{aligned}$$

$$\begin{aligned} & \times [w_h(\mathbf{x}_i; \mathbf{x})\{y_i - \hat{\pi}(\mathbf{x}_i; \mathbf{x})\}\mathbf{x}(\mathbf{x}_i; \mathbf{x}) - \lambda K \hat{\boldsymbol{\beta}}(\mathbf{x})]\mathbf{x}(\mathbf{x}_i; \mathbf{x})^T, \\ \hat{R}(\mathbf{x}) = & -\frac{1}{n} \sum_{i=1}^n [w_h(\mathbf{x}_i; \mathbf{x})\hat{\pi}(\mathbf{x}_i; \mathbf{x})\{1 - \hat{\pi}(\mathbf{x}_i; \mathbf{x})\}\mathbf{x}(\mathbf{x}_i; \mathbf{x})\mathbf{x}(\mathbf{x}_i; \mathbf{x})^T - \lambda K], \end{aligned}$$

where $\hat{\pi}(\mathbf{x}_i; \mathbf{x}) = \exp\{\hat{\boldsymbol{\beta}}(\mathbf{x})^T \mathbf{x}(\mathbf{x}_i; \mathbf{x})\} / [1 + \exp\{\hat{\boldsymbol{\beta}}(\mathbf{x})^T \mathbf{x}(\mathbf{x}_i; \mathbf{x})\}]$.

3. Multivariate regularized local likelihood method

In this section, we describe the regularized local likelihood method for the case of d explanatory variables.

3.1 Model and estimation

Suppose we have n independent observations $\{(y_i, \mathbf{x}_i); i = 1, 2, \dots, n\}$, where y_i are random response variables and $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$ are d -dimensional explanatory variable vectors. It is assumed that the conditional distribution of Y , given $\mathbf{x} = (x_1, x_2, \dots, x_d)^T$, belongs to the exponential family which is given by replacing x with \mathbf{x} in (2.1). The parameters $\theta(\cdot)$ and $\psi(\cdot)$ are related to the conditional mean and conditional variance similarly to the single covariate case.

Assume that the function $\theta(\mathbf{x})$ is differentiable at a point \mathbf{x}_0 . Then $\theta(\mathbf{x})$ can be approximated locally by a polynomial of degree 1:

$$(3.1) \quad \theta(\mathbf{x}) \approx \theta(\mathbf{x}_0) + \frac{\partial \theta(\mathbf{x}_0)}{\partial \mathbf{x}^T} (\mathbf{x} - \mathbf{x}_0) = \boldsymbol{\beta}^*(\mathbf{x}_0)^T \mathbf{x}^*(\mathbf{x}; \mathbf{x}_0),$$

where $\boldsymbol{\beta}^*(\mathbf{x}_0) = (\beta_0^*(\mathbf{x}_0), \beta_1^*(\mathbf{x}_0), \dots, \beta_d^*(\mathbf{x}_0))^T$, $\beta_0^*(\mathbf{x}_0) = \theta(\mathbf{x}_0)$, $\beta_j^*(\mathbf{x}_0) = \partial \theta(\mathbf{x}_0) / \partial x_j$ ($j = 1, 2, \dots, d$) and $\mathbf{x}^*(\mathbf{x}; \mathbf{x}_0) = (1, (\mathbf{x} - \mathbf{x}_0)^T)^T$. More generally, approximation (3.1) can be expanded to the higher order term (see Ruppert and Wand (1994)). But we use approximation (3.1) in order not to increase the number of unknown parameters.

Then the data $\{y_1, y_2, \dots, y_n\}$ were summarized by a model from a class of probability densities $f(y_i | \mathbf{x}_i; \boldsymbol{\beta}^*(\mathbf{x}_0), \psi(\mathbf{x}_0))$ which are expressed by (2.5). We propose choosing unknown parameters $\boldsymbol{\beta}^*(\mathbf{x}_0)$ and $\psi(\mathbf{x}_0)$ to maximize the regularized local log-likelihood function

$$(3.2) \quad \begin{aligned} & RL(\boldsymbol{\beta}^*(\mathbf{x}_0), \psi(\mathbf{x}_0); H, \lambda) \\ & = \sum_{i=1}^n w_H(\mathbf{x}_i; \mathbf{x}_0) \left\{ \frac{\boldsymbol{\beta}^*(\mathbf{x}_0)^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}_0) y_i - b(\boldsymbol{\beta}^*(\mathbf{x}_0)^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}_0))}{\psi(\mathbf{x}_0)} \right. \\ & \quad \left. + c(y_i, \psi(\mathbf{x}_0)) \right\} \\ & \quad - n\lambda \boldsymbol{\beta}^*(\mathbf{x}_0)^T K \boldsymbol{\beta}^*(\mathbf{x}_0) / 2. \end{aligned}$$

The weight function $w_H(\mathbf{x}_i; \mathbf{x}_0)$ is the Gaussian product kernel:

$$(3.3) \quad w_H(\mathbf{x}; \mathbf{x}_0) = \frac{K_d(H^{-1}(\mathbf{x} - \mathbf{x}_0))}{|\det(H)|}, \quad K_d(\mathbf{u}) = (2\pi)^{-d/2} \prod_{j=1}^d \exp\left(-\frac{1}{2}u_j^2\right),$$

where a matrix H is a $p \times p$ nonsingular bandwidth matrix. We use $H = hI_d$, where I_d is a d -dimensional identity matrix and h is a bandwidth. The regularization term is

given by equation:

$$(3.4) \quad \beta^*(\mathbf{x}_0)^T K \beta^*(\mathbf{x}_0) = \beta^*(\mathbf{x}_0)^T \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & I_d \end{bmatrix} \beta^*(\mathbf{x}_0),$$

where $\mathbf{0}_d$ is a d -dimensional zero vector.

The estimators $\hat{\beta}^*(\mathbf{x}_0)$ and $\hat{\psi}(\mathbf{x}_0)$ are obtained by the iterative algorithm. Replacing the unknown parameters $\beta^*(\mathbf{x}_0)$ and $\psi(\mathbf{x}_0)$ by $\hat{\beta}^*(\mathbf{x}_0)$ and $\hat{\psi}(\mathbf{x}_0)$ respectively and noting that the parameter $\theta(\mathbf{x}_i)$ in the exponential family $f(y_i | \mathbf{x}_i; \beta^*(\mathbf{x}_0), \psi(\mathbf{x}_0))$ is directly given by $\theta(\mathbf{x}_i) = \beta^*(\mathbf{x}_i)^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}_i) = \beta_0^*(\mathbf{x}_i)$, we obtain the estimated model:

$$(3.5) \quad f(y_i | \mathbf{x}_i, \hat{\beta}^*(\mathbf{x}_i), \hat{\psi}(\mathbf{x}_i)) = \exp \left\{ \frac{\hat{\beta}_0^*(\mathbf{x}_i) y_i - b(\hat{\beta}_0^*(\mathbf{x}_i))}{\hat{\psi}(\mathbf{x}_i)} + c(y_i, \hat{\psi}(\mathbf{x}_i)) \right\}.$$

Example 3.1. We illustrate the proposed regularized local likelihood modeling by fitting surface to the simulation data. The random samples $\{(y_i, \mathbf{x}_i); i = 1, 2, \dots, 300\}$ are generated from the true model $y_i = \sin(5\pi x_{i1}) + \cos(2\pi x_{i2}) + \varepsilon_i$, $\varepsilon_i \sim N(0, 0.1^2)$, where the design points x_{i1} and x_{i2} in $\mathbf{x}_i = (x_{i1}, x_{i2})^T$ are uniformly distributed in $[-1, 1]$. Figure 2 shows the true surface, where x_1, x_2 in $\mathbf{x} = (x_1, x_2)^T$ and y are expressed as X-, Y- and Z-axis respectively. We apply the local likelihood method and the regularized local likelihood method to the simulation data respectively, to examine the effectiveness of the regularization parameter $\zeta = \lambda\sigma^2$, where we use the Gaussian model. Figure 3(a) shows the estimated surface for the bandwidth matrix $H = h \times I_2$ ($h = 0.07$) without the help of regularization ($\lambda = \zeta = 0$). This surface is obviously undersmoothed in the boundary region. We use a large bandwidth $h = 0.1$ and the corresponding estimated surface is given in Fig. 3(b), but the fitting in the boundary region cannot be improved. This result implies that the estimated surface is unstable in the boundary region. Figure 3(c) shows the estimated surface for the bandwidth $h = 0.07$ with the help of regularization ($\zeta = 10^{-3}$). This surface gives a good representation of the underlying function over the boundary region. Figure 3(d) shows the four curves for the true and estimated surfaces in the boundary region ($x_2 = -1$), where the solid line in Fig. 3(d) is the true model, the broken line in Fig. 3(d) is the local likelihood model ($h = 0.07$), the dotted line in Fig. 3(d) is the local likelihood model ($h = 0.1$) and the

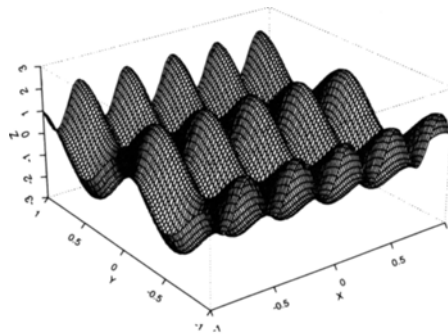


Fig. 2. True surface $y = \sin(5\pi x_1) + \cos(2\pi x_2)$.

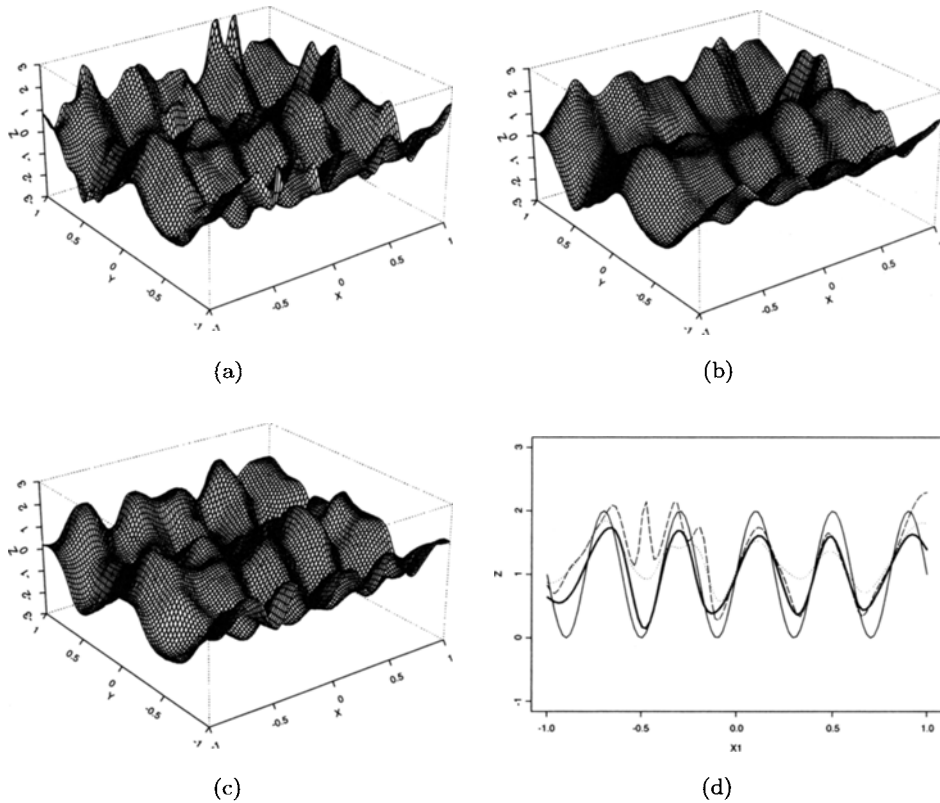


Fig. 3. Comparison of the estimated surfaces. (a) and (b) show the estimated surfaces for the local likelihood method with $h = 0.07$ and $h = 0.1$ respectively. (c) shows the estimated surface for the regularized local likelihood method with $h = 0.07$ and $\zeta = 10^{-3}$. (d) shows the four curves for the true and estimated surfaces in the boundary region ($x_2 = -1$).

thick line in Fig. 3(d) is the regularized local likelihood model ($h = 0.07, \zeta = 10^{-3}$). In this figure, the thick line gives the most appropriate estimate of the true model since this gives a stable sine curve.

In Section 2, the local likelihood method may cause the instability of estimator for large degrees of polynomial p . In multivariate case, this issue may be occurred for large number of explanatory variables d .

3.2 Model selection

The crucial issue is how to choose a bandwidth matrix H and a regularization parameter λ . We derive an information criterion to evaluate the estimated model in (3.5). Similarly to the way in Subsection 2.3, we obtain the information criterion based on the multivariate regularized local likelihood method by

$$(3.6) \quad \text{GIC}_{H,\lambda} = -2 \sum_{i=1}^n \left\{ \frac{\hat{\beta}_0^*(\mathbf{x}_i) y_i - b(\hat{\beta}_0^*(\mathbf{x}_i))}{\hat{\psi}(\mathbf{x}_i)} + c(y_i, \hat{\psi}(\mathbf{x}_i)) \right\}$$

$$+ \frac{2}{n} \sum_{i=1}^n \text{tr}\{\hat{R}(\mathbf{x}_i)^{-1} \hat{Q}(\mathbf{x}_i)\},$$

where

$$\hat{Q}(\mathbf{x}) = \frac{1}{n\hat{\psi}(\mathbf{x})} \begin{bmatrix} A^{(2)} \mathbf{X} / \hat{\psi}(\mathbf{x}) - \lambda K \hat{\boldsymbol{\beta}}^*(\mathbf{x}) \mathbf{1}_n^T \Lambda \mathbf{X} & A^{(1)} \mathbf{p} - \hat{\psi}(\mathbf{x}) \lambda K \hat{\boldsymbol{\beta}}^*(\mathbf{x}) \mathbf{1}_n^T \mathbf{p} \\ \mathbf{p}^T A^{(1)T} & \hat{\psi}(\mathbf{x}) \mathbf{p}^T \mathbf{W} \mathbf{p} \end{bmatrix},$$

$$\hat{R}(\mathbf{x}) = \frac{1}{n\hat{\psi}(\mathbf{x})} \begin{bmatrix} \mathbf{X}^T \mathbf{W} \Gamma \mathbf{X} + n\hat{\psi}(\mathbf{x}) \lambda K & A^{(1)} \mathbf{1}_n / \hat{\psi}(\mathbf{x}) \\ \mathbf{1}_n^T A^{(1)T} / \hat{\psi}(\mathbf{x}) & -\hat{\psi}(\mathbf{x}) \mathbf{q}^T \mathbf{W} \mathbf{1}_n \end{bmatrix},$$

$$A^{(k)} = \mathbf{X}^T \mathbf{W} \Lambda^k, \quad \mathbf{X} = (\mathbf{x}^*(\mathbf{x}_1; \mathbf{x}), \dots, \mathbf{x}^*(\mathbf{x}_n; \mathbf{x}))^T, \quad \mathbf{W} = \text{diag}\{w_H(\mathbf{x}_i; \mathbf{x})\}.$$

The notations Λ , Γ , \mathbf{p} and \mathbf{q} are given by (2.18), where we use $\hat{\boldsymbol{\beta}}^*(\mathbf{x})^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x})$ and $\hat{\psi}(\mathbf{x})$ instead of $\hat{\boldsymbol{\beta}}(\mathbf{x})^T \mathbf{x}_i$ and $\hat{\psi}(\mathbf{x}_i)$ respectively. We select a bandwidth matrix H and a regularization parameter λ which minimize the information criterion (3.6).

We give the information criteria in the case of Gaussian and logistic models.

3.2.1 Gaussian model

We use the Gaussian regression model $f_N(y_i | \mathbf{x}_i; \hat{\boldsymbol{\beta}}^*(\mathbf{x}_i), \hat{\sigma}^2(\mathbf{x}_i))$ in (2.13), where the estimators $\hat{\boldsymbol{\beta}}^*(\mathbf{x}_0)$ and $\hat{\sigma}^2(\mathbf{x}_0)$ are given by

$$(3.7) \quad \hat{\boldsymbol{\beta}}^*(\mathbf{x}_0) = (\mathbf{X}^T \mathbf{W} \mathbf{X} + n\zeta Q)^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y},$$

$$\hat{\sigma}^2(\mathbf{x}_0) = \{\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}^*(\mathbf{x}_0)\}^T \mathbf{W} \{\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}^*(\mathbf{x}_0)\} / \text{tr}(\mathbf{W}).$$

Then we obtain the information criterion for the estimated Gaussian model as follows;

$$(3.8) \quad \text{GIC}_{H,\lambda} = \sum_{i=1}^n \left[\log\{2\pi\hat{\sigma}^2(\mathbf{x}_i)\} + \frac{\{y_i - \hat{\beta}_0^*(\mathbf{x}_i)\}^2}{\hat{\sigma}^2(\mathbf{x}_i)} \right] + \frac{2}{n} \sum_{i=1}^n \text{tr}\{\hat{R}(\mathbf{x}_i)^{-1} \hat{Q}(\mathbf{x}_i)\},$$

$$\hat{Q}(\mathbf{x}) = \frac{1}{n\hat{\sigma}^4(\mathbf{x})} \times \begin{bmatrix} \frac{1}{\hat{\sigma}^2(\mathbf{x})} A_N^{(2)} \mathbf{X} - \lambda \hat{\sigma}^2(\mathbf{x}) K \hat{\boldsymbol{\beta}}^*(\mathbf{x}) \mathbf{1}_n^T \Lambda_N \mathbf{X} & \frac{1}{2\hat{\sigma}^2(\mathbf{x})} A_N^{(3)} \mathbf{1}_n - \frac{1}{2} A_N^{(1)} \mathbf{1}_n \\ \frac{1}{2\hat{\sigma}^2(\mathbf{x})} \mathbf{1}_n^T A_N^{(3)T} - \frac{1}{2} \mathbf{1}_n^T A_N^{(1)T} & \frac{1}{4\hat{\sigma}^4(\mathbf{x})} \mathbf{1}_n^T \Lambda_N^4 \mathbf{W} \mathbf{1}_n - \frac{1}{4} \text{tr}(\mathbf{W}) \end{bmatrix},$$

$$\hat{R}(\mathbf{x}) = \frac{1}{n\hat{\sigma}^2(\mathbf{x})} \begin{bmatrix} \mathbf{X}^T \mathbf{W} \mathbf{X} + n\hat{\sigma}^2(\mathbf{x}) \lambda K & \frac{1}{\hat{\sigma}^2(\mathbf{x})} A_N^{(1)} \mathbf{1}_n \\ \frac{1}{\hat{\sigma}^2(\mathbf{x})} \mathbf{1}_n^T A_N^{(1)T} & \frac{1}{2\hat{\sigma}^2(\mathbf{x})} \text{tr}(\mathbf{W}) \end{bmatrix},$$

where $A_N^{(k)} = \mathbf{X}^T \mathbf{W} \Lambda_N^k$ and $\Lambda_N = \text{diag}\{y_i - \hat{\boldsymbol{\beta}}^*(\mathbf{x})^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x})\}$.

3.2.2 Logistic model

In the case of logistic model, we can obtain the information criterion as follows

$$(3.9) \quad \text{GIC}_{H,\lambda} = -2 \sum_{i=1}^n \log f_L(y_i | \mathbf{x}_i, \hat{\boldsymbol{\beta}}^*(\mathbf{x}_i)) + \frac{2}{n} \sum_{i=1}^n \text{tr}\{\hat{R}(\mathbf{x}_i)^{-1} \hat{Q}(\mathbf{x}_i)\},$$

$$\hat{Q}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{\pi}(\mathbf{x}_i; \mathbf{x})\} \times [w_H(\mathbf{x}_i; \mathbf{x}) \{y_i - \hat{\pi}(\mathbf{x}_i; \mathbf{x})\} \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}) - \lambda K \hat{\boldsymbol{\beta}}^*(\mathbf{x})] \mathbf{x}^*(\mathbf{x}_i; \mathbf{x})^T,$$

$$\hat{R}(\mathbf{x}) = -\frac{1}{n} \sum_{i=1}^n [w_H(\mathbf{x}_i; \mathbf{x}) \hat{\pi}(\mathbf{x}_i; \mathbf{x}) \{1 - \hat{\pi}(\mathbf{x}_i; \mathbf{x})\} \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}) \mathbf{x}^*(\mathbf{x}_i; \mathbf{x})^T - \lambda K],$$

for the estimator $\hat{\beta}^*(\mathbf{x})$ which maximizes the regularized local likelihood function (3.2), where $\hat{\pi}(\mathbf{x}_i; \mathbf{x}) = \exp\{\hat{\beta}^*(\mathbf{x})^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x})\} / [1 + \exp\{\hat{\beta}^*(\mathbf{x})^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x})\}]$ and $f_L(y_i | \mathbf{x}_i, \hat{\beta}^*(\mathbf{x}_i)) = \hat{\pi}(\mathbf{x}_i; \mathbf{x}_i)^{y_i} \{1 - \hat{\pi}(\mathbf{x}_i; \mathbf{x}_i)\}^{1-y_i}$.

4. Real examples and numerical results

In this section we use a real data example and Monte Carlo simulations to investigate the performance of the regularized local likelihood modeling.

4.1 Summer rainfall data

We apply the proposed modeling procedure to the rainfall data in Kagoshima. In general, investigators in national meteorological observatories observe the sky at fixed times. For example, when they observe the clouds, they look out over the sky and observe the rate of the clouds which is called “cloud amount”. They classify the weather condition as “Blue Sky” or “Cloudy”. However, a weather condition such as “Rain” does not relate to cloud amount directly. Therefore, we investigate the relationship between cloud amount and “Rain” using our nonlinear modeling procedure.

We use the data $\{(y_i, \mathbf{x}_i); i = 1, 2, \dots, 92\}$, where y_i are the binary response variables having the value 1 (if the weather is “Rain”) or 0 (otherwise) and $\mathbf{x}_i \in [0, 10]$ are the average daily cloud amounts. These variables were observed in Kagoshima from June, 2001 to August, 2001. We apply our nonlinear logistic regression model in (2.16) to this data and select the bandwidth h , the degree of polynomial p and the regularization parameter λ by minimizing GIC in (2.20).

Figure 4(a) shows the minimum GIC with respect to the bandwidth h , where the degree of polynomial $p = 3$ and the regularization parameter $\lambda = 1.00 \times 10^{-7}$. We observe that the optimal bandwidth is $\hat{h} = 2.5$, and show the estimated curve in Fig. 4(b). In general, we tend to predict the probability of precipitation that is higher as the cloud amount increases. The estimated curve is not monotone increasing, since the probability of precipitation is influenced by cloud types in addition to the cloud amount.

We examine the relationship among cloud amount, humidity and “Rain” using our nonlinear regression modeling. We use the data $\{(y_i, \mathbf{x}_i); i = 1, 2, \dots, 92\}$, where $\mathbf{x}_i = (x_{i1}, x_{i2})$, x_{i1} are the average daily cloud amounts and x_{i2} are the average daily values

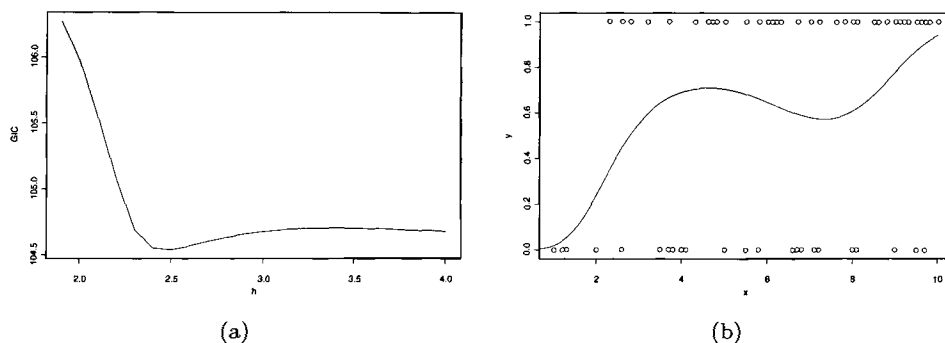


Fig. 4. (a) The relationship between the bandwidth h and GIC with the number of polynomial $p = 3$ and the regularization parameter $\lambda = 1.00 \times 10^{-7}$. (b) The smoothed curve based on the regularized local likelihood model and GIC ($\hat{h} = 2.5$, $\hat{\lambda} = 1.00 \times 10^{-7}$).

of humidity. The data were observed in Kagoshima from June, 2001 to August, 2001, respectively. We apply our nonlinear logistic regression model to this data and select the bandwidth matrix H and the regularization parameter λ by minimizing GIC in (3.9). We use $H = hI_2$ (I_2 is 2×2 identity matrix) as the bandwidth matrix.

Figure 5(a) shows the minimum GIC with respect to the regularization parameter λ , with the optimal bandwidth. We select the optimal bandwidth $\hat{h} = 0.11$ and regularization parameter $\hat{\lambda} = 3.16 \times 10^{-5}$, and show the corresponding estimated surface in Fig. 5(b). We observe that the probability of precipitation is influenced by the humidity rather than the cloud amount, however the relationship between the humidity and the probability of precipitation is not monotonic. We conclude that the other factors, such as the atmospheric pressure and the sort of clouds, may influence the probability of precipitation.

Figure 6 shows the result of the nonlinear logistic regression model for the data $\{(y_i, \mathbf{x}_i); i = 1, 2, \dots, 92\}$ observed in Niigata from June, 2001 to August, 2001. Figure 6(a) shows the minimum GIC with respect to the regularization parameter λ for the optimal bandwidth, and (b) shows the estimated surface ($\hat{h} = 0.1$, $\hat{\lambda} = 1.00 \times 10^{-4}$)

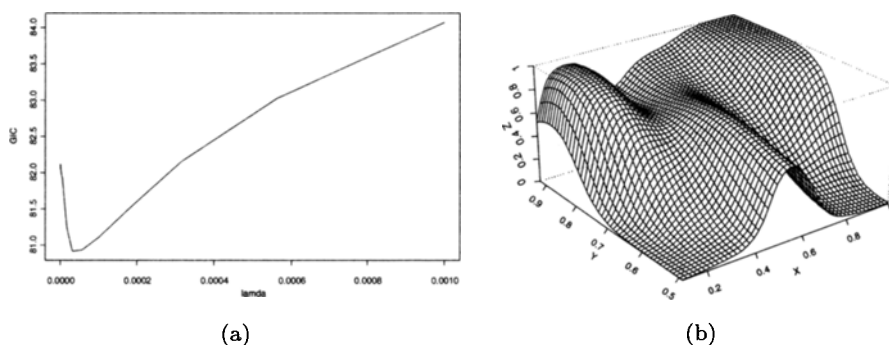


Fig. 5. (a) The relationship between the regularization parameter λ and GIC with the optimal bandwidth h . (b) The smoothed surface based on the regularized local likelihood model and GIC ($\hat{h} = 0.11$, $\hat{\lambda} = 3.16 \times 10^{-5}$).

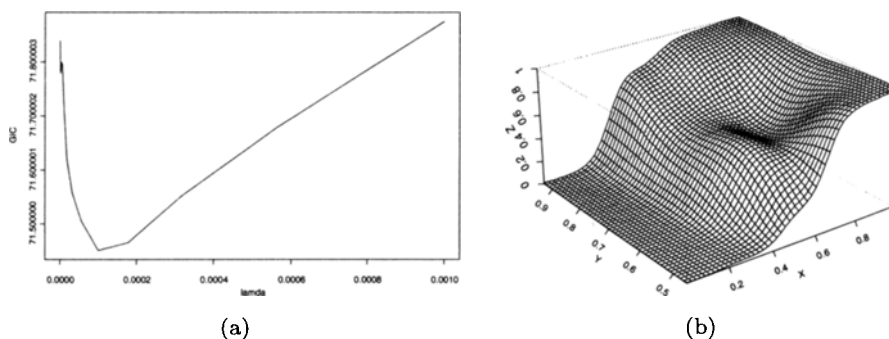


Fig. 6. (a) The relationship between the regularization parameter λ and GIC with the optimal bandwidth h (in Niigata). (b) The smoothed surface based on the regularized local likelihood model and GIC ($\hat{h} = 0.1$, $\hat{\lambda} = 1.00 \times 10^{-4}$).

based on the regularized local likelihood method and GIC. It is clear that the estimated surface varies according to the locality where the data were observed, by comparing Fig. 5(b) with Fig. 6(b). This difference of the estimated surfaces may be caused by the property of summer climate in each place. In this example, Kagoshima is located in the Pacific coast, while Niigata is located in the coast of the Sea of Japan. The probability of precipitation in Kagoshima is related to the humidity because of the large influence of rain shower and typhoon. Conversely, the probability of precipitation in Niigata is related to the cloud amount because of the small influence of rain shower and typhoon.

4.2 Monte Carlo simulations

Monte Carlo simulations were conducted to investigate the performance of the nonlinear modeling strategy based on the regularized local likelihood method. For the simulation study, repeated random samples $\{(y_i, \mathbf{x}_i); i = 1, 2, \dots, n\}$ were generated from the true regression model $y_i = m(\mathbf{x}_i) + \varepsilon_i$, where $\mathbf{x}_i = (x_{i1}, x_{i2})^T$ and the design points x_{i1}, x_{i2} are uniformly distributed in $[-1, 1]$. The errors ε_i are assumed to be independently distributed according to a mixture of normal distributions $rN(0, (0.5R_y)^2) + (1 - r)N(0, (0.1R_y)^2)$, where R_y is the range of $m(\mathbf{x})$ over \mathbf{x} and $r = 0.1, 0.5$. For the true curve $m(\mathbf{x})$ we consider the following function: $m(\mathbf{x}_i) = \sin(2\pi x_{i1})/2 + \cos(4\pi x_{i2})/2$.

We fit the nonlinear regression model with Gaussian noise defined by (2.13) to the simulated data. The model is estimated by the maximization of the local log-likelihood function and the regularized local log-likelihood function (2.6). The values of the bandwidth matrix $H = hI_2$ and the regularization parameter λ are chosen as the minimizers of the criteria GIC and Modified AIC given by

$$\text{MAIC}_{H,\lambda} = 2 \sum_{i=1}^n \left[\log(2\pi\hat{\sigma}^2(\mathbf{x}_i)) + \frac{\{y_i - \hat{\beta}^*(\mathbf{x}_i)^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}_i)\}^2}{\hat{\sigma}^2(\mathbf{x}_i)} \right] + 2 \text{tr}(S),$$

where S is the $n \times n$ smoother matrix (Hastie and Tibshirani (1990)) given by $(\mathbf{s}(\mathbf{x}_1), \dots, \mathbf{s}(\mathbf{x}_n))^T$ with $\mathbf{s}(\mathbf{x}_i) = [\mathbf{e}_i^T \mathbf{X} \{ \mathbf{X}^T \mathbf{W} \mathbf{X} + n\zeta \mathbf{Q} \}^{-1} \mathbf{X}^T \mathbf{W}]^T$ and \mathbf{e}_i are the n -dimensional vectors having the value 1 in the i -th entry and zero elsewhere. The cross-validation (CV) and the generalized cross-validation (GCV) were also used for the choice of the adjusted parameters. These criteria were examined by comparing the mean squared error:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n \{ \hat{\beta}^*(\mathbf{x}_i)^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}_i) - m(\mathbf{x}_i) \}^2$$

and the standard deviations (SD) of the selected h and λ .

Tables 1 and 2 compare the MSE between the true and estimated functions, in which the fitted functions are obtained by averaging over 100 repeated Monte Carlo trials so we set this as AMSE; h and λ in each table are the averages of optimal h and λ , and (SD) is the standard deviation of h and λ . In the case where $n = 100$ and $r = 0.1$ in Table 1, RLLM & GIC and RLLM & GCV are superior to other methods in the sense of decreasing MSE. In the other cases where $r = 0.5$ in Table 2, LLM & GIC and RLLM & GCV are superior to other methods in that sense. Moreover, the selected models based on GIC are more stable since the optimal h has smaller variance than the use of CV and GCV. This implies that the proposed modeling procedure gives the fitted functions that capture the true structure in practical applications.

Table 1. The result of the surface when $r = 0.1$.

LLM	GIC	MAIC	CV	GCV
h	0.119	0.063	0.126	0.096
(SD)	0.015	0.008	0.034	0.014
AMSE $\times 10^4$	395	331	454	314
(SD)	1.03×10^{-2}	5.72×10^{-3}	2.28×10^{-2}	5.95×10^{-3}
RLLM	GIC	MAIC	CV	GCV
h	0.076	0.054	0.088	0.080
(SD)	0.009	0.002	0.024	0.014
λ	4.58×10^{-1}	8.40×10^{-3}	3.66×10^{-2}	8.56×10^{-2}
(SD)	5.05×10^{-2}	2.51×10^{-2}	7.46×10^{-2}	1.08×10^{-1}
AMSE $\times 10^4$	291	355	299	292
(SD)	4.02×10^{-3}	5.36×10^{-3}	5.57×10^{-3}	4.98×10^{-3}

Table 2. The result of the surface when $r = 0.5$.

LLM	GIC	MAIC	CV	GCV
h	0.123	0.064	0.205	0.170
(SD)	0.017	0.007	0.065	0.055
AMSE $\times 10^4$	1183	1971	1348	1283
(SD)	2.21×10^{-2}	3.34×10^{-2}	2.51×10^{-2}	2.84×10^{-2}
RLLM	GIC	MAIC	CV	GCV
h	0.084	0.051	0.154	0.134
(SD)	0.014	0.002	0.062	0.042
λ	3.53×10^{-2}	7.22×10^{-5}	1.89×10^{-2}	2.20×10^{-2}
(SD)	7.65×10^{-3}	1.50×10^{-4}	1.59×10^{-2}	1.49×10^{-2}
AMSE $\times 10^4$	1232	2185	1241	1186
(SD)	2.34×10^{-2}	3.45×10^{-2}	3.03×10^{-2}	3.24×10^{-2}

Figure 7 shows the comparison between GIC and MAIC, where we use the regularization parameter $\zeta = 10^{-5}$. Figure 7(a) indicates the relationship between the bandwidth h and the log-likelihood term given by

$$2 \sum_{i=1}^n \left[\log(2\pi\hat{\sigma}^2(\mathbf{x}_i)) + \frac{\{y_i - \hat{\beta}^*(\mathbf{x}_i)^T \mathbf{x}^*(\mathbf{x}_i; \mathbf{x}_i)\}^2}{\hat{\sigma}^2(\mathbf{x}_i)} \right],$$

which is a monotone increasing function with respect to the bandwidth h . Figure 7(b) shows the comparison between GIC and MAIC with respect to the corrected bias terms, where the dotted line indicates the bias of MAIC and the solid line the bias of GIC. The biases present a monotone decrease function with respect to the bandwidth h . Figure 7(c) shows the results for GIC and MAIC. In the case of MAIC, the optimal bandwidth tends to be too small to select as shown in Fig. 7(c), since the estimator may be incomputable for smaller values of h . On the contrary, GIC always selects the optimal bandwidth as in Fig. 7(c), and has smaller variance than the cases of CV and GCV.

Nonaka *et al.* (2003) examined the efficiency of our proposed method through a Monte Carlo simulation for the one-dimensional explanatory variable x .

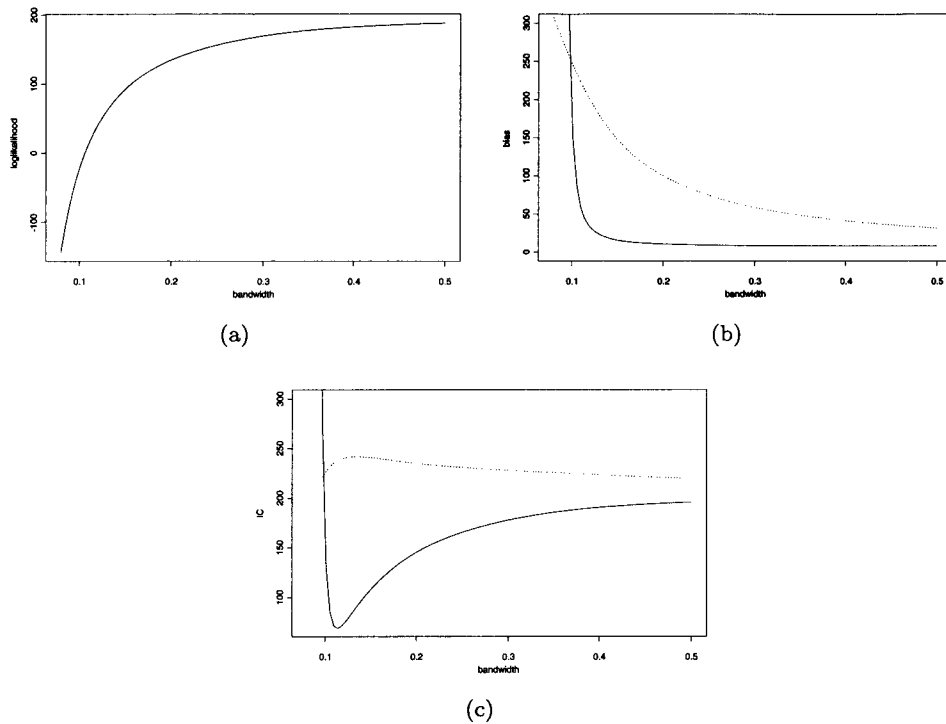


Fig. 7. (a) The relationship between the bandwidth h and the log-likelihood term with the regularization parameter $\zeta = 10^{-5}$. (b) The relationship between the bandwidth h and the corrected bias term. The dotted line expresses the bias of MAIC and the solid line expresses the bias of GIC. (c) The relationship between the bandwidth h and the information criterion (IC). The dotted line expresses the value of MAIC and the solid line expresses the value of GIC.

5. Concluding remarks

The main aim of this paper was to introduce the regularized local likelihood method in constructing a nonlinear regression model; determining a set of kernel function with bandwidth, estimating the unknown parameters by regularization and then evaluating the constructed model to select a suitable one. The estimated curve based on the local likelihood method tends to be unstable for a small bandwidth h and a higher degree of polynomial p . In particular, the estimated surface can not be calculated for a small $H = hI_d$ in the multivariate case. We proposed the nonlinear regression modeling procedure based on the regularized local likelihood method in order to obtain a stable estimator, and derived a model selection criterion for evaluating constructed models from an information-theoretic point of view.

We applied the regularized local likelihood method to summer rainfall data and simulated data. We observed that our method is effective in constructing nonlinear regression models for the multivariate data, and that the proposed strategy using the information criterion GIC yields stable parameter estimates. It may be applied to construct Gaussian, logistic and Poisson nonlinear regression models, and provides a tool to draw information about the system under consideration from a finite and noisy data

set. We would recommend implementing nonlinear regression modeling based on the regularized local likelihood method, using the information criterion GIC.

Acknowledgements

The authors would like to thank the referees and the Editor for their helpful comments and suggestions.

Appendix: Derivation of the information criterion

We derive an information criterion to evaluate models estimated by the regularized local likelihood method.

Suppose that z_1, z_2, \dots, z_n are future observations for the response variable Y drawn from $g(y | x)$. Let $f(z | X; \hat{\theta}(X)) = \prod_{i=1}^n f(z_i | x_i; \hat{\beta}(x_i), \hat{\psi}(x_i))$ and $g(z | X) = \prod_{i=1}^n g(z_i | x_i)$. An information criterion may be derived as an estimator of the Kullback-Leibler information (Kullback and Leibler (1951))

$$(A.1) \quad KL\{g, f\} = E_{G(z|X)}[\log g(z | X)] - E_{G(z|X)}[\log f(z | X; \hat{\theta}(X))]$$

conditional on $\hat{\theta}(X) = (\hat{\beta}(X)^T, \hat{\psi}(X)^T)^T$.

The first term in the right-hand side of equation (A.1) is constant over all models and only the second term

$$(A.2) \quad E_{G(z|X)}[\log f(z | X; \hat{\theta}(x_0))] = \int \log f(z | x; \hat{\theta}(x_0)) dG(z | x),$$

is relevant. Hence, instead of minimizing the Kullback-Leibler information (A.1), we maximize the expected log-likelihood (A.2) that depends on the unknown true distribution $G(z | X)$. An estimate of the expected log-likelihood is the log-likelihood

$$(A.3) \quad \sum_{i=1}^n \log f(y_i | x_i; \hat{\theta}(x_0)),$$

obtained by replacing the unknown distribution $G(z | X)$ by the empirical distribution. Then the bias of the log-likelihood in estimating the expected log-likelihood is given by

$$b(G) = E_{G(y|X)}[\log f(y | X; \hat{\theta}(x_0)) - E_{G(z|X)}[\log f(z | X; \hat{\theta}(x_0))]].$$

Konishi and Kitagawa (1996) considered an asymptotic bias for a statistical model with functional estimators and gave the bias by a function of the empirical influence function of estimators and the score function of a specified parametric model. It may be seen that the regularized local likelihood estimator $\hat{\theta}(x_0) = (\hat{\beta}(x_0)^T, \hat{\psi}(x_0)^T)^T$ can be expressed as $\hat{\theta}(x_0) = T(\hat{G})$ for the functional $T(\cdot)$ defined by

$$(A.4) \quad \int \frac{\partial}{\partial \theta} \{w_h(x; x_0) \log f(z | x; \theta(x_0)) - \lambda \beta(x_0)^T K \beta(x_0) / 2\} \Big|_{T(G)} dG(z) = \mathbf{0},$$

where G and \hat{G} are respectively the joint distribution of (x, y) and the empirical distribution function based on the observed data. Replacing G in (A.4) by $G_\epsilon = (1-\epsilon)G + \epsilon\delta_{(y,x)}$

with $\delta_{(y,x)}$ being a point of mass at (y, x) and differentiating with respect to ε yield the influence function of the regularized estimator $\hat{\theta}(x_0) = T(\hat{G})$ in the form

$$(A.5) \quad T^{(1)}(z | x; G) = R(G)^{-1} \frac{\partial}{\partial \theta} \{w_h(x; x_0) \log f(z | x; \theta(x_0)) - \lambda \beta(x_0)^T K \beta(x_0) / 2\} \Big|_{T(G)},$$

where

$$(A.6) \quad R(G) = - \int \left\{ \frac{\partial^2 A(z | x; \theta(x_0))}{\partial \theta \partial \theta^T} \right\} \Big|_{T(G)} dG(z),$$

with $A(z | x; \theta(x_0)) = w_h(x; x_0) \log f(z | x; \theta(x_0)) - \lambda \beta(x_0)^T K \beta(x_0) / 2$.

It follows from Theorem 2.1 in Konishi and Kitagawa (1996) (see also Konishi (1999)) that the bias is asymptotically given by $b(G) = \text{tr}\{R(G)^{-1}Q(G)\} + o(1/n)$, where

$$Q(G) = \int \left\{ \frac{\partial A(z | x; \theta(x_0))}{\partial \theta} \frac{\partial \log f(z | x; \theta(x_0))}{\partial \theta^T} \right\} \Big|_{T(G)} dG(z).$$

By replacing the unknown distribution G by the empirical distribution \hat{G} , we have an information criterion

$$(A.7) \quad \text{GIC}_{h,p,\lambda}(x_0) = -2 \sum_{i=1}^n \log f(y_i | x_i; \hat{\theta}(x_0)) + 2 \text{tr}\{R(\hat{G})^{-1}Q(\hat{G})\}.$$

where

$$Q(\hat{G}) = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\partial A(y_i | x_i; \theta(x_0))}{\partial \theta} \frac{\partial \log f(y_i | x_i; \theta(x_0))}{\partial \theta^T} \right\} \Big|_{\hat{\theta}(x_0)},$$

$$R(\hat{G}) = -\frac{1}{n} \sum_{i=1}^n \left\{ \frac{\partial^2 A(y_i | x_i; \theta(x_0))}{\partial \theta \partial \theta^T} \right\} \Big|_{\hat{\theta}(x_0)}.$$

It might be noticed here that the information criterion based on the local method has two problems: the information criterion $\text{GIC}_{h,p,\lambda}(x_0)$ in (A.7) depends on the point x_0 , and this method assesses the closeness of $f(y | x; \hat{\theta}(x_0))$ to the model $g(y | x)$ for a fixed point x_0 . In order to assess the closeness of $f(y | x; \theta(x))$ to the model $g(y | x)$, we modify the information criterion (A.7) in the following:

$$(A.8) \quad \text{GIC}_{h,p,\lambda} = -2 \sum_{i=1}^n \log f(y_i | x_i; \hat{\theta}(x_i)) + \frac{2}{n} \sum_{i=1}^n \text{tr}\{\hat{R}(x_i)^{-1}\hat{Q}(x_i)\},$$

where we replace $Q(\hat{G})$ and $R(\hat{G})$ with $\hat{Q}(x_0)$ and $\hat{R}(x_0)$, respectively. For the problem of choosing among different models, we select the model for which the value of the information criterion $\text{GIC}_{h,p,\lambda}$ is smallest.

Irizarry (2001) proposed the use of a weighted version of the Kullback-Leibler information and derived the model selection criterion WAIC. That assesses the closeness of $f(y | x; \hat{\theta}(x_0))$ to the model $g(y | x)$ for a fixed point x_0 .

REFERENCES

- Bartlett, M. S. (1954). A note on some multiplying factors for various χ approximations, *Journal of the Royal Statistical Society, Series, B*, **16**, 296–298.
- Cleveland, W. S. (1979). Robust locally weighted regression and smoothing scatterplots, *Journal of the American Statistical Association*, **74**, 829–836.
- Copas, B. J. and Eguchi, S. (1998). Sensitivity approximations for selectivity bias in observational data analysis, Research Memo., No. 660, The Institute of Statistical Mathematics, Tokyo.
- Eguchi, S. and Copas, J. (1998). A class of local likelihood methods and near-parametric asymptotics, *Journal of the Royal Statistical Society, Series, B*, **60**, 709–724.
- Eguchi, S. and Kim, T. Y. (2001). Local likelihood method and theory for a bridge of parametric and nonparametric regression, Research Memo., No. 806, The Institute of Statistical Mathematics, Tokyo.
- Eguchi, S., Kim, T. Y. and Park, B. U. (2003). Local likelihood method: A bridge over parametric and nonparametric regression, *Journal of Nonparametric Statistics*, **15**, 665–683.
- Fan, J. and Gijbels, I. (1996). *Local Polynomial Modelling and Its Applications*, Chapman & Hall, London.
- Fan, J., Heckman, N. E. and Wand, M. P. (1995). Local polynomial kernel regression for generalized linear models and quasi-likelihood functions, *Journal of the American Statistical Association*, **90**, 141–150.
- Green, P. J. and Silverman, B. W. (1994). *Nonparametric Regression and Generalized Linear Models*, Chapman & Hall, London.
- Hastie, T. and Tibshirani, R. (1990). *Generalized Additive Models*, Chapman & Hall, London.
- Hjort, N. L. and Jones, M. L. (1996). Locally parametric nonparametric density estimation, *Annals of Statistics*, **24**, 1667–1678.
- Irizarry, R. A. (2001). Information and posterior probability criteria for model selection in local likelihood estimation, *Journal of the American Statistical Association*, **96**, 303–316.
- Konishi, S. (1999). Statistical model evaluation and information criteria, *Multivariate Analysis, Design of Experiments and Survey Sampling* (ed. S. Ghosh), 369–399, Marcel Dekker, New York.
- Konishi, S. and Kitagawa, G. (1996). Generalised information criteria in model selection, *Biometrika*, **83**, 875–890.
- Kullback, S. and Leibler, R. A. (1951). On information and sufficiency, *Annals of Mathematical Statistics*, **22**, 79–86.
- Loader, C. R. (1999). *Local Regression and Likelihood*, Springer, New York.
- McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*, Chapman & Hall, London.
- Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized linear models, *Journal of the Royal Statistical Society, Series, A*, **135**, 370–384.
- Nonaka, Y., Ando, T. and Konishi, S. (2003). Nonlinear regression modeling using regularized local likelihood (in Japanese), *Journal of the Japanese Society of Computational Statistics*, **16**, 43–57.
- Ruppert, D. and Wand, M. P. (1994). Multivariate locally weighted least squares regression, *Annals of Statistics*, **22**, 1346–1370.
- Seifert, B. and Gasser, T. (1996). Variance properties of local polynomials and ensuing modifications, *Statistical Theory and Computational Aspects of Smoothing* (eds. W. Härdle and M. G. Schimek), 50–127, Physica, Heidelberg.
- Seifert, B. and Gasser, T. (2000). Data adaptive ridging in local polynomial regression, *Journal of the Computational and Graphical Statistics*, **9**, 338–360.
- Simonoff, J. S. (1996). *Smoothing Methods in Statistics*, Springer, New York.
- Stone, C. J. (1977). Consistent nonparametric regression (with discussion), *Annals of Statistics*, **5**, 595–645.
- Tibshirani, R. J. and Hastie, T. J. (1987). Local likelihood estimation, *Journal of the American Statistical Association*, **82**, 559–567.
- Wand, M. P. and Jones, M. C. (1995). *Kernel Smoothing*, Chapman & Hall, London.