

On regression model selection for the data with correlated errors

Wen Hsiang Wei

Received: 18 November 2005 / Revised: 1 May 2007 / Published online: 11 August 2007
© The Institute of Statistical Mathematics, Tokyo 2007

Abstract A class of regression model selection criteria for the data with correlated errors is proposed. The proposed class of selection criteria is an estimator of weighted prediction risk. In addition, the proposed selection criteria are the generalizations of several commonly used criteria in statistical analysis. The theoretical and asymptotic properties for the class of criteria are established. Further, in the medium-sample case, the results based on a simulation study are quite consistent with the theoretical ones. The proposed criteria perform well in the simulations. Several applications are also given for a variety of statistical models.

Keywords Generalized cross-validation · Model selection · Nonparametric regression · Penalized likelihood · Smoothing splines

1 Introduction

Consider first the model

$$y_i = f(\mathbf{t}_i) + \epsilon_i, \quad i = 1, \dots, n,$$

where y_i are observations at design points $\mathbf{t}_i = (t_{i1}, t_{i2}, \dots, t_{id})$, $f(t)$ is a function and ϵ_i are zero mean, uncorrelated random errors with common variance σ^2 . Let the fitted values $\hat{\mathbf{f}}(\boldsymbol{\lambda}) = \mathbf{H}(\boldsymbol{\lambda})\mathbf{y}$, where $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_k)$ is a set of parameters associated with the selection of the model, $\mathbf{H}(\boldsymbol{\lambda})$ is an $n \times n$ matrix and $\mathbf{y} = (y_1, y_2, \dots, y_n)^t$. The parameter λ_j could be the subset of the discrete index set $\{1, 2, \dots, p_j\}$ (see Li 1987) or the selection parameter in multivariate nonparametric regression, for examples, the

W. H. Wei (✉)
Department of Statistics, Tung Hai University, Taichung, Taiwan, ROC
e-mail: wenwei@thu.edu.tw

bandwidth in kernel-based method or the smoothing parameter in smoothing splines. Therefore, λ_j is assumed to be non-negative.

Obtaining a good selection parameter estimate is very crucial in the fitting process. One approach is to use the prediction risk

$$P(\lambda) = E \left\{ \frac{\| \mathbf{y}^* - \hat{\mathbf{f}}(\lambda) \|^2}{n} \right\} = \frac{E \left\{ [\mathbf{y}^* - \hat{\mathbf{f}}(\lambda)]' [\mathbf{y}^* - \hat{\mathbf{f}}(\lambda)] \right\}}{n}, \tag{1}$$

as an object function (see Eubank 1988, p. 17), where $\mathbf{y}^* = (y_1^*, y_2^*, \dots, y_n^*)^t$ is a vector of n new observations, $y_i^* = f(t_i) + \epsilon_i^*$, and where ϵ_i^* , uncorrelated with $\epsilon_1, \dots, \epsilon_n$, are zero mean random errors. The estimator of the prediction risk can be used as a criterion for selecting a sensible value of λ . For instance, the commonly used GCV (generalized cross validation) criterion (Craven and Wahba 1979),

$$\text{GCV}(\lambda) = \frac{\hat{\sigma}^2(\lambda)}{[1 - \mu_1(\lambda)]^2},$$

is nearly an unbiased estimator of the prediction risk in some cases (see Eubank 1988, Theorem 2.1), where $\hat{\sigma}^2(\lambda) = \mathbf{y}'[\mathbf{I} - \mathbf{H}(\lambda)]'[\mathbf{I} - \mathbf{H}(\lambda)]\mathbf{y}/n$ is the variance estimate, $[1 - \mu_1(\lambda)]^2$ is a penalty function for the smoothness of the fit, $\mu_1(\lambda) = \text{Tr}[\mathbf{H}(\lambda)]/n$, and $\text{Tr}(\mathbf{A})$ is the trace of the matrix \mathbf{A} . The other predictive mean square error (see Wahba 1990, p. 55), closely related to the previous one, is

$$T(\lambda) = \frac{\| \hat{\mathbf{f}}(\lambda) - \mathbf{f} \|^2}{n} = \frac{[\hat{\mathbf{f}}(\lambda) - \mathbf{f}]' [\hat{\mathbf{f}}(\lambda) - \mathbf{f}]}{n}, \tag{2}$$

where $\mathbf{f} = [f(t_1), f(t_2), \dots, f(t_n)]^t$. As indicated by the weak GCV theorem in the paper of Craven and Wahba, the minimizers of expected values of the predictive mean square error $T(\lambda)$ and GCV criterion are asymptotically equal in term of the prediction risk $E[T(\lambda)]$. Since $E[T(\lambda)] = P(\lambda) - \sigma^2$, the minimizers of the two prediction risks, $E[T(\lambda)]$ and $P(\lambda)$, are the same provided that σ^2 is known. In addition to GCV, other related criteria have the form, $\hat{\sigma}^2(\lambda)/\phi[\mu_1(\lambda)]$, where $\phi(\cdot)$ is a penalty function for the smoothness of the fit. Commonly used selection criteria can be obtained by employing different choices of ϕ , including

1. (Craven and Wahba 1979) GCV: $\phi(\mu_1) = (1 - \mu_1)^2$, (3)

2. (Akaike 1974) AIC: $\phi(\mu_1) = \exp(-2\mu_1)$, (4)

3. (Rice 1984) T : $\phi(\mu_1) = 1 - 2\mu_1$, (5)

4. (Akaike 1970) FPE: $\phi(\mu_1) = \frac{1 - \mu_1}{1 + \mu_1}$, (6)

$$5. \text{ (Shibata 1981) } nS(\lambda): \phi(\mu_1) = \frac{1}{1 + 2\mu_1}, \tag{7}$$

$$6. \text{ (Hocking 1976) } U(\lambda): \phi(\mu_1) = \frac{(1 - \mu_1)(n - 1 - n\mu_1)}{n - 1}, \tag{8}$$

(also see Eubank 1988, pp. 38–40). The selection parameter might be sensitive to the presence of correlation in the errors. The breakdown of several popular data-driven smoothing parameter selection methods was indicated by Opsomer et al. (2001) in non-parametric regression. Further, the prediction risk given in the expression (1) might not be sensible for random errors with unequal variances. For example, suppose the variance of ϵ_i^* is much larger than the one of ϵ_j^* . This implies the error in predicting y_j^* by $f(t_j)$ might be more “predictable” than the one in predicting y_i^* by $f(t_i)$, owing to the smaller variation. Thus, it seems to be sensible to assign different weights, which are related to the variances of the random errors, to the predictive errors. In next section, a weighted prediction risk and its estimator are proposed for cases where the errors of assumed model are correlated. The theoretical properties of the proposed estimator are also provided. Several examples, including a simulation study to justify the theoretical results and the applications of the proposed estimator to a variety of statistical models, are presented in Sects. 3 and 4, respectively. Finally, a concluding discussion is given in Sect. 5. The computational details for the proposed criteria and some supplementary material can be found at <http://web.thu.edu.tw/wenwei/www/papers/aismSupplement.pdf/>.

2 Weighted selection criteria

2.1 Weighted predictive mean square error and weighted selection criteria

Assume the $n \times n$ variance–covariance matrix of the correlated errors $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^t$ is $\text{Var}(\epsilon) = \sigma^2 V^{-1}(\alpha)$, where $\alpha = (\alpha_1, \dots, \alpha_m)$ is a set of correlation parameters. Let the fitted values $\hat{f}_v = H_v(\mathbf{h})y$, where $\mathbf{h} = (\alpha, \lambda) = (h_1, \dots, h_{m+k})$. Define the weighted prediction risk,

$$\begin{aligned} \text{WP}(\mathbf{h}) &= E \left\{ \frac{\|y^* - \hat{f}_v\|_{V(\alpha)}^2}{\psi[V(\alpha)]} \right\} \\ &= \frac{E \left\{ [y^* - H_v(\lambda)y]^t V(\alpha) [y^* - H_v(\lambda)y] \right\}}{\psi[V(\alpha)]}, \end{aligned} \tag{9}$$

where $\psi[V(\alpha)]$ is a positive function. The choice of $\psi(\cdot)$ reflects the effect of the matrix $V(\alpha)$ on the weighted prediction risk. Except the sample size, another sensible choice is $\psi[V(\alpha)] = \text{Tr}[V(\alpha)]$. Thus, when $V(\alpha)$ is an identity matrix, the weighted prediction risk is also the prediction risk given in the expression (1). When $V(\alpha)$ is a diagonal matrix with all diagonal elements equal to 1 except that the i th element is

equal to 0, i.e., no contribution from observation i to the weighted prediction risk, $\psi[V(\alpha)] = n - 1$ in this situation might be a better choice than $\psi[V(\alpha)] = n$. The other possible choice of $\psi(\cdot)$ is the retained number of principal components of $V(\alpha)$. The weighted predictive mean square error, which generalizes the one given in the expression (2), is

$$WT(\mathbf{h}) = \frac{\| \mathbf{f} - \hat{\mathbf{f}}_v \|^2_{V(\alpha)}}{\psi[V(\alpha)]} = \frac{(\mathbf{f} - \hat{\mathbf{f}}_v)^t V(\alpha) (\mathbf{f} - \hat{\mathbf{f}}_v)}{\psi[V(\alpha)]}. \tag{10}$$

When $V(\alpha)$ is a diagonal matrix and $\psi[V(\alpha)] = n$, $WT(\mathbf{h})$ is the weighted mean square error discussed in O’Sullivan et al. (1986). If only $\psi[V(\alpha)] = n$, $WT(\mathbf{h})$ is the weighted mean square error in Wang (1998) with the order of the matrix $V(\alpha)$ equal to one.

The proposed estimator of $WP(\mathbf{h})$ is

$$W(\mathbf{h}) = \frac{\hat{\sigma}_v^2(\mathbf{h})}{\phi_W[\mu_{1v}(\mathbf{h})]}, \tag{11}$$

where

$$\hat{\sigma}_v^2(\mathbf{h}) = \frac{\mathbf{y}^t [\mathbf{I} - \mathbf{H}_v(\mathbf{h})]^t V(\alpha) [\mathbf{I} - \mathbf{H}_v(\mathbf{h})] \mathbf{y}}{\psi[V(\alpha)]},$$

is the variance estimate,

$$\mu_{1v}(\mathbf{h}) = \frac{\text{Tr}[\mathbf{H}_v(\mathbf{h})]}{\psi[V(\alpha)]},$$

is a bounded function, and $\phi_W(\cdot)$ is a penalty function for the smoothness of the fit satisfying

$$\phi_W(x) = 1 - 2x + p(x), \quad \lim_{x \rightarrow 0} \frac{p(x)}{x^2} = c,$$

and where c is a finite constant. When $V(\alpha)$ is an identity matrix and $\psi[V(\alpha)] = n$, the selection criteria given in expressions from (3) to (8) take the form of the proposed class of weighted selection criteria $W(\mathbf{h})$. In addition, when $\psi[V(\alpha)] = n$ and $\phi_W(x) = (1 - x)^2$, $W(\mathbf{h})$ is the GCV function considered in Wang’s paper, which corresponds to the one used in Altman (1990).

2.2 Properties of weighted selection criteria

The following theorem and corollary provide theoretical supports for the proposed selection criteria. The proofs largely follow the GCV theorem (Eubank 1988, Theorem 2.1) and the weak GCV theorem given by Craven and Wahba.

Theorem 1 Let $\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon}$, $\mathbf{f} = [f(t_1), \dots, f(t_n)]^t$, and $\text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{V}^{-1}(\boldsymbol{\alpha})$ is an $n \times n$ variance–covariance matrix. Denote a positive function

$$\mu_{2v}(\mathbf{h}) = \frac{\text{Tr} [\mathbf{H}_v^t(\mathbf{h}) \mathbf{V}(\boldsymbol{\alpha}) \mathbf{H}_v(\mathbf{h}) \mathbf{V}^{-1}(\boldsymbol{\alpha})]}{\psi[\mathbf{V}(\boldsymbol{\alpha})]}.$$

Then,

$$\frac{|WP(\mathbf{h}) - \{E[W(\mathbf{h})] - r(\mathbf{h})\}|}{E[WT(\mathbf{h})]} \leq e(\mathbf{h}),$$

where

$$e(\mathbf{h}) = \frac{1}{\phi_W[\mu_{1v}(\mathbf{h})]} \left\{ \left(\left| \frac{p[\mu_{1v}(\mathbf{h})]}{\mu_{1v}(\mathbf{h})} \right| + 2 \right) |\mu_{1v}(\mathbf{h})| + \frac{1_v p[\mu_{1v}(\mathbf{h})]}{\mu_{2v}(\mathbf{h})} \right\},$$

$r(\mathbf{h}) = 2(1_v - 1)\mu_{1v}(\mathbf{h})\sigma^2/\phi_W[\mu_{1v}(\mathbf{h})]$, and where $1_v = n/\psi[\mathbf{V}(\boldsymbol{\alpha})]$.

If the values of $e(\mathbf{h})$ and $r(\mathbf{h})$ are small, the theorem implies that the difference between $E[W(\mathbf{h})]$ and $WP(\mathbf{h})$ is small relative to the other weighted prediction risk $E[WT(\mathbf{h})]$. Further, if $E[WT(\mathbf{h})]$ is bounded, $W(\mathbf{h})$ is nearly an unbiased estimator of $WP(\mathbf{h})$ in this situation. Also, the minimum values of $E[WT(\mathbf{h})]$ and $E[W(\mathbf{h})]$ have similar properties in term of minimizing the weighted prediction risk $E[WT(\mathbf{h})]$, as indicated by the following corollary.

Corollary 1 If $E[WT(\mathbf{h})]$ has (at least) one minimizer \mathbf{h}_n^* and $\psi[\mathbf{V}(\boldsymbol{\alpha})] = O(n)$, there exists a sequence of minimizers $\hat{\mathbf{h}}_n$ of $E[W(\mathbf{h})]$ such that

$$\lim_{n \rightarrow \infty} \frac{E[WT(\hat{\mathbf{h}}_n)]}{E[WT(\mathbf{h}_n^*)]} = 1,$$

under the assumptions that

$$\lim_{n \rightarrow \infty} \mu_{1v}(\mathbf{h}_n) = 0, \quad \lim_{n \rightarrow \infty} \frac{\mu_{1v}^2(\mathbf{h}_n)}{\mu_{2v}(\mathbf{h}_n)} = 0, \quad \lim_{n \rightarrow \infty} \frac{(1_v - 1)\mu_{1v}(\mathbf{h}_n)}{\mu_{2v}(\mathbf{h}_n^*)} = 0,$$

where $\mathbf{h}_n = \hat{\mathbf{h}}_n$ or $\mathbf{h}_n = \mathbf{h}_n^*$.

Both the proofs of Theorem 1 and Corollary 1 are given in the Appendices A and B, respectively. The condition

$$\lim_{n \rightarrow \infty} \frac{(1_v - 1)\mu_{1v}(\mathbf{h}_n)}{\mu_{2v}(\mathbf{h}_n^*)} = 0,$$

is not required if either $\psi[\mathbf{V}(\boldsymbol{\alpha})] = n$ or $E[WT(\mathbf{h}_n^*)]$ does not tend to zero. The corollary indicates the selection criteria $W(\mathbf{h})$ given in the expression (11) are weighted

predictive mean-square error (the one given in the expression (10)) criteria. When $V(\alpha)$ is an identity matrix and $\phi_W(x) = (1 - x)^2$, the corollary is the weak GCV theorem. Since $E[\text{WT}(\mathbf{h})] = \text{WP}(\mathbf{h}) - 1_v\sigma^2$, the minimizers of the two prediction risks, $E[\text{WT}(\mathbf{h})]$ and $\text{WP}(\mathbf{h})$, are the same given known values of σ^2 and correlation parameters.

To have the optimal properties, it is not necessary to define $\mu_{1v}(\mathbf{h})$ as a function of $\mathbf{H}_v(\mathbf{h})$ through its trace. For instance, if the function $\mu_{1v}(\mathbf{h})$ satisfies $\mu_{1v}(\mathbf{h}) - \text{Tr}[\mathbf{H}_v(\mathbf{h})]/\psi[V(\alpha)]$ converging uniformly to 0 or more relaxed condition $\lim_{n \rightarrow \infty} \{\mu_{1v}(\mathbf{h}_n) - \text{Tr}[\mathbf{H}_v(\mathbf{h}_n)]/\psi[V(\alpha_n)]\} = 0$, the theorem and corollary still hold, where $\mathbf{h}_n = (\alpha_n, \lambda_n)$. Thus, more flexible choices of $\mu_{1v}(\mathbf{h})$ can be made than only through the trace of $\mathbf{H}_v(\mathbf{h})$.

3 Simulations

The purpose of the following simulations is to illustrate that all the selection criteria given in the expressions from (3) to (8) perform well. In addition, the numerical results are consistent with the theoretical ones, even in medium-sample case. A range of scenarios, including different choices of $V(\alpha)$ and noise levels, have been set up for the simulation study.

3.1 Weighted linear regression

In the simulation, the values of four input variables, X_1, X_2, X_3, X_4 , were in $[0, 1]$ and 100 observations were generated from the model

$$y_i = 1 + 2x_{i1} + 4x_{i3} + \epsilon_i, \quad i = 1, \dots, 100,$$

where ϵ_i are zero mean random errors. The errors were generated from both Gaussian AR(1) and MA(1) processes with the standard deviations of uncorrelated Gaussian errors, σ_g , equal to 0.2, 1, and 2. The autocorrelation values at lag 1, $\rho(1)$, for the Gaussian AR(1) process were $-0.8, -0.2, 0.2$, and 0.8 , while $-0.4, -0.2, 0.2$, and 0.4 for the Gaussian MA(1) process. For simplicity, assume $V(\alpha) = V$ is known. Two choices for $\psi(V)$ were the sample size and retained number of principal components of the matrix V by including just enough components to explain 90% amount of the variance. 500 replicates of random errors were generated. For each sample, the proposed selection criteria and associated weighted predictive mean-square error given in the expression (9) were computed for all possible models with at least one input variable, i.e., total 15 possible models. The averages of these quantities can be used to estimate the expected values of the proposed selection criteria and weighted predictive mean-square error. The results for the averages corresponding to the true model, i.e., the estimates of $\text{WP}(\hat{\mathbf{h}}_n)$ and $E[W(\hat{\mathbf{h}}_n)]$, can be obtained. The selection criteria and weighted predictive mean-square error produce very similar results, which provide numerical support for Theorem 1. Since the results under different settings are quite consistent, Table 1 summarizes parts of these results for Gaussian AR(1) process. The first number in the parenthesis is the average with $\psi(V)$ equal to the sample size and

Table 1 Weighted linear regression and smoothing with Gaussian AR(1) and MA(1) random errors, respectively

	AR(1)		MA(1)	
	$\rho(1) = -0.8$ $\sigma_g = 1$	$\rho(1) = 0.8$ $\sigma_g = 1$	$\rho(1) = -0.4$ $\sigma_g = 1$	$\rho(1) = 0.4$ $\sigma_g = 1$
WP(\mathbf{h})	(1.030,1.688)	(1.030,1.688)		
GCV	(1.036,1.768)	(1.032,1.760)	(1.001,1.057)	(1.002,1.106)
AIC	(1.035,1.764)	(1.031,1.756)	(1.000,1.023)	(1.002,1.106)
T	(1.037,1.773)	(1.033,1.765)	(1.001,1.057)	(1.002,1.106)
FPE	(1.035,1.764)	(1.031,1.756)	(1.000,1.023)	(1.002,1.106)
nS	(1.034,1.756)	(1.029,1.748)	(1.000,1.023)	(1.054,1.106)
U	(1.037,1.769)	(1.032,1.761)	(1.001,1.057)	(1.002,1.106)

the second number otherwise. On the other hand, since the averages of the proposed selection criteria and weighted predictive mean-square error attain their minimums as the postulated model being the true model, the estimate of $E[\text{WT}(\hat{\mathbf{h}}_n)]/E[\text{WT}(\mathbf{h}_n^*)]$ is equal to 1.

3.2 Weighted smoothing

In the simulation, 100 observations were generated from the model

$$y_i = \sin(2\pi i/n), \quad i = 1, \dots, 100,$$

where ϵ_i are zero mean random errors. The function was utilized in the simulation study (Wang 1998, p. 344). The errors were generated from both Gaussian AR(1) and MA(1) processes with the same variances of uncorrelated Gaussian errors and autocorrelation values at lag 1 as the ones in the simulation study given in Sect. 3.1. Also, there were two choices for $\psi(\mathbf{V})$, as given in the previous simulation study. 100 replicates of random errors were generated. For each sample, the smoothing spline fit using the B-spline of degree 3 and a second order penalty (see Eilers and Marx 1996) was computed. The chosen knots divided the domain of t ($0.05\pi - 2\pi$) into 20 intervals of equal width. Then, the proposed selection criteria and weighted predictive mean-square error given in the expression (9) were computed for the smoothing parameter with values 10, 20, ..., 1000. The averages of the proposed selection criteria and weighted predictive mean-square error were computed. Evaluated at their minimums, the results can be obtained. Similar to Sect. 3.1, the numbers in the parenthesis are the averages based on different choices of $\psi(\mathbf{V})$. The values for different selection criteria and the weighted predictive mean-square error are quite close. This provides numerical support for Theorem 1. In addition, the estimated values of $E[\text{WT}(\hat{\mathbf{h}}_n)]/E[\text{WT}(\mathbf{h}_n^*)]$ are quite close to 1. Note that similar results can be obtained when the domain of t was divided into 50 intervals of equal width. Since these results under different settings are quite

consistent, parts of these results, the estimated values of $E[\text{WT}(\hat{\mathbf{h}}_n)]/E[\text{WT}(\mathbf{h}_n^*)]$ for Gaussian MA(1) process, are provided in Table 1.

4 Applications

In the following examples, the proposed weighted selection criteria are applied to a variety of statistical models.

4.1 Weighted linear regression

Consider the linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ is an $n \times p$ design matrix and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^t$ is a $p \times 1$ coefficient vector. For simplicity, assume $\mathbf{V}(\boldsymbol{\alpha}) = \mathbf{V}$ is known. Let the discrete index set $\Lambda = \{1, 2, \dots, p\}$ and $\mathbf{X}_\lambda = (\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_l})$, where $\lambda = \{i_1, \dots, i_l\}$ is a subset of the index set Λ . Thus, $\mathbf{H}_v(\lambda) = \mathbf{X}_\lambda[\mathbf{X}_\lambda^t \mathbf{V} \mathbf{X}_\lambda]^{-1} \mathbf{X}_\lambda^t \mathbf{V}$ and $\mu_{1v}(\lambda) = \mu_{2v}(\lambda) = l/\psi(\mathbf{V})$. Further, both $e(\lambda)$ and $r(\lambda)$ tend to 0 as n tends to infinity. Therefore, $W(\lambda)$ is nearly an unbiased estimator of $\text{WP}(\lambda)$ implied by Theorem 1. As $\lim_{n \rightarrow \infty} (1_v - 1) = 0$, Corollary 1 also holds.

4.2 Weighted smoothing

Suppose the weighted (or generalized) smoothing spline estimate \hat{f}_v is the minimizer of

$$(\mathbf{y} - \mathbf{f})^t \mathbf{V}(\mathbf{y} - \mathbf{f}) + \sum_{j=1}^k \lambda_j \|J_j(\mathbf{f})\|^2, \tag{12}$$

over the class of all twice differentiable functions (also see Opsomer et al. 2001, p. 148), where J_j are some operators, for examples, the orthogonal projectors given in chapter 10 of Wahba (1990) or $\|J_j(\mathbf{f})\|^2 = \int [\frac{\partial^2 f(\mathbf{t})}{\partial t_j^2}]^2 dt_j$ for additive splines (see Hastie and Tibshirani 1990). Let $f(t) = \sum_{j=1}^{p_n} a_j B_j(t)$, where p_n is the number of suitably chosen basis functions, usually at least large enough to ensure the accuracy of the approximation and $B_j(t)$ are basis functions, for example, the commonly used B-splines (see Green and Silverman 1994, pp. 155–159). Thus, when $\|J_j(\mathbf{f})\|^2 = \mathbf{a}^t \mathbf{P}_j \mathbf{a}$, the expression (12) reduces to

$$(\mathbf{y} - \mathbf{B}\mathbf{a})^t \mathbf{V}(\mathbf{y} - \mathbf{B}\mathbf{a}) + \sum_{j=1}^k \lambda_j \mathbf{a}^t \mathbf{P}_j \mathbf{a},$$

where $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n)^t = [B_j(t_i)]_{ij}$ is an $n \times p_n$ matrix, $\mathbf{a} = (a_1, \dots, a_{p_n})^t$ and \mathbf{P}_j are $p_n \times p_n$ penalty matrices. The estimate of the coefficient \mathbf{a} is $\hat{\mathbf{a}}(\boldsymbol{\lambda}) = (\mathbf{B}^t \mathbf{V} \mathbf{B} + \sum_{j=1}^k \lambda_j \mathbf{P}_j)^{-1} \mathbf{B}^t \mathbf{V} \mathbf{y}$. Then, the minimizer is

$$\hat{\mathbf{f}}_v = \left[\hat{f}_v(t_1, \boldsymbol{\lambda}), \dots, \hat{f}_v(t_n, \boldsymbol{\lambda}) \right]^t = \mathbf{B} \left(\mathbf{B}^t \mathbf{V} \mathbf{B} + \sum_{j=1}^k \lambda_j \mathbf{P}_j \right)^{-1} \mathbf{B}^t \mathbf{V} \mathbf{y} = \mathbf{H}_v(\boldsymbol{\lambda}) \mathbf{y},$$

where $\mathbf{H}_v(\boldsymbol{\lambda}) = [h_{is}(\boldsymbol{\lambda})]_{is} = \mathbf{B}(\mathbf{B}^t \mathbf{V} \mathbf{B} + \sum_{j=1}^k \lambda_j \mathbf{P}_j)^{-1} \mathbf{B}^t \mathbf{V}$ is the hat matrix, $i = 1, \dots, n, s = 1, \dots, n$.

Suppose \mathbf{P}_j are positive-definite matrices with eigenvalues $0 < h_{j,p_n}^* \leq \dots \leq h_{j,2}^* \leq h_{j,1}^*$ and $0 \leq h_{\min(n,p_n)} \leq \dots \leq h_2 \leq h_1$ are the eigenvalues of $\mathbf{B}^t \mathbf{V} \mathbf{B}$. If $h_i = O(i^{-q_1})$, $q_1 > 1$ and the fastest decay rate of the eigenvalues $h_{j,\cdot}^*$ is $h_{l,p_n}^* = O(n^{-q_2})$, $q_2 < 1/2$, Theorem 1 and Corollary 1 hold. The justifications are given in Appendix C. When $\mathbf{B}^t \mathbf{V} \mathbf{B}$ is nonsingular and there exist generalized eigenvalues $h_{l,i}$ of $(\mathbf{B}^t \mathbf{V} \mathbf{B})^{-1} \mathbf{x} = \lambda^* \mathbf{P}_l^{-1} \mathbf{x}$ with decay rate $O(i^{-q_3})$, $q_3 > 1$ (see Golub and Van Loan 1993, pp. 466–472), Theorem 1 and Corollary 1 also hold. The detailed justifications were delegated to the supplementary material. When the matrix \mathbf{V} involves a set of correlation parameters, Theorem 1 and Corollary 1 can be justified alone the lines given in Appendix C by specifying the decay rate of the eigenvalues of $\mathbf{B}^t \mathbf{V}(\boldsymbol{\alpha}) \mathbf{B}$.

4.3 Model selection in multivariate linear regression

Consider the linear model

$$\mathbf{y}_j = \mathbf{X}_j \boldsymbol{\beta}_j + \boldsymbol{\epsilon}_j, \quad j = 1, \dots, k,$$

where $\mathbf{y}_j = (y_{j1}, \dots, y_{jn_j})^t$, $\mathbf{X}_j = (\mathbf{x}_{j1}, \dots, \mathbf{x}_{jp_j})$ is an $n_j \times p_j$ design matrix, $\boldsymbol{\beta}_j = (\beta_{j1}, \dots, \beta_{jp_j})^t$ and $\boldsymbol{\epsilon}_j = (\epsilon_{j1}, \epsilon_{j2}, \dots, \epsilon_{jn_j})^t$. Let the discrete index sets $\Lambda_j = \{1, 2, \dots, p_j\}$ and $\mathbf{X}_{\lambda_j} = (\mathbf{x}_{ji_1}, \dots, \mathbf{x}_{ji_{l_j}})$, where $\lambda_j = \{i_1, \dots, i_{l_j}\}$ is a subset of Λ_j . Thus,

$$\mathbf{H}_v(\mathbf{h}) = \mathbf{X}_{\boldsymbol{\lambda}} \left[\mathbf{X}_{\boldsymbol{\lambda}}^t \mathbf{V}(\boldsymbol{\alpha}) \mathbf{X}_{\boldsymbol{\lambda}} \right]^{-1} \mathbf{X}_{\boldsymbol{\lambda}}^t \mathbf{V}(\boldsymbol{\alpha}),$$

and $\mu_1(\mathbf{h}) = \text{Tr}[\mathbf{H}_v(\mathbf{h})] / \psi[V(\boldsymbol{\alpha})]$, where $\mathbf{X}_{\boldsymbol{\lambda}} = \text{Diag}(\mathbf{X}_{\lambda_1}, \dots, \mathbf{X}_{\lambda_k})$ is a $\sum_{j=1}^k n_j \times \sum_{j=1}^k l_j$ matrix and $\sigma^2 \mathbf{V}^{-1}(\boldsymbol{\alpha})$ is the variance–covariance matrix of $(\boldsymbol{\epsilon}_1^t, \dots, \boldsymbol{\epsilon}_k^t)^t$.

4.4 Generalized ridge regression

Let $\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}$, where \mathbf{X} is an $n \times p_n$ design matrix and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_{p_n})^t$, and where $p_n \leq n$. The weighted (or generalised) ridge regression estimate of $\boldsymbol{\beta}$, the

minimizer of

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^t \mathbf{V}(\boldsymbol{\alpha})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \sum_{j=1}^k \lambda_j \boldsymbol{\beta}^t \mathbf{P}_j \boldsymbol{\beta},$$

is $[\mathbf{X}^t \mathbf{V}(\boldsymbol{\alpha})\mathbf{X} + \sum_{j=1}^k \lambda_j \mathbf{P}_j]^{-1} \mathbf{X}^t \mathbf{V}(\boldsymbol{\alpha})\mathbf{y}$, where \mathbf{P}_j are nonsingular matrices. Thus,

$$\mathbf{H}_v(\mathbf{h}) = \mathbf{X} \left[\mathbf{X}^t \mathbf{V}(\boldsymbol{\alpha})\mathbf{X} + \sum_{j=1}^k \lambda_j \mathbf{P}_j \right]^{-1} \mathbf{X}^t \mathbf{V}(\boldsymbol{\alpha}).$$

4.5 Partial splines

Let

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + f(t_i) + \epsilon_i, \quad i = 1, \dots, n,$$

(see Eubank 1988, pp. 292–293; Wahba 1990, Chapt. 6) and $f(t) = \sum_{j=1}^{p_n} a_j B_j(t)$.

Then, the weighted (or generalized) partial spline estimate \hat{f}_v is the minimizer of

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{f})^t \mathbf{V}(\boldsymbol{\alpha})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{f}) + \sum_{j=1}^k \lambda_j \|J_j(f)\|^2, \tag{13}$$

where J_j are some operators. When $\|J_j(f)\|^2 = \mathbf{a}^t \mathbf{P}_j \mathbf{a}$, the expression (13) reduces to

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{B}\mathbf{a})^t \mathbf{V}(\boldsymbol{\alpha})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{B}\mathbf{a}) + \sum_{j=1}^k \lambda_j \mathbf{a}^t \mathbf{P}_j \mathbf{a},$$

where $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n)^t = [B_j(t_i)]_{ij}$ is an $n \times p_n$ matrix, $\mathbf{a} = (a_1, \dots, a_{p_n})^t$, and \mathbf{P}_j are $p_n \times p_n$ penalty matrices. The estimates of the coefficients \mathbf{a} and $\boldsymbol{\beta}$ are

$$\hat{\mathbf{a}}(\mathbf{h}) = \left(\mathbf{B}^t \mathbf{V}(\boldsymbol{\alpha})\mathbf{B} + \sum_{j=1}^k \lambda_j \mathbf{P}_j \right)^{-1} \mathbf{B}^t \mathbf{V}(\boldsymbol{\alpha}) [\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}(\mathbf{h})],$$

and

$$\hat{\boldsymbol{\beta}}(\mathbf{h}) = \{ \mathbf{X}^t \mathbf{V}(\boldsymbol{\alpha}) [\mathbf{I} - \mathbf{H}_2(\mathbf{h})] \mathbf{X} \}^{-1} \mathbf{X}^t \mathbf{V}(\boldsymbol{\alpha}) [\mathbf{I} - \mathbf{H}_2(\mathbf{h})] \mathbf{y},$$

respectively, where

$$H_2(\mathbf{h}) = \mathbf{B} \left(\mathbf{B}^t \mathbf{V}(\alpha) \mathbf{B} + \sum_{j=1}^k \lambda_j \mathbf{P}_j \right)^{-1} \mathbf{B}^t \mathbf{V}(\alpha).$$

Then, the minimizer is

$$\begin{aligned} \hat{\mathbf{f}}_v &= [\hat{f}_v(t_1, \mathbf{h}), \dots, \hat{f}_v(t_n, \mathbf{h})]^t = \mathbf{H}_v(\mathbf{h}) \mathbf{y} \\ &= [\mathbf{H}_1(\mathbf{h}) + \mathbf{H}_2(\mathbf{h}) - \mathbf{H}_1(\mathbf{h})\mathbf{H}_2(\mathbf{h})] \mathbf{y}, \end{aligned}$$

where

$$\mathbf{H}_1(\mathbf{h}) = \mathbf{X} \{ \mathbf{X}^t \mathbf{V}(\alpha) [\mathbf{I} - \mathbf{H}_2(\mathbf{h})] \mathbf{X} \}^{-1} \mathbf{X}^t \mathbf{V}(\alpha) [\mathbf{I} - \mathbf{H}_2(\mathbf{h})].$$

4.6 Spline smoothing in generalized linear models

Consider the standard generalized linear model in which each component of the response vector has a distribution taking the form

$$f(y_i; \theta_i, \phi) = \exp \left[\frac{y_i \theta_i - m(\theta_i)}{u(\phi)} + c(y_i, \phi) \right],$$

where θ_i and ϕ are scalar parameters, and $m(\cdot)$, $u(\cdot)$ and $c(\cdot)$ are specific functions. The dependence of the response y_i on the associated explanatory variable t_i can be modeled through the link function $d(\cdot)$, where $\theta_i = d(\alpha + \beta^t t_i)$, and where α and β are some parameters. The natural link and $u(\phi) = 1$ are assumed hereafter. In addition, let $\theta_i = f(t_i)$ and $\mathbf{f} = \mathbf{B}\mathbf{a}$. The estimate of f is the minimizer of the penalized negative logarithm of the likelihood,

$$\begin{aligned} &\sum_{i=1}^n \{m[f(t_i)] - y_i f(t_i)\} + \frac{1}{2} \sum_{j=1}^k \lambda_j \|J_j(f)\|^2 \\ &= \sum_{i=1}^n m[(\mathbf{B}\mathbf{a})_i] - \mathbf{y}^t \mathbf{B}\mathbf{a} + \frac{1}{2} \sum_{j=1}^k \lambda_j \mathbf{a}^t \mathbf{P}_j \mathbf{a}, \end{aligned}$$

where $(\mathbf{B}\mathbf{a})_i$ is the i th element of the vector $\mathbf{B}\mathbf{a}$. The estimate of the coefficient vector \mathbf{a} can be written as a weighted penalized least squares estimate,

$$\hat{\mathbf{a}}(\lambda) = \left(\mathbf{B}^t \ddot{\mathbf{M}} \mathbf{B} + \sum_{j=1}^k \lambda_j \mathbf{P}_j \right)^{-1} \mathbf{B}^t \ddot{\mathbf{M}} \mathbf{z},$$

where $\ddot{\mathbf{M}} = \text{diag}\{m''[\hat{\mathbf{a}}^t(\lambda)\mathbf{b}_1], \dots, m''[\hat{\mathbf{a}}^t(\lambda)\mathbf{b}_n]\}$ and $\mathbf{z} = \mathbf{B}\hat{\mathbf{a}}(\lambda) + \ddot{\mathbf{M}}^{-1}(\mathbf{y} - \hat{\mathbf{m}})$, and where $\hat{\mathbf{m}} = \{m'[\hat{\mathbf{a}}^t(\lambda)\mathbf{b}_1], \dots, m'[\hat{\mathbf{a}}^t(\lambda)\mathbf{b}_n]\}^t$. Let

$$\mathbf{H}_v(\lambda) = \ddot{\mathbf{M}}^{1/2} \mathbf{B} \left(\mathbf{B}^t \ddot{\mathbf{M}} \mathbf{B} + \sum_{j=1}^k \lambda_j \mathbf{P}_j \right)^{-1} \mathbf{B}^t \ddot{\mathbf{M}}^{1/2},$$

and $\mu_{1v}(\lambda) = \text{Tr}[\mathbf{H}_v(\lambda)]/\psi(\ddot{\mathbf{M}})$. The smoothing parameter estimate is the minimizer of the following function, $\|\ddot{\mathbf{M}}^{-1/2}(\mathbf{y} - \hat{\mathbf{m}})\|^2/\phi_W[\mu_{1v}(\lambda)]$. When $\phi_W(x) = (1 - x)^2$ and $\psi(\ddot{\mathbf{M}}) = n$, the above function is the GCV function (see Wahba 1990, p. 113; O'Sullivan et al. 1986).

4.7 Nonparametric regression incorporating the information provided by derivatives

The derivatives of a function can provide useful information for data analysis. To incorporate the information provided by the derivatives, consider first the model

$$y_i^{(d)} = f^{(d)}(t_i) + \epsilon_i^{(d)}, \quad i = 1, \dots, n; \quad d = 0, \dots, k,$$

where the data, $y_i^{(d)}$, associated with the first k 'th derivatives of the function $f(\cdot)$ are available, $f^{(d)}(\cdot)$ is the d th order of derivative of the function with $f(\cdot) = f^{(0)}(\cdot)$, and where $\epsilon_i^{(d)}$ are zero mean random errors. Let $f(t) = \sum_{j=1}^{p_n} a_j B_j^{(0)}(t)$. The weighted sum of squares are

$$\sum_{d=0}^k \lambda_d \left(\mathbf{y}^{(d)} - \mathbf{B}^{(d)} \mathbf{a} \right)^t \mathbf{V}^{(d)}(\boldsymbol{\alpha}) \left(\mathbf{y}^{(d)} - \mathbf{B}^{(d)} \mathbf{a} \right),$$

where $\lambda_0 = 1$, $\mathbf{y}^{(d)} = (y_1^{(d)}, \dots, y_n^{(d)})^t$, $\mathbf{B}^{(d)} = \{B_j^{(d)}(t_i)\}_{ij}$, $\sigma^2 \mathbf{V}^{(d)}(\boldsymbol{\alpha})$ is the variance–covariance matrix of the data $\mathbf{y}^{(d)}$, and where $B_j^{(d)}$ is the d th order of derivative of the basis function $B_j^{(0)}$ with respect to t . The parameters λ_d play a key role in controlling the trade-off between zero order derivative information represented by the weighted residual sum of squares and the information provided by other order of derivatives. Then,

$$\begin{aligned} & \mathbf{H}_v(\mathbf{h}) \mathbf{y} \\ &= \left(\oplus_{d=0}^k \mathbf{B}^{(d)} \right) \left\{ \oplus_{i=0}^k \left\{ \mathbf{B}^t \mathbf{V}^{(0)}(\boldsymbol{\alpha}) \mathbf{B} + \left[\sum_{d=1}^k \lambda_d (\mathbf{B}^{(d)})^t \mathbf{V}^{(d)}(\boldsymbol{\alpha}) \mathbf{B}^{(d)} \right] \right\}^{-1} \right\} \\ & \quad \times \left[\oplus_{d=0}^k \mathbf{B}^{(d)} \mathbf{V}^{(d)}(\boldsymbol{\alpha}) \right] \mathbf{y}, \end{aligned}$$

where $\mathbf{y} = [(y^{(0)})^t, \dots, (y^{(k)})^t]^t$ and $\oplus_{i=1}^k \mathbf{A}_i$ is the Kronecker sum of the matrices $\mathbf{A}_1, \dots, \mathbf{A}_k$.

5 Concluding discussion

As indicated by Theorem 1, all the selection criteria given in the expressions from (3) to (8) have similar optimal properties in term of the difference between their expected values and the prediction risk $WP(\mathbf{h})$. Further, Corollary 1 indicates that the minimizers of $E[W(\mathbf{h})]$ and $E[WT(\mathbf{h})]$ also have similar properties in term of minimizing the weighted prediction risk. In Sect. 3.1, the correct model can be selected by both the averages of these selection criteria and the prediction risk $WP(\mathbf{h})$. As illustrated in Table 1, the numerical results based on the simulated data are quite consistent with the theoretical results, even in the medium-sample case. In addition, the proposed criteria have wide applications, as presented in Sect. 4.

When the correlation is parametrically specified, commonly used structures, such as the ARIMA model, might be employed to fit the data. Thus, the proposed criteria can be used to estimate both the selection parameters and correlation parameters. For the correlation not parametrically specified, the analogue criteria can be obtained by replacing $V(\alpha)$ in the proposed criteria given in the expression (11) with a sensible non-parametric estimate \hat{V} . The minimum values of $E[WT(\lambda)]$ and the expected value of the analogue criteria still have similar properties in term of minimizing the weighted prediction risk $E[WT(\lambda)]$, i.e., a result analogous to Corollary 1. The arguments in proving the result are similar to the ones given in Corollary 1 of Wei (2005).

There is still room for future research in regression model selection with correlated errors. Nonparametric modeling and semi-parametric modelling have been widely used techniques in recent years (see Härdle et al. 2004). The proposed criteria could be applied to some of these models. The proposed criteria have different sensitivities to the changes of the selection parameters or correlation parameters. A thorough robust and sensitivity analysis for different selection criteria could be helpful. The other issue is about the comparison of different criteria. The characterization of efficiency associated with a selection criterion is still unclear.

6 Appendix: Proofs

6.1 Appendix A: Proofs of Theorem 1

Based on the formula for the mean of a quadratic form, the expected value of the weighted predictive mean-square error is

$$E [WT(\mathbf{h})] = b^2(\mathbf{h}) + \sigma^2 \mu_{2v}(\mathbf{h}),$$

where $b^2(\mathbf{h}) = \{1/\psi[V(\alpha)]\} \mathbf{f}^t [\mathbf{I} - \mathbf{H}_v(\mathbf{h})]^t V(\alpha) [\mathbf{I} - \mathbf{H}_v(\mathbf{h})] \mathbf{f}$ is the bias term and $\sigma^2 \mu_{2v}(\mathbf{h})$ is the variance term. Similarly,

$$E [W(\mathbf{h})] = \frac{b^2(\mathbf{h}) + \sigma^2 [1_v - 2\mu_{1v}(\mathbf{h}) + \mu_{2v}(\mathbf{h})]}{\phi_W [\mu_{1v}(\mathbf{h})]}.$$

The following shows that the difference between the expected values of the selection criteria and $WP(\lambda)$ is close to 0. Denote

$$r(\mathbf{h}) = 2(1_v - 1)\mu_{1v}(\mathbf{h})\sigma^2 / \phi_W[\mu_{1v}(\mathbf{h})].$$

Then,

$$\begin{aligned} E[W(\mathbf{h})] - 1_v\sigma^2 - r(\mathbf{h}) &= \frac{b^2(\mathbf{h}) + \sigma^2 [1_v - 2\mu_{1v}(\mathbf{h}) + \mu_{2v}(\mathbf{h})] - 1_v\sigma^2\phi_W[\mu_{1v}(\mathbf{h})]}{\phi_W[\mu_{1v}(\mathbf{h})]} - r(\mathbf{h}) \\ &= \frac{b^2(\mathbf{h}) + \sigma^2\mu_{2v}(\mathbf{h}) - 1_v\sigma^2 p[\mu_{1v}(\mathbf{h})]}{\phi_W[\mu_{1v}(\mathbf{h})]}. \end{aligned}$$

Further, since $WP(\mathbf{h}) = E[WT(\mathbf{h})] + 1_v\sigma^2$,

$$\begin{aligned} \frac{WP(\mathbf{h}) - \{E[W(\mathbf{h})] - r(\mathbf{h})\}}{E[WT(\mathbf{h})]} &= \frac{b^2(\mathbf{h}) + \sigma^2\mu_{2v}(\mathbf{h}) - \{b^2(\mathbf{h}) + \sigma^2\mu_{2v}(\mathbf{h}) - 1_v\sigma^2 p[\mu_{1v}(\mathbf{h})]\} / \phi_W[\mu_{1v}(\mathbf{h})]}{b^2(\mathbf{h}) + \sigma^2\mu_{2v}(\mathbf{h})} \\ &= \frac{-2\mu_{1v}(\mathbf{h}) + p[\mu_{1v}(\mathbf{h})]}{\phi_W[\mu_{1v}(\mathbf{h})]} + \frac{1_v\sigma^2 p[\mu_{1v}(\mathbf{h})]}{\phi_W[\mu_{1v}(\mathbf{h})][b^2(\mathbf{h}) + \sigma^2\mu_{2v}(\mathbf{h})]}. \end{aligned}$$

Finally, since $|\mu_{1v}(\mathbf{h})| \leq M$ and $b^2(\mathbf{h}) \geq 0$,

$$\begin{aligned} &\frac{|WP(\mathbf{h}) - \{E[W(\mathbf{h})] - r(\mathbf{h})\}|}{E[WT(\mathbf{h})]} \\ &\leq \frac{1}{\phi_W[\mu_{1v}(\mathbf{h})]} \left\{ \left(\left| \frac{p[\mu_{1v}(\mathbf{h})]}{\mu_{1v}(\mathbf{h})} \right| + 2 \right) |\mu_{1v}(\mathbf{h})| + \frac{1_v p[\mu_{1v}(\mathbf{h})]}{\mu_{2v}(\mathbf{h})} \right\} \\ &= e(\mathbf{h}), \end{aligned}$$

where M is some constant.

6.2 Appendix B: Proofs of Corollary 1

By Theorem 1, the following inequality can be obtained,

$$E[WT(\mathbf{h})][1 - e(\mathbf{h})] \leq E[W(\mathbf{h})] - 1_v\sigma^2 - r(\mathbf{h}) \leq E[WT(\mathbf{h})][1 + e(\mathbf{h})],$$

for all \mathbf{h} . Note that $r(\mathbf{h})$ tends to 0 as $\mu_{1v}(\mathbf{h})$ tends to 0 and $\psi[V(\alpha)] = O(n)$. Thus, the difference between the expected values of $W(\mathbf{h})$ and the predictive mean-square error $WT(\mathbf{h})$ is approximately $1_v\sigma^2$. Intuitively, this implies the minimizers of $E[WT(\mathbf{h})]$ and $E[W(\mathbf{h})]$ should be very close and further the corollary holds. The rigorous justifications are as follows. Since $E[W(\mathbf{h}_n^*)] - 1_v\sigma^2 - r(\mathbf{h}_n^*) \leq E[WT(\mathbf{h}_n^*)][1 + e(\mathbf{h}_n^*)]$

and $r(\mathbf{h}_n^*)$ tends to 0, there exists N such that at least one minimizer $\hat{\mathbf{h}}_n$ of $E[W(\mathbf{h})]$ is in the nonempty set

$$\left\{ \mathbf{h} : E[W(\mathbf{h})] - 1_v\sigma^2 - r(\mathbf{h}) \leq E[W(\mathbf{h}_n^*)] - 1_v\sigma^2 - r(\mathbf{h}_n^*) \right\},$$

for $n > N$. Thus,

$$\begin{aligned} & E \left[\text{WT}(\hat{\mathbf{h}}_n) \right] \left[1 - e(\hat{\mathbf{h}}_n) \right] \\ & \leq E \left[W(\hat{\mathbf{h}}_n) \right] - 1_v\sigma^2 - r(\hat{\mathbf{h}}_n) \\ & \leq E \left[W(\mathbf{h}_n^*) \right] - 1_v\sigma^2 - r(\hat{\mathbf{h}}_n) \\ & \leq E \left[\text{WT}(\mathbf{h}_n^*) \right] \left[1 + e(\mathbf{h}_n^*) \right] + \left[r(\mathbf{h}_n^*) - r(\hat{\mathbf{h}}_n) \right]. \end{aligned}$$

Further,

$$\begin{aligned} & \frac{E \left[\text{WT}(\hat{\mathbf{h}}_n) \right]}{E \left[\text{WT}(\mathbf{h}_n^*) \right]} \\ & \leq \left[\frac{1 + e(\mathbf{h}_n^*)}{1 - e(\hat{\mathbf{h}}_n)} \right] + \left\{ \frac{1}{E \left[\text{WT}(\mathbf{h}_n^*) \right]} \right\} \left[\frac{r(\mathbf{h}_n^*) - r(\hat{\mathbf{h}}_n)}{1 - e(\hat{\mathbf{h}}_n)} \right] \\ & \leq \left[\frac{1 + e(\mathbf{h}_n^*)}{1 - e(\hat{\mathbf{h}}_n)} \right] + \frac{|1_v - 1| (|\mu_{1_v}(\mathbf{h}_n^*)| + |\mu_{1_v}(\hat{\mathbf{h}}_n)|)}{|\min \{ \phi_W[\mu_{1_v}(\mathbf{h}_n^*)], \phi_W[\mu_{1_v}(\hat{\mathbf{h}}_n)] \} \mu_{2_v}(\mathbf{h}_n^*)|} \left[\frac{2}{1 - e(\hat{\mathbf{h}}_n)} \right]. \end{aligned}$$

As $\mu_{1_v}(\mathbf{h}_n)$ and $\mu_{1_v}^2(\mathbf{h}_n)/\mu_{2_v}(\mathbf{h}_n)$ tend to 0, $e(\mathbf{h}_n)$, $r(\mathbf{h}_n)$ and $p[\mu_{1_v}(\mathbf{h}_n)]/\mu_{2_v}(\mathbf{h}_n)$ tend to 0. Since $(1_v - 1)\mu_{1_v}(\mathbf{h}_n)/\mu_{2_v}(\mathbf{h}_n^*)$ tends to 0,

$$\frac{E \left[\text{WT}(\hat{\mathbf{h}}_n) \right]}{E \left[\text{WT}(\mathbf{h}_n^*) \right]} \rightarrow 1.$$

6.3 Appendix C: Theorem 1 and Corollary 1 applied to weighted smoothing

Denote $\mathbf{Z} = \mathbf{V}^{1/2}\mathbf{B}$ and let $\mathbf{Z} = \mathbf{U}\mathbf{D}\mathbf{Q}$ by singular value decomposition, where \mathbf{U} and \mathbf{Q} are orthogonal matrices and \mathbf{D} is an $n \times p_n$ matrix whose diagonal entries are the square roots of the eigenvalues $h_1, \dots, h_{\min(n, p_n)}$ of $\mathbf{Z}^t\mathbf{Z} = \mathbf{B}^t\mathbf{V}\mathbf{B}$ and with all other entries equal to 0.

As given in Sect. 6.1., $E[\text{WT}(\lambda)] = b^2(\lambda) + \sigma^2\mu_{2_v}(\lambda)$, where

$$\psi(\mathbf{V})b^2(\lambda) = \mathbf{f}^t [\mathbf{I} - \mathbf{H}_v(\lambda)]^t \mathbf{V} [\mathbf{I} - \mathbf{H}_v(\lambda)] \mathbf{f},$$

and $\psi(\mathbf{V})\mu_{2_v}(\lambda) = \text{Tr} [\mathbf{H}_v^t(\lambda)\mathbf{V}\mathbf{H}_v(\lambda)\mathbf{V}^{-1}]$.

Let s_{ij} be the (i, j) th element of the matrix $(\mathbf{D}^t \mathbf{D} + \sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t)^{-1}$. Then,

$$\begin{aligned} \max_{i,j} |s_{ij}| &\leq \left\| \left(\mathbf{D}^t \mathbf{D} + \sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right)^{-1} \right\|_2 \\ &= \sigma_{p_n}^{-1} \left(\mathbf{D}^t \mathbf{D} + \sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right) \\ &\leq \left[\sigma_{p_n}(\mathbf{D}^t \mathbf{D}) + \sigma_{p_n} \left(\sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right) \right]^{-1} \\ &\leq \frac{1}{\sum_{j=1}^k \lambda_j h_{j,p_n}^*}, \end{aligned}$$

and

$$\begin{aligned} s_{ii} &\geq \sigma_{p_n} \left[\left(\mathbf{D}^t \mathbf{D} + \sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right)^{-1} \right] \\ &= \sigma_1^{-1} \left(\mathbf{D}^t \mathbf{D} + \sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right) \\ &\geq \left[\sigma_1(\mathbf{D}^t \mathbf{D}) + \sigma_1 \left(\sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right) \right]^{-1} \\ &\geq \frac{1}{h_1 + \sum_{j=1}^k \lambda_j h_{j,1}^*}, \end{aligned}$$

where $\|\mathbf{M}_1\|_2$ is the matrix 2-norm of the matrix \mathbf{M}_1 and $\sigma_n(\mathbf{M}_1) \leq \sigma_{n-1}(\mathbf{M}_1) \leq \dots \leq \sigma_2(\mathbf{M}_1) \leq \sigma_1(\mathbf{M}_1)$ are the singular values (or eigenvalues) of the $n \times n$ symmetric matrix \mathbf{M}_1 . Then,

$$\begin{aligned} \psi(V)\mu_{1v}(\lambda) &= \text{Tr} \left[\mathbf{U} \mathbf{D} \left(\mathbf{D}^t \mathbf{D} + \sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right)^{-1} \mathbf{D}^t \mathbf{U}^t \right] \\ &= \sum_{i=1}^{\min(n,p_n)} h_i s_{ii} \end{aligned}$$

$$\begin{aligned} &\leq \frac{\sum_{i=1}^{\min(n, p_n)} h_i}{\sum_{j=1}^k \lambda_j h_{j, p_n}^*} \\ &\leq c_0 n^{q_2} \sum_{i=1}^{\min(n, p_n)} i^{-q_1} \\ &\approx c_0 c_1 n^{q_2}, \end{aligned}$$

and

$$\begin{aligned} \psi(\mathbf{V})\mu_{2v}(\boldsymbol{\lambda}) &= \text{Tr} \left\{ \left[\mathbf{D} \left(\mathbf{D}^t \mathbf{D} + \sum_{j=1}^k \lambda_j \mathbf{Q} \mathbf{P}_j \mathbf{Q}^t \right)^{-1} \mathbf{D}^t \right]^2 \right\} \\ &= \sum_{i=1}^{\min(n, p_n)} \sum_{j=1}^{\min(n, p_n)} h_i h_j s_{ij}^2 \\ &\geq \sum_{i=1}^{\min(n, p_n)} h_i^2 s_{ii}^2 \\ &\geq \sum_{i=1}^{\min(n, p_n)} \left(\frac{h_i}{h_1 + \sum_{j=1}^k \lambda_j h_{j,1}^*} \right)^2 \\ &\approx c_2 \sum_{i=1}^{\min(n, p_n)} i^{-2q_1} \\ &\approx c_2 c_3, \end{aligned}$$

where c_0, c_1, c_2 and c_3 are some constants. Thus, both $\mu_{1v}(\boldsymbol{\lambda})$ and $0 \leq \mu_{2v}(\boldsymbol{\lambda})/\mu_{2v}(\boldsymbol{\lambda}) \leq c_0^2 c_1^2 n^{2q_2} / [c_2 c_3 \psi(\mathbf{V})]$ tend to 0 as n tends to infinity. Further, both $e(\boldsymbol{\lambda})$ and $r(\boldsymbol{\lambda})$ tend to 0. Therefore, $W(\boldsymbol{\lambda})$ is nearly an unbiased estimator of $WP(\boldsymbol{\lambda})$ implied by Theorem 1. By differentiating $\psi(\mathbf{V})b^2(\boldsymbol{\lambda})$ and $\psi(\mathbf{V})\mu_{2v}(\boldsymbol{\lambda})$,

$$\left[\frac{\psi(\mathbf{V})\partial b^2(\boldsymbol{\lambda})}{\partial \lambda_j} \right]_{\boldsymbol{\lambda}=0} = -2\mathbf{f}^t [\mathbf{I} - \mathbf{H}_v(0)]^t \mathbf{V} \left[\frac{\partial \mathbf{H}_v(\boldsymbol{\lambda})}{\partial \lambda_j} \right]_{\boldsymbol{\lambda}=0} \mathbf{f} = 0,$$

and $\psi(\mathbf{V})\partial \mu_{2v}(\boldsymbol{\lambda})/\partial \lambda_j < 0$. Thus, $E[\text{WT}(\boldsymbol{\lambda})]$ has strictly negative gradient at $\boldsymbol{\lambda} = 0$ and

$$\lim_{\boldsymbol{\lambda} \rightarrow \infty} E[\text{WT}(\boldsymbol{\lambda})] = [1/\psi(\mathbf{V})] \sum_{i=1}^n \tilde{f}_i^2 > 0,$$

where $\mathbf{V}^{1/2} \mathbf{f} = \tilde{\mathbf{f}} = (\tilde{f}_1, \dots, \tilde{f}_n)^t$. These implies that $E[\text{WT}(\boldsymbol{\lambda})]$ has at least one minimizer $\boldsymbol{\lambda}_n^* > 0$. Finally, if $1_v - 1 = O(n^{-q})$ and $q \geq 1/2$, Corollary 1 holds.

Acknowledgments The author would like to thank two referees and the editor for helpful suggestions that led to a substantial improvement in the paper. This research is partly supported by Taiwan NSC Grant (Project: NSC 95-2118-M-029-005).

References

- Akaike, H. (1970). Statistical predictor identification. *Annals of the Institute of Statistical Mathematics*, 22, 203–217.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19, 716–723.
- Altman, N. S. (1990). Kernel smoothing of data with correlated errors. *Journal of the American Statistical Association*, 85, 749–759.
- Craven, P., Wahba, G. (1979). Smoothing noisy data with spline functions. *Numerische Mathematik*, 31, 377–403.
- Eilers, P. H. C., Marx, B. D. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science*, 11, 89–121.
- Eubank, R. L. (1988). *Spline smoothing and nonparametric regression*. New York: Dekker.
- Golub, G. H., Van Loan, C. F. (1993). *Matrix computations*. Baltimore and London: The Johns Hopkins University Press.
- Green, P. J., Silverman, B. W. (1994). *Nonparametric regression and generalized linear models: a roughness penalty approach*. London: Chapman and Hall.
- Härdle, W., Müller, M., Sperlich, S., Werwatz, A. (2004). *Nonparametric and semiparametric models*. Berlin: Springer.
- Hastie, T. J., Tibshirani, R. J. (1990). *Generalized additive models*. New York: Chapman and Hall.
- Hocking, R. R. (1976). The analysis and selection of variables in linear regression. *Biometrics*, 32, 1–49.
- Li, K. C. (1987). Asymptotic optimality for C_p , C_L , cross-validation and generalized cross-validation: discrete index set. *Annals of Statistics*, 15, 958–975.
- Opsomer, J., Wang, Y., Yang, Y. (2001). Nonparametric regression with correlated errors. *Statistical Science*, 16, 134–153.
- O’Sullivan, F., Yandell, B., Raynor, W. (1986). Automatic smoothing of regression function in generalized linear models. *Journal of the American Statistical Association*, 81, 96–103.
- Rice, J. (1984). Bandwidth choice for nonparametric regression. *Annals of Statistics*, 12, 1215–1230.
- Shibata, R. (1981). An optimal selection of regression variables. *Biometrika*, 68, 45–54.
- Wahba, G. (1990). *Spline models for observational data*. Philadelphia: Society for Industrial and Applied Mathematics.
- Wang, Y. (1998). Smoothing spline models with correlated random errors. *Journal of the American Statistical Association*, 93, 341–348.
- Wei, W. H. (2005). The smoothing parameter, confidence interval and robustness for smoothing splines. *Journal of Nonparametric Statistics*, 17, 613–642.