

On hybrid methods of inverse regression-based algorithms

Li-Xing Zhu^{a, b, *}, Megu Ohtaki^c, Yingxing Li^d

^a*Hong Kong Baptist University, Hong Kong, China*

^b*Renmin University of China, Beijing, China*

^c*Research Institute of Radiation Biology, Hiroshima, Japan*

^d*Cornell University, New York, USA*

Received 10 March 2005; received in revised form 4 August 2005; accepted 20 December 2005

Available online 3 February 2006

Abstract

This paper is two-fold. First, we present a further investigation for the hybrid methods of inverse regression-based algorithms. This investigation provides the evidence of how the hybrids gain the advantages to become more powerful methods than the existing methods when the central dimension reduction (CDR) space is estimated. Second, a Bayes Information Criterion (BIC)-type algorithm is recommended to estimate the dimension of the CDR space. Differing from the popularly used sequential test methods, the new algorithm does not require the asymptotic normality of the estimator of the inverse regression-based matrix. The BIC-based estimator is proven to be consistent. A set of simulations for several typical models were carried out to guide the selection of coefficient in the hybrids.

© 2006 Published by Elsevier B.V.

Keywords: Dimension reduction; Sliced inverse regression; Sliced average variance estimation; Hybrid methods

1. Introduction

Consider regression problem with the response Y and p -dimensional covariable $\mathbf{x} = (x_1, \dots, x_p)^T$. In sufficient dimension reduction, the main object of interest is the intersection $S_{Y|\mathbf{x}}$ of all subspaces $S \subseteq R^p$ such that Y is independent of \mathbf{x} when $P_S \mathbf{x}$ is given, where $P_S \mathbf{x}$ is the projection of \mathbf{x} onto the subspace S . $S_{Y|\mathbf{x}}$ is called central dimension reduction (CDR) space. An example is the following model proposed by Li (1991):

$$Y = f(B^T \mathbf{x}, \varepsilon), \quad (1.1)$$

where $f(\cdot)$ is an unknown function and ε is the error term independent of \mathbf{x} , and $B = (\beta_1, \dots, \beta_k)$ is an unknown $p \times k$ matrix whose columns are of unitary length in L^2 -norm and orthogonal to one another. This model includes many popularly used models, for example, the single-index model (Härdle et al., 1993) and the projection pursuit regression model (Friedman and Stuetzle, 1981).

To estimate the CDR space, two inverse regression-based algorithms are popularly used: sliced inverse regression (SIR) (Li, 1991) and sliced average variance estimate (SAVE) (Cook and Weisberg, 1991; Cook, 2000). These

* Corresponding author. Department of Mathematics, Hong Kong Baptist University, Kowloon Tong, Hong Kong, China. Tel.: +86 3411 7016; fax: +86 3411 5811.

E-mail address: lzhu@hkbu.edu.hk (L.-X. Zhu).

two approaches are, respectively, based on the inverse mean of \mathbf{x} given Y and the inverse conditional variance of \mathbf{x} given Y . To be precise, denote $\mathbf{z} = \Sigma_{\mathbf{x}}^{-1/2}(\mathbf{x} - E(\mathbf{x}))$ and the conditional covariance matrix of \mathbf{z} , given y by $Cov(\mathbf{z}|y) = E\{[(\mathbf{z} - E(\mathbf{z}|Y))(\mathbf{z} - E(\mathbf{z}|Y))^T | Y = y]\}$. Adopting the notion of Cook (1998, pp. 105, 188, 197), we have $S_{Y|\mathbf{z}} = \Sigma_{\mathbf{x}}^{1/2} S_{Y|\mathbf{x}}$. This is the dimension-reduction subspace for the regression of Y on \mathbf{z} . Furthermore, let $S_{E(\mathbf{z}|y)}$ be the inverse regression (IR) subspace spanned by $E(\mathbf{z}|y) := E(\mathbf{z}|Y = y)$ for all y , and let $S_{I_p - Cov(\mathbf{z}|y)}$ be the inverse conditional variance (IRV) subspace spanned by $I_p - Cov(\mathbf{z}|y)$ for all y . *SIR* is used to estimate $S_{E(\mathbf{z}|y)}$, and *SAVE* is used to estimate $S_{I_p - Cov(\mathbf{z}|y)}$. As Li (1991) proved, if the following linear condition is satisfied:

$$E(\mathbf{z} | P_{S_{Y|\mathbf{z}}} \mathbf{z}) = P_{S_{Y|\mathbf{z}}} \mathbf{z}, \quad (1.2)$$

where $P_{(\cdot)}$ stands for the projection operator in the standard inner product (see Cook, 1998), then for all y , $E(\mathbf{x}|y)$ lies within the space $S_{Y|\mathbf{z}}$. Hence, $S_{E(\mathbf{z}|y)}$ is a subspace of $S_{Y|\mathbf{z}}$ and $\dim(S_{E(\mathbf{z}|y)}) \leq \dim(S_{Y|\mathbf{z}})$, where $\dim(\cdot)$ stands for dimension. Similarly, $S_{I_p - Cov(\mathbf{z}|y)}$ is also contained in $S_{Y|\mathbf{z}}$ (see Cook and Weisberg, 1991) if, together with condition (1.2),

$$Cov(\mathbf{z} | P_{S_{Y|\mathbf{z}}} \mathbf{z}) = I_p - P_{S_{Y|\mathbf{z}}}. \quad (1.3)$$

Both *SIR* and *SAVE* have pros and cons. *SIR* is relatively insensitive to the choice of slice number. Li (1991) demonstrated this through the numerical study. Zhu and Ng (1995) proved that under some regularity conditions, the asymptotic normality of *SIR* holds in a large range of number of slices, that is, from $O(\sqrt{n})$ to $n/2$. Zhu and Fang (1996) studied kernel estimation, in which the bandwidth selection is also flexible in a fairly large range. When the regression is an odd function, *SIR* works well, as indicated by the simulation in Li (1991). However, when the regression function is symmetric or the CDR space relates to the variance of error, *SIR* does not work well. A typical example is $y = (\beta_1^T \mathbf{z})^2 + (\beta_2^T \mathbf{z})^2 \varepsilon$, where \mathbf{z} and ε have the respective normal distributions $N(0, I_p)$ and $N(0, 1)$, $\beta_1 = (1, 0, \dots, 0)^T$ and $\beta_2 = (0, 1, \dots, 0)^T$, that is, $y = z_1^2 + z_2^2 \varepsilon$. We can easily prove that $E(\mathbf{z}|y) \equiv 0$ through $E(\mathbf{z} | z_1^2, z_2^2, \varepsilon) \equiv 0$ and $E(\mathbf{z}|y) = E(E(\mathbf{z} | z_1^2, z_2^2, \varepsilon) | y)$, therefore, $\dim(S_{E(\mathbf{z}|y)}) = 0$. This issue was also investigated by Li (1991, Remark 4.5) and the discussants of Li's paper. If $y = f_1(\beta_1^T \mathbf{z}) + f_2(\beta_2^T \mathbf{z}) \varepsilon$, and $f_i(\cdot)$ are even functions with respect to $\beta_i^T \mathbf{z}$, then $\dim(S_{E(\mathbf{z}|y)})$ may be zero. As to *SAVE*, it can handle the above model well. The IRV space may not be a null space as has been discussed by Li (1991, Remark 4.5 and Rejoinder) and the discussants of Li's paper (see, e.g. Cook and Weisberg, 1991). Cook (2000) investigated *SAVE* in more detail. Furthermore, Cook and Critchley (2000) and Ye and Weiss (2003) showed that under the regularity conditions IR space is a subspace of IRV space. Therefore, theoretically, under regularity conditions, *SAVE* should be a more powerful method to estimate CDR space than *SIR*. However, *SAVE* is less robust because it uses higher moments, and is sensitive to the choice of the number of slices. Cook (2000) heuristically stated that, like a tuning parameter in non-parametric estimation, the number of data per slice should be large enough to allow a reasonable estimation of the intra-slice conditional covariance matrices of \mathbf{z} given y . Our empirical studies also indicate this (see Section 5 below).

Hence, how to gain the advantages from these two algorithms is of interest. When investigating the relationship among *SIR*, *SIR_{II}* and *SAVE* in the rejoinder to the discussions on his paper, Li (1991) suggested a hybrid algorithm in an ad hoc manner, that is, a hybrid algorithm of *SIR*² and *SIR_{II}*, that is, $SIR_{II_{a_1}} = (1 - a_1) SIR^2 + a_1 SIR_{II}$ for $0 \leq a_1 \leq 1$. In this hybrid, *SIR* is indirectly used. Gannoun and Saracco (2003a) studied the asymptotic properties of $SIR_{II_{a_1}}$ and Gannoun and Saracco (2003b) considered a single-index model when $SIR_{II_{a_1}}$ is used. Ye and Weiss (2003) considered the hybrid method of *SIR* and pHd, and briefly discussed the hybrid of *SIR* and *SAVE*. That is, the algorithms $(1 - a)SIR + apHd$, and $(1 - a)SIR + aSAVE$ are used to estimate CDR space. Although *SAVE* is more comprehensive than *SIR*, Cook and Critchley (2000) and Ye and Weiss (2003) pointed out that the increased flexibility comes with a price that relatively straightforward structure that is manifest through the mean $E(\mathbf{z}|y)$ is harder to detect with *SAVE* than with *SIR*. Therefore, properly selecting the coefficient may make a more powerful algorithm to estimate CDR space.

Hence with the hybrid idea, there are two issues: how to properly select the coefficient, and how to estimate the dimension of CDR space. Ye and Weiss (2003) proposed a bootstrap selection for the coefficient when the hybrid of *SIR* and pHd is used, and estimation for the dimension. However, when slicing estimation is used, the consistency of the bootstrap estimators is still an open problem. This is because the concomitants \mathbf{z}_i associated with the ordered y_i are

Download English Version:

<https://daneshyari.com/en/article/417357>

Download Persian Version:

<https://daneshyari.com/article/417357>

[Daneshyari.com](https://daneshyari.com)