

Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Journal of the Korean Statistical Society

journal homepage: www.elsevier.com/locate/jkss

Fused sliced average variance estimation

Hyoin An, Sungmin Won, Jae Keun Yoo *

Department of Statistics, Ewha Womans University, Seoul 03760, Republic of Korea

ARTICLE INFO

Article history:

Received 9 March 2017

Accepted 30 June 2017

Available online xxxx

AMS 2000 subject classifications:

62G08

62H05

Keywords:

Fusing

Inverse regression

Sliced average variance estimation

Sufficient dimension reduction

ABSTRACT

In this paper, we propose an approach to combine the kernel matrices constructed by sliced average variance estimation (SAVE) with various numbers of slices. The proposed approach is called fused sliced average variance estimation (FSAVE). By fusing the information by usual SAVE applications with different slice numbers, the sensitivity to slices can be reduced, so the structural dimension estimation can be improved. Numerical studies confirm this, and a real data analysis is presented.

© 2017 The Korean Statistical Society. Published by Elsevier B.V. All rights reserved.

1. Introduction

The main goal of sufficient dimension reduction (SDR) is to replace the original predictor \mathbf{X} with a lower-dimensional linearly transformed predictor $\mathbf{M}^T\mathbf{X}$ without losing any information about the distribution of $Y|\mathbf{X}$. This is equivalently expressed as $Y \perp\!\!\!\perp \mathbf{X}|\mathbf{M}^T\mathbf{X}$, where $\perp\!\!\!\perp$ indicates stochastic independence. The central subspace $\mathcal{S}_{Y|\mathbf{X}}$ is the intersection of all subspaces $\mathcal{S}(\mathbf{M})$ satisfying $Y \perp\!\!\!\perp \mathbf{X}|\mathbf{M}^T\mathbf{X}$, where $\mathcal{S}(\mathbf{A})$ stands for a subspace spanned by the columns of $\mathbf{A} \in \mathbb{R}^{p \times q}$. Then SDR pursues the estimation of $\mathcal{S}_{Y|\mathbf{X}}$.

Two methods of sliced inverse regression (SIR; Li, 1991) and sliced average variance estimation (SAVE; Cook and Weisberg, 1991) are classical but most popular among other SDR methods. We will discuss SAVE in detail in later section. To implement SIR and SAVE in practice, a categorization of Y , called *slicing*, is crucial. Critical drawback of the slicing is that there are no optimal numbers which can be used as a thumb rule. Many studies show that improper numbers of the slices mislead the results. This problem is more severe in SAVE than in SIR, because the former utilizes the information of the second moment of $\mathbf{X}|Y$, but the latter does the first moment.

Cook and Zhang (2014) proposed a simple solution to this problem with the application to SIR. They combine sample kernel matrices of SIR from many different numbers of slices. According to Cook and Zhang (2014), the combining procedure to SIR results in equally good basis estimates for the different numbers of slices, varying from 3 to 15.

It should be noted that the Cook–Zhang approach is applicable to SAVE. The primary interest of the paper is to develop a fused approach for SAVE. It is not guaranteed that the fused SAVE always yields robust basis estimates to various numbers of slices like SIR. However, more information on $\mathcal{S}_{Y|\mathbf{X}}$ can be obtained by combining the sample kernel matrices by SAVE, so it can be expected to provide robust basis estimates for mild numbers of slices, say 3 to 6 and to have better asymptotic behaviors in the dimension estimation of $\mathcal{S}_{Y|\mathbf{X}}$ than usual SAVE application.

* Corresponding author.

E-mail address: peter.yoo@ewha.ac.kr (J.K. Yoo).

The organization of the paper is as follows. In Section 2, the method of SAVE and required conditions are briefly discussed. A fused approach of SAVE is proposed in Section 3. Numerical studies and a real data example are presented in Section 4. We summarize our work in Section 5.

2. Sliced average variance estimation

To explain SAVE, the original predictor \mathbf{X} is standardized as follows: $\mathbf{Z} = \Sigma^{-1/2}(\mathbf{X} - E(\mathbf{X}))$, where $\Sigma = \text{cov}(\mathbf{X})$. Define $S_{Y|Z}$ as the central subspace for a regression of $Y|Z$. Denote d as the true structured dimension of $S_{Y|X}$ and $S_{Y|Z}$. Let η and η_z be $p \times d$ orthonormal basis matrix for $S_{Y|X}$ and $S_{Y|Z}$, respectively. Then, we have $S_{Y|X} = \Sigma^{-1/2}S_{Y|Z}$, equivalently, $\eta = \Sigma^{-1/2}\eta_z$. For rigorous proofs of the relationship between $S_{Y|X}$ and $S_{Y|Z}$, one can read Cook (1998, proposition 6.3).

Consider the following two conditions called linearity and constant variance condition, respectively:

A1. $E(\mathbf{Z}|\mathbf{P}_{S_{\eta_z}}\mathbf{Z})$ is linear in $\mathbf{P}_{\eta_z}\mathbf{Z}$.

A2. $\text{cov}(\mathbf{Z}|\mathbf{P}_{\eta_z}\mathbf{Z})$ is constant.

If \mathbf{X} has an elliptically contoured distribution, condition A1 is satisfied and A2 approximately holds. If \mathbf{X} follows a multivariate normal distribution, A1 and A2 are guaranteed to hold. In the case that the two conditions fail, \mathbf{X} can often be one-to-one transformed to induce these conditions.

Then conditions A1 and A2 force the following relationship to hold for $\text{cov}(\mathbf{Z}|Y)$:

$$\mathbf{I}_p - \text{cov}(\mathbf{Z}|Y) = \mathbf{P}_{\eta_z} \{ \mathbf{I}_p - \text{cov}(\mathbf{Z}|Y) \} \mathbf{P}_{\eta_z} \in S_{Y|Z}. \tag{1}$$

Rigorous discussion on Eq. (1) is given in Cook and Weisberg (1991). Therefore, a proper subset of $S_{Y|Z}$ can be produced by varying y in $\mathbf{I}_p - \text{cov}(\mathbf{Z}|Y = y)$:

$$S\{\mathbf{I}_p - \text{cov}(\mathbf{Z}|Y)\} = S[E\{\mathbf{I}_p - \text{cov}(\mathbf{Z}|Y)\}^2] \subseteq S_{Y|Z}.$$

An approach to estimate $S_{Y|Z}$ through $\mathbf{M}_{\text{SAVE}} = E\{\mathbf{I}_p - \text{cov}(\mathbf{Z}|Y)\}^2$ is called sliced average variance estimation (SAVE; Cook and Weisberg, 1991). According to Ye and Weiss (2003, section 2.2) SAVE provides a more comprehensive estimation of $S_{Y|X}$ than SIR, which utilizes the information of $E(\mathbf{Z}|Y)$, in sense that a subspace constructed by the former contains that by the latter. The estimation of $\text{cov}(\mathbf{Z}|Y)$ by using the exact distribution of $\mathbf{Z}|Y$ is not realistic in practice, not only because the regression is for studying $Y|Z$, not $\mathbf{Z}|Y$, but also because a strong parametric assumption for $\mathbf{Z}|Y$ is required. Note that $\text{cov}(\mathbf{Z}|Y)$ can be nonparametrically replaced by usual sample moment estimators of the covariance of \mathbf{Z} within each category of Y , if Y is categorical. Following this, the estimation of $\text{cov}(\mathbf{Z}|Y)$ would be simple through a categorization of Y , called *slicing*, if Y is many-valued or continuous. First obtain \tilde{Y} by slicing Y with h levels. Hereafter, h stands for the number of slices in the SAVE application. Then compute the sample inverse covariance within each slice $s, s = 1, \dots, h$:

$$\widehat{\text{cov}}(\mathbf{Z}|\tilde{Y} = s) = \frac{1}{n_s} \sum_{\tilde{Y}_i = s} (\hat{\mathbf{Z}}_{i \in s} - \bar{\mathbf{Z}}_s)(\hat{\mathbf{Z}}_{i \in s} - \bar{\mathbf{Z}}_s)^T,$$

where n_s is the sample size of the category s and $\bar{\mathbf{Z}}_s = \frac{1}{n_s} \sum_{\tilde{Y}_i = s} \hat{\mathbf{Z}}_i$.

Then construct a sample kernel matrix $\hat{\mathbf{M}}_{\text{SAVE}}$:

$$\hat{\mathbf{M}}_{\text{SAVE}} = \sum_{s=1}^h \frac{n_s}{n} \left\{ \mathbf{I}_p - \widehat{\text{cov}}(\mathbf{Z}|\tilde{Y} = s) \right\} \left\{ \mathbf{I}_p - \widehat{\text{cov}}(\mathbf{Z}|\tilde{Y} = s) \right\}.$$

Then, by spectral decomposition of $\hat{\mathbf{M}}_{\text{SAVE}}$, the eigenvectors corresponding to non-zero eigenvalues form an orthonormal basis of $S_{Y|Z}$.

3. Fused approach of SAVE

In this section, we propose a fused approach of SAVE. We start the section with defining that

$$\mathbf{M}_{\text{FSAVE}}^{(h)} = (\mathbf{M}_{\text{SAVE}}^{(2)}, \dots, \mathbf{M}_{\text{SAVE}}^{(h)}), \tag{2}$$

where $\mathbf{M}_{\text{SAVE}}^{(h)}$ stands for the kernel matrix of SAVE with h slices. The case of $h = 1$ is naturally excluded, because it simply returns null matrix. Also the case of $h = 2$ will not be considered, because $\mathbf{M}_{\text{FSAVE}}^{(2)} = \mathbf{M}_{\text{SAVE}}^{(2)}$.

Theoretically, we have $S(\mathbf{M}_{\text{SAVE}}^{(k)}) \subseteq S_{Y|Z}$ for $k = 2, 3, \dots, h$, which directly implies that

$$S(\mathbf{M}_{\text{SAVE}}^{(k)}) \subseteq S(\mathbf{M}_{\text{FSAVE}}^{(h)}) \subseteq S_{Y|Z}, \quad k = 2, \dots, h.$$

Therefore, $\mathbf{M}_{\text{FSAVE}}^{(h)}$ is a possible kernel matrix to estimate $S_{Y|Z}$, so we have

$$\Sigma^{-1/2} S(\mathbf{M}_{\text{FSAVE}}^{(h)}) \subseteq S_{Y|X}.$$

Download English Version:

<https://daneshyari.com/en/article/7546314>

Download Persian Version:

<https://daneshyari.com/article/7546314>

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