

Representative approach for big data dimension reduction with binary responses

Xuelong Wang, Jie Yang

May 16, 2019

- 1 Background and Challenges
- 2 A novel approach: Representative
- 3 Simulation results
- 4 Aysmptotic property

Sufficient Dimension Reduction

Model Settings

$$Y = f(X, \epsilon) = f(\beta_1 X, \dots, \beta_k X, \epsilon)$$

x is explanatory variable, column vectors on \mathcal{R}^p ,

β' 's are unknown row vectors,

ϵ is independent of X ,

f is an **arbitrary unknown** function on \mathcal{R}^{k+1}

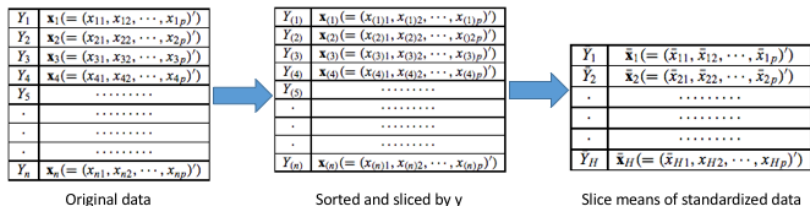
- $(\beta_1 X, \dots, \beta_k X)'$ is the projection of the $X \in \mathcal{R}^p$ into \mathcal{R}^K ,
 $K \ll p$
- **Lower dimension projection of X contains most of the information**

Sufficient Dimension Reduction

The space spanned by β 's

- ① Effective dimension-reduction direction (e.d.r)
 - A Linear combination of β 's
 - ② A Linear space \mathcal{B} :
 - Spanned by β 's ($Span(\beta)$) \Leftrightarrow All the possible linear combination of β 's
-
- Since f is arbitrary, f and β 's are not Estimable
 - Only the \mathcal{B} can be identified
 - Inverse Regression is one of the methods of estimating the Effective dimension-reduction space (\mathcal{B})

Sliced Inverse Regression method



$$\hat{V} = n^{-1} \sum_{h=1}^H n_h \bar{x}_h \bar{x}_h^T \quad \xrightarrow{\text{Conduct PCA on } \hat{V} \text{ Find the first Kth eigenvectors } \hat{\eta}} \quad \hat{\beta}_k = \hat{\eta}_k \Sigma_{xx}^{-1/2}$$

Estimated Covariance matrix

Binary response

The curse of Binary response for the sliced-based methods

Since the responses only have to values, the number slices

- For the method using only first moment, only one direction can be recovered
- For the method using more than moment(SAVE), there are situations they cannot recover all the directions

Existing approaches

- Using SVM to estimate the pseudo conditional probability

Representative approach

Data generating models for binary response

Let Y^* as the latent response and Y as the observed binary response,

$$Y = \begin{cases} 0 & Y^* < \theta \\ 1 & Y^* \geq \theta \end{cases}$$

Where θ is the cutoff value.

$$Y^* = f(X^T \beta_1, \dots, X^T \beta_p, \epsilon),$$

Where ϵ is a random variable. Thus, we have

$$\mathcal{P}(Y = 1|X) = \Pr(Y^* > \theta|X).$$

Note that different distribution ϵ will eventually affect the distribution of Y .

Data generating models

Conditional probability model

$$\lambda(\mathcal{P}(Y = 1|X)) = h(\beta_1^T X, \dots, \beta_p X)$$

$$\begin{aligned}\mathcal{P}(Y = 1|X) &= \lambda^{-1} \circ h(\beta_1^T X, \dots, \beta_p X) \\ &= g(\beta_1^T X, \dots, \beta_p X)\end{aligned}$$

Troubles for traditional methods

Simulation results

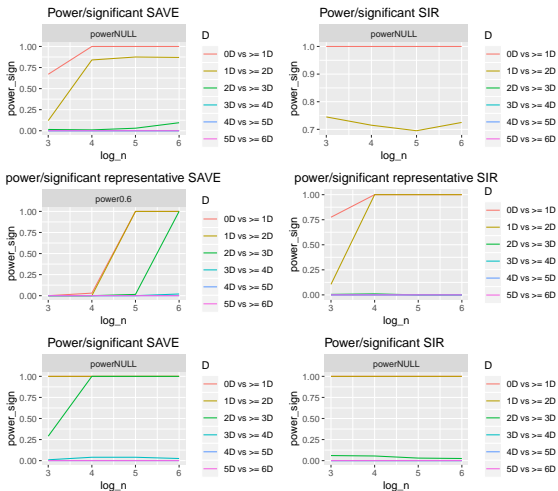


Figure 2:

Simulation results

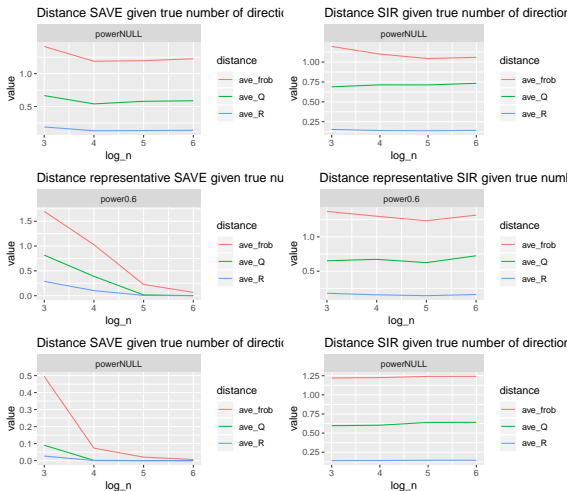


Figure 2:

The conditional expection and representative

The choice of number of cluster K