

# Sufficient reductions in regression with mixed predictors

Efstathia Bura

ASTAT

Institute of Statistics and Mathematical Methods in Economics

joint with Liliana Forzani, Rodrigo García Arancibia, Pamela Llop and Diego Tomassi at  
Universidad Nacional del Litoral.



TECHNISCHE  
UNIVERSITÄT  
WIEN  
Vienna | Austria

1. The Problem
2. Our approach
3. Sufficient reductions for the regression of  $\mathbf{Y}$  on  $\mathbf{X}$  and  $\mathbf{H}$
4. Governance index application
5. References

1. The Problem
2. Our approach
3. Sufficient reductions for the regression of  $\mathbf{Y}$  on  $\mathbf{X}$  and  $\mathbf{H}$
4. Governance index application
5. References

- ▶ Most data sets comprise of measurements on continuous and categorical variables.

- ▶ Most data sets comprise of measurements on continuous and categorical variables.
  - ▶ **Facebook** is collecting data on *media mix*, i.e., factors that may have influence over sales, both continuous and discrete, and is using them to quantify the weight of each factor to create a model to predict marketing results for future strategy.

- ▶ Most data sets comprise of measurements on continuous and categorical variables.
  - ▶ **Facebook** is collecting data on *media mix*, i.e., factors that may have influence over sales, both continuous and discrete, and is using them to quantify the weight of each factor to create a model to predict marketing results for future strategy.
- ▶ Little work has been previously done, mostly for the unsupervised problem (*location model* of Olkin and Tate **1961**, Gaussian Graphical Model (GGM))

- ▶ Most data sets comprise of measurements on continuous and categorical variables.
  - ▶ **Facebook** is collecting data on *media mix*, i.e., factors that may have influence over sales, both continuous and discrete, and is using them to quantify the weight of each factor to create a model to predict marketing results for future strategy.
- ▶ Little work has been previously done, mostly for the unsupervised problem (*location model* of Olkin and Tate **1961**, Gaussian Graphical Model (GGM))

1. The Problem
2. Our approach
3. Sufficient reductions for the regression of  $\mathbf{Y}$  on  $\mathbf{X}$  and  $\mathbf{H}$
4. Governance index application
5. References



Suppose the response  $Y$  is either continuous or categorical. We consider the conditional distribution

$$Y \mid (\mathbf{X}, \mathbf{H}), \quad (1)$$

where  $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$  is a vector of  $p$  continuous, and  $\mathbf{H} = (H_1, H_2, \dots, H_q)^T$  is a vector of  $q$  binary predictor variables.

- ▶ We want to find a lower dimensional function  $\mathbf{R} : \mathbb{R}^p \times \mathbb{R}^q \rightarrow \mathbb{R}^d$  of the mixed predictor vector  $\mathbf{Z} = (\mathbf{X}^T, \mathbf{H}^T)^T$  such that

$$F(Y \mid \mathbf{Z}) = F(Y \mid \mathbf{R}(\mathbf{Z})),$$

where  $F(\cdot \mid \cdot)$  denotes the conditional cumulative distribution function of the response given the predictors.

- ▶  $\mathbf{R}$  is called a **sufficient reduction** of the regression of  $Y$  on  $\mathbf{Z}$ .

1. The Problem
2. Our approach
3. Sufficient reductions for the regression of  $\mathbf{Y}$  on  $\mathbf{X}$  and  $\mathbf{H}$
4. Governance index application
5. References

## A Minimal Sufficient Reduction

for the regression of  $Y$  on  $(\mathbf{X}, \mathbf{H})$  is given by

$$\mathbf{R}(\mathbf{X}, \mathbf{H}) = \boldsymbol{\alpha}_{\mathbf{b}}^T (\mathbf{t}(\mathbf{X}, \mathbf{H}) - \mathbf{E}(\mathbf{t}(\mathbf{X}, \mathbf{H})))$$

where

$$\mathbf{t}(\mathbf{X}, \mathbf{H}) = \left( \mathbf{X}^T, \mathbf{H}^T, (\mathbf{J}_q \text{vech}(\mathbf{H}\mathbf{H}^T))^T \right)^T$$

and  $\boldsymbol{\alpha}_{\mathbf{b}}$  is a basis for  $\mathcal{S}_{\mathbf{b}} = \text{span}\{\mathbf{b}\}$  with

$$\mathbf{b} = \begin{pmatrix} \Delta^{-1} \mathbf{A} \\ \mathbf{L}_q \boldsymbol{\tau} - \boldsymbol{\beta}^T \Delta^{-1} \mathbf{A} \\ \mathbf{J}_q \boldsymbol{\tau} \end{pmatrix} : \text{to be estimated}$$

1. We assume a random sample  $(y_i, \mathbf{x}_i, \mathbf{h}_i)$ ,  $i = 1, \dots, n$ , is drawn from the joint distribution of  $(Y, \mathbf{X}, \mathbf{H})$

1. We assume a random sample  $(y_i, \mathbf{x}_i, \mathbf{h}_i)$ ,  $i = 1, \dots, n$ , is drawn from the joint distribution of  $(Y, \mathbf{X}, \mathbf{H})$
2. We derive
  - ▶ the MLEs of the parameters  $\Delta, \mu, \mu_{\mathbf{H}}, \mathbf{A}, \beta, \tau_0, \tau$ , and the MLE of the reduction,

1. We assume a random sample  $(y_i, \mathbf{x}_i, \mathbf{h}_i)$ ,  $i = 1, \dots, n$ , is drawn from the joint distribution of  $(Y, \mathbf{X}, \mathbf{H})$
2. We derive
  - ▶ the MLEs of the parameters  $\Delta, \mu, \mu_{\mathbf{H}}, \mathbf{A}, \beta, \tau_0, \tau$ , and the MLE of the reduction,
  - ▶ the asymptotic normality of  $\hat{\alpha}_{\mathbf{b}}$  and asymptotic tests for dimension, and

1. We assume a random sample  $(y_i, \mathbf{x}_i, \mathbf{h}_i)$ ,  $i = 1, \dots, n$ , is drawn from the joint distribution of  $(Y, \mathbf{X}, \mathbf{H})$
2. We derive
  - ▶ the MLEs of the parameters  $\Delta, \mu, \mu_{\mathbf{H}}, \mathbf{A}, \beta, \tau_0, \tau$ , and the MLE of the reduction,
  - ▶ the asymptotic normality of  $\hat{\alpha}_{\mathbf{b}}$  and asymptotic tests for dimension, and
  - ▶ a regularized estimator for simultaneous variable (feature) selection and dimension reduction (feature extraction).
3. Details in Bura et al. 2022 and code at [https://github.com/lforzani/SDR\\_mixed\\_predictions](https://github.com/lforzani/SDR_mixed_predictions)

1. The Problem
2. Our approach
3. Sufficient reductions for the regression of  $Y$  on  $X$  and  $H$
4. Governance index application
5. References



We want to build a *Composite Governance index* (CG),

$$\mathbf{R}(\mathbf{Z}) = \boldsymbol{\omega}^T \mathbf{Z} \in \mathbb{R},$$

to predict  $Y$  = the logarithm of per capita *Gross Domestic Product* (GDP), measured in 2010 US dollars, over the period 1996 to 2018.

- ▶ The World Bank considers the following six aggregate indicators of governance that combine the views of a large number of enterprise, citizen and expert survey respondents:
  - ▶ control of corruption ( $X_1$ );
  - ▶ rule of law ( $X_2$ );
  - ▶ regulatory quality ( $X_3$ );
  - ▶ government effectiveness ( $X_4$ );
  - ▶ political stability ( $X_5$ );
  - ▶ voice and accountability ( $X_6$ ).

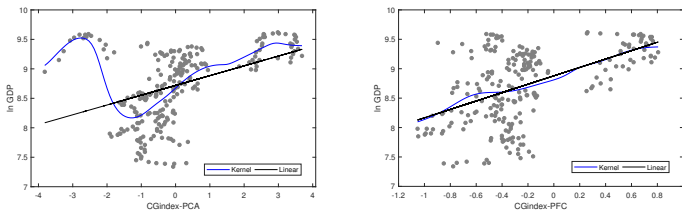


Figure: Log of per capita GDP versus PCA and PFC (Y-targeted PCA) based composite governance indexes.

We add country effect by introducing eleven binary variables  $\mathbf{H}$ . We plot the log of GDP versus the CG index constructed by PCA for mixed variables (PCAMIX) in the left panel and by our mixed OPTIMAL SDR approach in the right panel.

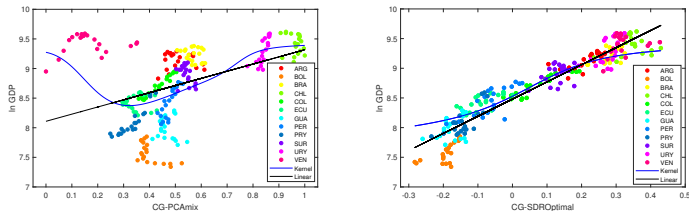


Figure: Log of per capita GDP versus Composite Governance index with country effect.

1. The Problem
2. Our approach
3. Sufficient reductions for the regression of  $\mathbf{Y}$  on  $\mathbf{X}$  and  $\mathbf{H}$
4. Governance index application
5. References



Bura, E., L. Forzani, R. G. Arancibia, P. Llop, and D. Tomassi (2022). “Sufficient reductions in regression with mixed predictors”. In: *Journal of Machine Learning Research* 23.102, pp. 1–47.



Olkin, I. and R. Tate (1961). “Multivariate Correlation Models with Mixed Discrete and Continuous Variables”. In: *Ann. Math. Statist.* 32.2, pp. 448–465.

**Thanks!**