

Federal Committee on  
Statistical Methodology Research Conference  
November, 2013

Prediction Performance of Single Index Principal  
Fitted Component Models

Jia-Ern Pai

&

Kofi Placid Adragi, PhD



# Outline

- Introduction
- Inverse Reduction Models
  - Principal Component model
  - Principal Fitted Component model
- Principal Fitted Component Regression
- Prediction Comparisons
- Simulation Study
  - Simulation Setup
  - Simulation Results
- Conclusion



# Introduction

- Frequently encountered problems in regression analyses
  - Large  $p$  and small  $n$  problems
    - When  $n < p$ , the OLS regression cannot provide stable parameter estimates
    - Example:  
Modeling the effect of upgraded Fuel System Integrity:  
age of vehicles, vehicle types, manufactures, model years, impact speeds, and driver's ages...etc.
  - Collinearity problems
    - Inflated variances of estimates and predictions
    - Example:  
age of vehicles and model year



Dimension reduction is necessary

# Introduction (cont.)

- $\mathbf{X}$  is a  $p$ -vector random predictors,  $Y$  is an univariate response variable
  - Forward regression methods are usually adopted
    - Denoted as  $Y | \mathbf{X}$
    - Examples:  
Partial Least Squares regression, LASSO regression, and principal component analysis, ..., etc.
- Utilizing the randomness property of  $\mathbf{X}$ 
  - Inverse regression models would be another options.
    - Denoted as  $\mathbf{X} | Y$
- Inverse reduction methods often provide more regression information between  $\mathbf{X}$  and  $Y$

Research interest is on prediction performances of the Signal-Index-Isotropic Principal Fitted Component models



# Introduction (cont.)

- Why are the inverse reduction methods more informative?
- Example: Principal Component Analysis
- Given a design matrix  $\mathbf{X} \in \mathbf{R}^{n \times p}$ 
  - From the sample  $\text{cov}(\mathbf{X})$ , we obtain:
    - Eigenvalues:  $\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_p$
    - Eigenvectors:  $\hat{\gamma}_1, \hat{\gamma}_2, \dots, \hat{\gamma}_p$
    - Principal components:  $\{\hat{\gamma}_1^T \mathbf{X}, \hat{\gamma}_2^T \mathbf{X}, \dots, \hat{\gamma}_p^T \mathbf{X}\}$
- Main drawbacks of principal component analysis:
  - The response variable is not involved
  - The dimension reduction cannot be conducted when  $\hat{\lambda}_i \approx \hat{\lambda}_j, \forall i, j$ , such that  $i \neq j$



# Introduction (cont.)

- Goal of Dimension Reduction
  - $\mathbf{X} \in \mathbf{R}^p$  and  $\mathbf{R}(\mathbf{X}) \in \mathbf{R}^d$ , such that  $d \leq p$
- What do we expect?
  - $\mathbf{R}(\mathbf{X})$  carries as much regression information as  $\mathbf{X}$  on  $Y$

- Original regression model:

$$Y | \mathbf{X} = \boldsymbol{\alpha}^T \mathbf{X} + \boldsymbol{\varepsilon} \quad (1)$$

- Replace  $\mathbf{X}$  by  $\mathbf{R}(\mathbf{X})$  without losing any regression information

$$Y | \mathbf{X} = \boldsymbol{\beta}^T \mathbf{R}(\mathbf{X}) + \boldsymbol{e} \quad (2)$$



# Introduction (cont.)

- Definition of the sufficient reduction (Cook, 2007)
- $\mathbf{X} \in \mathbf{R}^p$  and  $\mathbf{R}(\mathbf{X}) \in \mathbf{R}^d$ , such that  $d \leq p$ 
  - Inverse reduction,  $\mathbf{X} \mid (Y, \mathbf{R}(\mathbf{X})) \sim \mathbf{X} \mid \mathbf{R}(\mathbf{X})$
  - Forward reduction,  $Y \mid \mathbf{X} \sim Y \mid \mathbf{R}(\mathbf{X})$
  - Joint reduction,  $\mathbf{X}$  is independent of  $Y \mid \mathbf{R}(\mathbf{X})$
- If any condition holds, then  $\mathbf{R}(\mathbf{X})$  is a sufficient reduction



# Inverse Reduction Models

- Principal Component models
- Suppose  $\mathbf{X} \in \mathbf{R}^p$ , we regress  $\mathbf{X}$  on  $Y$

$$\mathbf{X} | y = \boldsymbol{\mu} + \boldsymbol{\Gamma} \boldsymbol{\nu}_y + \boldsymbol{\varepsilon} \quad (3)$$

- $\boldsymbol{\Gamma}$  is a semi-orthogonal matrix:

$$\boldsymbol{\Gamma}^T \boldsymbol{\Gamma} = \mathbf{I}_d, \text{ such that } d \leq p$$

- $\boldsymbol{\nu}_y$  is an unknown function of  $y$
- $\boldsymbol{\Gamma}^T \mathbf{X}$  is a sufficient reduction in the Principal Component model





# Inverse Reduction Models (cont.)

- Cook (2007) assumed  $\mathbf{v}_y = \boldsymbol{\beta} \mathbf{f}_y$ 
  - $\mathbf{f}_y$  is a known flexible basis function of  $y$
  - In practice,  $\mathbf{f}_y$  can be determined by polynomial basis functions or piecewise polynomial basis functions
- Principal Fitted Component (PFC) models

$$\mathbf{X} | y = \boldsymbol{\mu} + \boldsymbol{\Gamma} \boldsymbol{\beta} \mathbf{f}_y + \boldsymbol{\varepsilon} \quad (4)$$

- $\boldsymbol{\Gamma}^T \mathbf{X}$  is still a sufficient reduction
- PFC model is model-based in this research
- We assume  $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_p)$ 
  - $\text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_p$ : Isotropic error term



# Inverse Reduction Models (cont.)

- To have fair and straightforward prediction performance comparisons with the OLS and LASSO regressions
  - Only one principal fitted component is used
  - Set  $\mathbf{f}_y = y$
- Single-Index-Isotropic Principal Fitted Component model

$$\mathbf{X} | y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\beta y + \boldsymbol{\varepsilon} \quad (6)$$

- $\boldsymbol{\Gamma} \in \mathbf{R}^{p \times 1}$  and  $\boldsymbol{\Gamma}^T \boldsymbol{\Gamma} = 1$
- $\beta \in \mathbf{R}$  and  $E[Y] = 0$
- $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_p)$

We only concern the isotropic PFC model in this research, so the term “single index PFC model” is adopted in the rest presentation content



# Inverse Reduction Models (cont.)

- $\mathbf{X} | y = \boldsymbol{\mu} + \Gamma\beta y + \boldsymbol{\varepsilon}$  (6)

- A sufficient reduction in the single index PFC model is not unique

- Example:

- $\Gamma^T \mathbf{X} \equiv a\Gamma^T \mathbf{X}$  in the single index PFC model, if  $a \neq 0$

- However,  $\text{span}(\Gamma)$  is unique

- $\text{span}(\Gamma) = \text{span}(a\Gamma)$

- We should estimate  $\text{span}(\Gamma)$  instead of  $\Gamma$

- We still need to have a  $\Gamma$  before finding  $\text{span}(\Gamma)$

Parameter space in the single index PFC model:

$(\boldsymbol{\mu}, \Gamma, \beta, \sigma^2)$

- Estimated by MLE



# Principal Fitted Component Regression

- Consider a forward linear regression model

$$Y | \mathbf{X} = \boldsymbol{\alpha}^T \mathbf{X} + \boldsymbol{e} \quad (6)$$

- $\boldsymbol{\Gamma}^T \mathbf{X}$  is a sufficient reduction in the single index PFC model

$$Y | \mathbf{X} = \beta(\boldsymbol{\Gamma}^T \mathbf{X}) + \boldsymbol{\varepsilon} \quad (7)$$

- $\hat{\boldsymbol{\Gamma}}$  is obtained from the single index PFC model

- Denote  $\mathbf{Z}$  as  $\hat{\boldsymbol{\Gamma}}^T \mathbf{X}$

$$Y | \mathbf{Z} = \gamma \mathbf{Z} + \boldsymbol{\varepsilon}^* \quad (8)$$

- Like model (6), model (7) a simple linear regression model
  - $\hat{\boldsymbol{\Gamma}}^T \mathbf{X}$  is proxy of  $\mathbf{X}$

The procedure of replacing  $\mathbf{X}$  by a sufficient reduction in a forward regression is called the Principal Fitted Component Regression (PFCR)



# Prediction Comparisons

- How to make predictions with PFCR?
- Given two data set  $(\mathbf{X}, Y)$  and  $(\mathbf{X}^*, Y^*)$ 
  - $(\mathbf{X}, Y)$  is used for the model building
  - $(\mathbf{X}^*, Y^*)$  is used for making predictions
- $(\mathbf{X}, Y)$  and  $(\mathbf{X}^*, Y^*)$  are generated in the same way
- How to assess the prediction performance?
- The sample mean squared prediction error (PE) is adopted

$$- \text{PE} = \frac{1}{n} \sum_i (Y_i^* - \hat{E}(Y | \hat{\Gamma}^T \mathbf{X}_i^*))^2 \quad (9)$$



# Simulation Study

- Purpose

Compare the prediction performances of the single index PFC model with other forward methods, such as the OLS, Ridge, LASSO, and Partial Least Square (PLS) regressions

- We use single index PLS model to make fair comparisons

- Scenarios

- $n > p$  problem

- Large  $n$  case: All the predictors are response-related

- $n < p$  problems

- Dense case: All the predictors are response-related

- Sparse case: Only some of predictors are response-related



# Simulation Study (cont.)

- Data generation:
- We make  $\mathbb{X}$  as a linear function of  $\mathbf{Y}$

$$\mathbb{X} = \beta(\mathbf{\Gamma}\mathbf{Y})^T + \boldsymbol{\varepsilon} \quad (10)$$

- $\mathbf{\Gamma} \in \mathbf{R}^{p \times 1}$ ,  $\beta \in \mathbf{R}$
- $\mathbf{Y} = (y_1, y_2, \dots, y_n)$ , such that  $y_i \sim N(0, \sigma_y^2)$
- $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_p)$
- $\beta$  determines the strength of association between  $\mathbb{X}$  and  $\mathbf{Y}$
- $\beta$  is large,  $\mathbb{X}$  and  $\mathbf{Y}$  can reveal sufficient regression information to each other
  - Inverse and forward dimension reduction models should be able to find sufficient reductions more easily

Changing different values of  $\beta$ ,  $p$ , and  $n$ , distinct scenarios are created



# Simulation Study (cont.)

- Iterating 100 times data generations, model buildings, PE calculations for every model in each distinct scenario
- From 100 PE's
  - $\overline{PE}$ 's and  $SE(PE)$ 's can be calculated
- We present  $\overline{PE}$ 's and  $SE(PE)$ 's as simulation results





# Simulation Study (cont.)

- Large  $n$  case
- Simulation set up

	level 1	level 2	level 3
$n$	100	300	600
$\beta$	0.1	0.4	1
$p$	25		
$\sigma^2$	1		
$\sigma_y^2$	1		

- $\Gamma \in \mathbf{R}^{25 \times 1}$
- $\Gamma = \left( \frac{1}{\sqrt{25}}, \frac{1}{\sqrt{25}}, \dots, \frac{1}{\sqrt{25}} \right)^T$



- Large  $n$  case

Beta = 0.1					
	OLS	LASSO	Ridge	PLS	PFC
$n = 100$	1.32(0.020)	1.00(0.014)	1.31(0.020)	1.17(0.016)	1.15(0.015)
$n = 300$	1.08(0.010)	1.00(0.009)	1.08(0.010)	1.06(0.009)	1.05(0.008)
$n = 600$	1.04(0.006)	1.00(0.004)	1.03(0.006)	1.03(0.005)	1.03(0.006)

(1)

Beta = 0.4					
	OLS	LASSO	Ridge	PLS	PFC
$n = 100$	1.15(0.018)	1.00(0.013)	1.14(0.018)	1.02 (0.019)	1.03(0.021)
$n = 300$	0.97(0.008)	0.94(0.008)	0.97(0.008)	0.94(0.009)	0.95(0.009)
$n = 600$	0.93(0.006)	0.91(0.006)	0.92(0.005)	0.90(0.006)	0.91(0.005)

(2)

Beta = 1					
	OLS	LASSO	Ridge	PLS	PFC
$n = 100$	0.66(0.018)	0.63(0.007)	0.65(0.015)	0.60(0.010)	0.60(0.009)
$n = 300$	0.56(0.006)	0.53(0.006)	0.55(0.005)	0.53(0.006)	0.53(0.005)
$n = 600$	0.53(0.003)	0.52(0.003)	0.53(0.003)	0.52(0.003)	0.52(0.003)

(3)



# Simulation Study (cont.)

- Dense case
- Simulation set up

	level 1	level 2	level 3
$p$	100	200	400
$\beta$	0.1	0.4	1
$n$	100		
$\sigma^2$	1		
$\sigma_y^2$	1		

- $\mathbf{\Gamma} \in \mathbf{R}^{p \times 1}$
- $\mathbf{\Gamma} = (1, 1, \dots, 1)^T$



- Dense case

	Beta=0.1		
	LASSO	PLS	PFC
$p = 100$	5.60(2.000)	0.73(0.010)	0.74(0.011)
$p = 200$	0.82(0.015)	0.61(0.010)	0.61(0.011)
$p = 400$	0.72(0.012)	0.49(0.009)	0.50(0.009)

(4)

	Beta=0.4		
	LASSO	PLS	PFC
$p = 100$	133.82(130.234)	0.07(0.001)	0.07(0.001)
$p = 200$	0.10(0.002)	0.04(0.001)	0.04(0.001)
$p = 400$	0.11(0.002)	0.02(0.000)	0.02(0.000)

(5)

	Beta=1		
	LASSO	PLS	PFC
$p = 100$	0.70(0.152)	0.01(0.000)	0.01(0.000)
$p = 200$	0.02(0.000)	0.01(0.000)	0.01(0.000)
$p = 400$	0.02(0.000)	0.003(0.000)	0.003(0.000)

(6)



# Simulation Study (cont.)

- Sparse case
- We only compare the prediction performances of the sparse single index PFCR to LASSO regression
  - PLS regression does not have coefficient shrinkage procedure
- Simulation set up

	level 1	level 2	level 3
$p$	100	200	400
$\beta$	0.1	0.4	1
$p_0$	10		
$n$	100		
$\sigma^2$	1		
$\sigma_y^2$	1		

- $P_0$  is the number of active predictors
- $\Gamma \in \mathbf{R}^{p \times 1}$
- $\Gamma = (1, \dots, 1, 0, \dots, 0)^T$



- Sparse case
- PLS model is not considered , because it does not have a threshold procedure

	Beta=0.1	
	LASSO	Sparse PFC
$p = 100$	1.41(0.281)	1.15(0.020)
$p = 200$	1.05(0.020)	1.20(0.019)
$p = 400$	1.04(0.017)	1.16(0.017)

(7)

	Beta=0.4	
	LASSO	Sparse PFC
$p = 100$	1.92(0.595)	0.48(0.008)
$p = 200$	0.54(0.010)	0.53(0.009)
$p = 400$	0.57(0.010)	0.62(0.010)

(8)

	Beta=1	
	LASSO	Sparse PFC
$p = 100$	0.61(0.100)	0.10(0.002)
$p = 200$	0.12(0.002)	0.11(0.002)
$p = 400$	0.12(0.002)	0.13(0.003)

(9)



# Conclusion

- In some large  $n$  case, not all predictors are active
  - The PFCR is preferred
- The prediction performances of the signal index PFCR and signal index PLS regression are almost the same
  - Only under the assumption  $f_y = y$
  - The PFCR is more flexible
- It seems that the LASSO regression provides unstable prediction performances when  $p$  is close to  $n$ 
  - Sparse PFCR is recommended
  - “lars” is used when simulating the prediction performance of the LASSO regression



# Reference

- Cook, R. D. (2007). Fisher Lecture: Dimension Reduction in Regression. *Statistical Science*, Vol. 22, No. 1, 1-26.
- Cook, R. D. & Forzani, L. (2008). Principal Fitted Components for Dimension Reduction in Regression. *Statistical Science*, Vol. 23, No. 4, 485-501.
- Adraghi, K. P. & Cook, R. D. (2009). Sufficient dimension reduction and prediction in regression. *Philosophical Transactions of the Royal Society A*, Vol. 367, No. 1906, 4385-4405.
- Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society Series B (Methodological)*, Vol. 58, No. 1, 267-288.
- Hoerl, A. E. & Kennard, R. W. (1970). Ridge Regression: Applications to Nonorthogonal Problems. *Technometrics*, Vol. 12, No. 1, 69-82.
- Refaeilzadeh, P., Tang L., & Liu H. (2009) Cross-Validation. *Encyclopedia of Database Systems*, Part 3, 532-538.
- Anderson T. W. (2003). Principal Components. *An Introduction to Multivariate Statistical Analysis*, Chapter 11, 495-482.



# Thank you

[Jia-ern.pai@dot.org](mailto:Jia-ern.pai@dot.org)

