

Model-Based Sufficient Dimension Reduction for Regression

R.D. Cook

School of Statistics

University of Minnesota



Isaac Newton Institute

January 2008



In collaboration with Francesca Chiaromonte, Liliana Forzani, Bing Li and Lexin Li.

Fisher (1924), “The Influence of Rainfall on the Yield of Wheat at Rothamsted”

Opening issue: $n < p$.

Standard methodology of the time: Convert $n < p$ to $n > p^*$ by using marginal regressions.

Fisher’s conclusion: “The meteorological variables to be employed must be chosen without reference to the actual crop record.”

Principal Components in Regression

Goal : Reduce the dimension of the predictor vector $\mathbf{X} \in \mathbb{R}^p$ prior to studying the regression of Y on \mathbf{X} .

PC's : Replace \mathbf{X} with $\hat{\gamma}_1^T \mathbf{X}, \dots, \hat{\gamma}_d^T \mathbf{X}$ for some $d < p$, $\hat{\gamma}_1, \dots, \hat{\gamma}_p$ ordered eigenvectors of $\hat{\Sigma} = \widehat{\text{var}}(\mathbf{X})$.

Historical reasons :

- mitigate the impact of collinearity
- facilitate computation and exploratory analysis
- facilitate analysis when $n < p$ (Chiaromonte & Martinelli 2002; Li & Li, 2004; Bair, et al. 2006, supervised principal components)

Cox (1968):

A difficulty seems to be that there is no logical reason why the dependent variable should not be closely tied to the least important principal component.

Some agreed with Cox, providing examples of relevant “least important principal components.”

Mosteller and Tukey (1977) asked

...how can we find linear combinations of the [predictors] that will be likely, or unlikely, to pick up regression from some as yet unspecified y ?

and answered

A malicious person who knew our x 's and our plan for them could always invent a y to make our choices look horrible. But we don't believe nature works that way – more nearly that nature is, as Einstein put it (in German), “tricky, but not downright mean.”

On balance, the role for PC's in reg. seems elusive: (a) they are computed from the marginal of X , and (b) are not usefully invariant or equivariant.

Today they are ubiquitous in the applied sciences ... but not widely represented in Statistics – IR methods, lasso, MAVE, PLS, projection pursuit, contour regression, directional regression.

Sufficient Reductions

Variables: $Y \in \mathbb{R}^1, \mathbf{X} \in \mathbb{R}^p, (Y, \mathbf{X}) \sim F$

Data: (Y_i, \mathbf{X}_i) iid $F, i = 1, \dots, n.$

Goal: Reduce the dimension of \mathbf{X} without loss of relevant information on $Y|\mathbf{X}$.

Sufficient Reduction: A reduction $R : \mathbb{R}^p \rightarrow \mathbb{R}^q, q \leq p,$ is sufficient if any of the following hold:

1. $\mathbf{X}|(Y, R(\mathbf{X})) \sim \mathbf{X}|R(\mathbf{X})$ (Inverse reg.)
2. $Y|\mathbf{X} \sim Y|R(\mathbf{X})$ (Forward reg.)
3. $Y \perp\!\!\!\perp \mathbf{X}|R(\mathbf{X})$ (Joint)

We will pursue reductions via #1.

1. $\mathbf{X}|(Y, R(\mathbf{X})) \sim \mathbf{X}|R(\mathbf{X})$, 2. $Y|\mathbf{X} \sim Y|R(\mathbf{X})$, 3. $Y \perp\!\!\!\perp \mathbf{X}|R(\mathbf{X})$

“Model-free” inverse regression $\mathbf{X}|Y$, #1. .

The central subspace $\mathcal{S}_{Y|\mathbf{X}}: \alpha \in \mathbb{R}^{p \times d}$

$$R(\mathbf{X}) = \alpha^T \mathbf{X} \rightarrow \text{span}(\alpha) \rightarrow \mathcal{S}_{Y|\mathbf{X}} = \cap \text{span}(\alpha).$$

**SIR, SAVE, PHD, IHT, IR methods, contour reg.,
directional reg., Fourier methods, moment methods,
partial methods, kernel methods ...**

Parametric inverse regression $\mathbf{X}|Y$, #1.

Going back. $\mathbf{X}|Y \rightarrow Y|\mathbf{X}$ when (\mathbf{X}, Y) has a density,

$$\hat{E}(Y|\mathbf{X}_{\text{new}}) = \frac{\sum_{i=1}^n y_i \hat{f}\{R(\mathbf{X}_{\text{new}})|Y = y_i\}}{\sum_{i=1}^n \hat{f}\{R(\mathbf{X}_{\text{new}})|Y = y_i\}}$$

1. $\mathbf{X}|(Y, R(\mathbf{X})) \sim \mathbf{X}|R(\mathbf{X})$, 2. $Y|\mathbf{X} \sim Y|R(\mathbf{X})$, 3. $Y \perp\!\!\!\perp \mathbf{X}|R(\mathbf{X})$

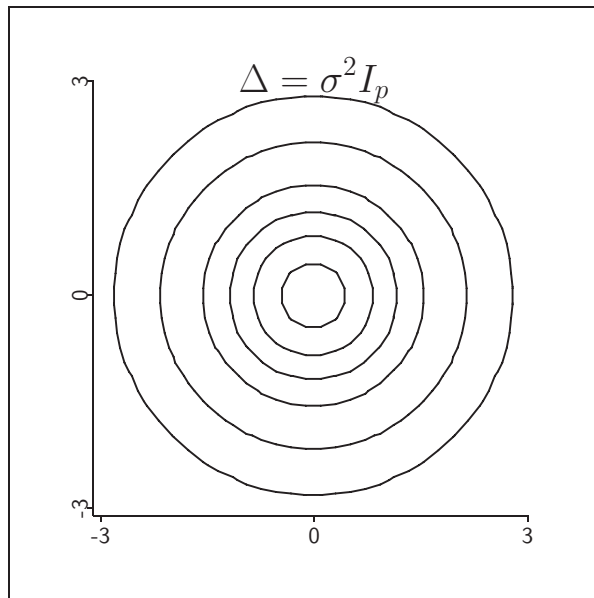
A First Nonlinear Inverse Model

$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\nu}_y + \sigma\boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, I_p)$$

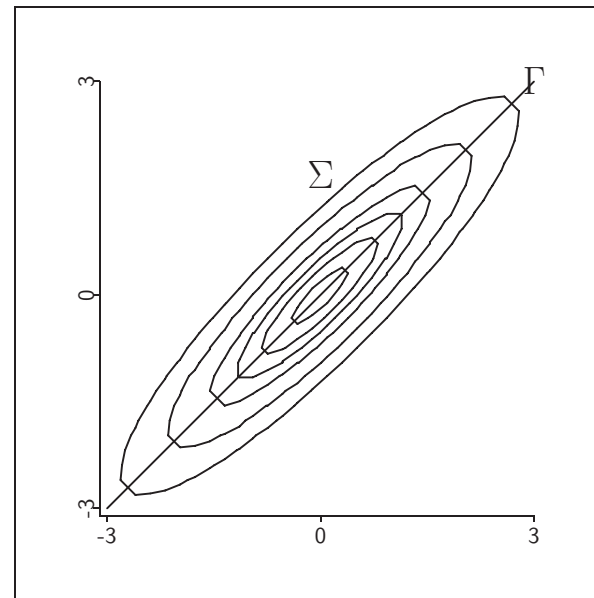
- $\mathbf{X}_y \sim \mathbf{X}|(Y = y)$
- $\boldsymbol{\Gamma} \in \mathbb{R}^{p \times d}$, $\boldsymbol{\nu}_y \in \mathbb{R}^d$, and d all unknown.
- $R(\mathbf{X}) = \boldsymbol{\alpha}^T \mathbf{X}$, $\boldsymbol{\alpha}$ any basis for $\mathcal{S}_{Y|\mathbf{X}} \equiv \text{span}(\boldsymbol{\Gamma})$ is a *minimal sufficient reduction*

$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\nu}_y + \sigma\boldsymbol{\varepsilon}$$

Using MLE, $\hat{\mathcal{S}}_{Y|\mathbf{X}} = \text{span}\{\hat{\boldsymbol{\gamma}}_1, \dots, \hat{\boldsymbol{\gamma}}_d\}$; $R(\mathbf{X})$ is estimated by 1st d sample PC's, $\hat{R}(\mathbf{X}) = (\hat{\boldsymbol{\gamma}}_1^T \mathbf{X}, \dots, \hat{\boldsymbol{\gamma}}_d^T \mathbf{X})$; role of Y .

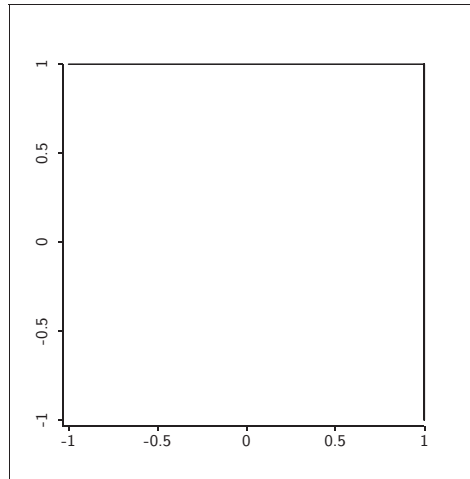


$$\Delta = \text{Var}(\mathbf{X}|Y)$$



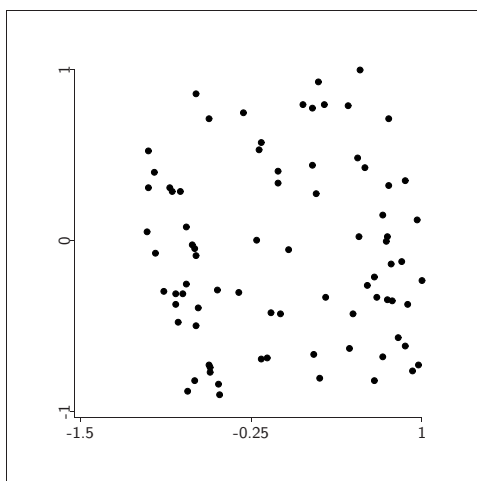
$$\Sigma = \text{Var}(\mathbf{X})$$

$\boldsymbol{\nu}_y, n = 80, d = 2$

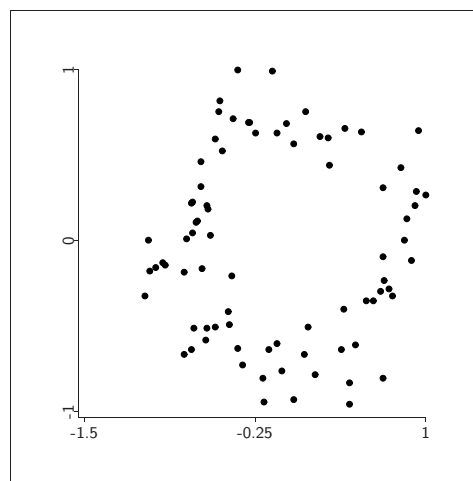


$$\begin{aligned}\mathbf{X}_y &= \mathbf{\Gamma}\boldsymbol{\nu}_y + .25N(0, I_p) \\ \text{vec}(\mathbf{\Gamma}) &\sim N(0, I_{2p})\end{aligned}$$

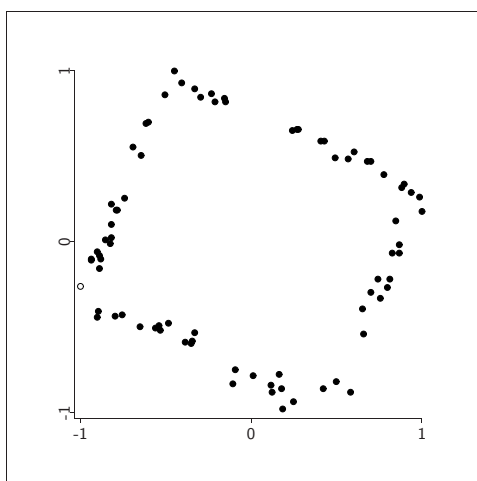
Next slide: Plots of $\hat{R}(\mathbf{X}_i) : 2 \times 1, i = 1, \dots, 80$.



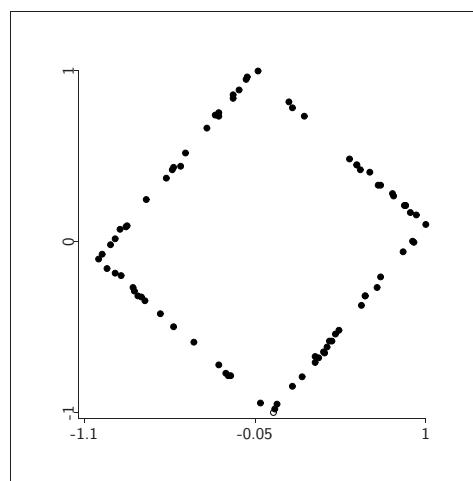
$p = 3$



$p = 5$



$p = 25$



$p = 500$

Extensions of $\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\nu}_y + \sigma\boldsymbol{\varepsilon}$

- **Model $\boldsymbol{\nu}_y = \boldsymbol{\beta}\mathbf{f}_y$; $\boldsymbol{\beta} \in \mathbb{R}^{d \times r}$, $d \leq r$, $\mathbf{f}_y \in \mathbb{R}^r$ known.**
- **More flexibility in $\boldsymbol{\Delta} = \text{Var}(\mathbf{X}_y)$.**

Principal Fitted Components

$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_y + \sigma\boldsymbol{\varepsilon}$$

- $R(\mathbf{X}) = \boldsymbol{\alpha}^T \mathbf{X}$ still a minimal sufficient reduction, $\boldsymbol{\alpha}$ any basis for $\mathcal{S}_{Y|\mathbf{X}} = \text{span}(\boldsymbol{\Gamma})$.

- **ML Estimation:** Let $\hat{\boldsymbol{\Sigma}}_{\text{fit}}$ sample covariance matrix of fitted values from the linear regression fit of \mathbf{X}_y on \mathbf{f}_y .

Then $\hat{\mathcal{S}}_{Y|\mathbf{X}}$ is the span of the first d eigenvectors of $\hat{\boldsymbol{\Sigma}}_{\text{fit}}$, leading to **Principal Fitted Components**.

- **Robustness:** \sqrt{n} consistency under
 - Normal $\boldsymbol{\varepsilon} \rightarrow t_5, \chi_5, U(0, 1)$.
 - Approximation of $\mathbf{f}_y \in \mathbb{R}^r$ with \mathbf{g}_y provided $\text{rank}\{\text{cov}(\mathbf{f}_Y, \mathbf{g}_Y)\} = r$.

Application: Gene Expression (Leek & Storey, Sept. 2007):

$$\mathbf{X}_{y,\mathbf{u}} = \mu + \beta \mathbf{f}_y + \Gamma_2 \boldsymbol{\nu}_{\mathbf{u}} + \sigma \boldsymbol{\varepsilon}$$

$$\mathbf{X}_{y,\mathbf{c},\mathbf{u}} = \mu + \Gamma \beta \mathbf{f}_y + \Gamma_1 \beta_1 \mathbf{g}_{\mathbf{c}} + \Gamma_2 \boldsymbol{\nu}_{\mathbf{u}} + \sigma \boldsymbol{\varepsilon}$$

\mathbf{c} = known confounders

\mathbf{u} = unknown confounders

$\mathbf{f}_y, \mathbf{g}_{\mathbf{c}}$ = known functions of y, \mathbf{c} .

Adding $\Delta \equiv \text{Var}(\mathbf{X}_y) > 0$

$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_y + \boldsymbol{\Delta}^{1/2}\boldsymbol{\varepsilon}, \quad \boldsymbol{\Delta} > 0$$

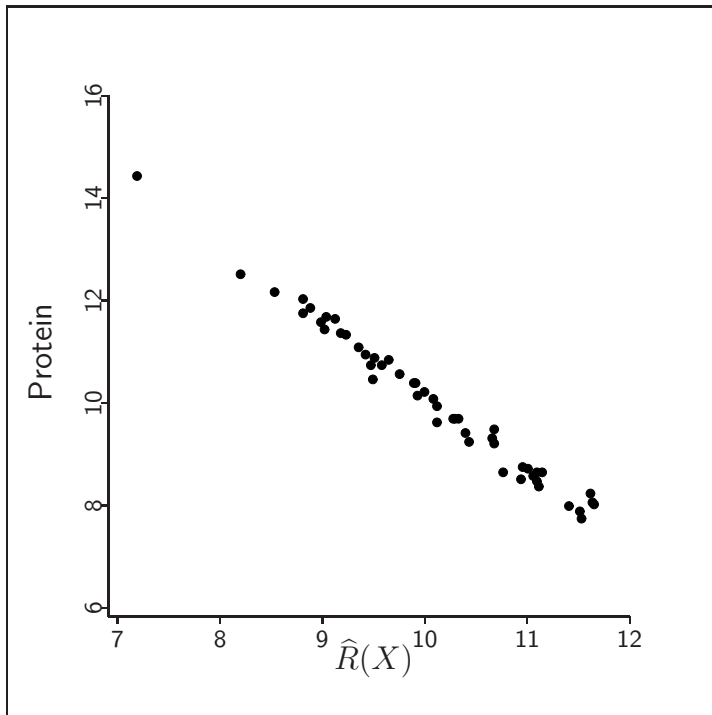
- $R(\mathbf{X}) = \boldsymbol{\alpha}^T \mathbf{X}$, $\boldsymbol{\alpha}$ any basis for $\mathcal{S}_{Y|\mathbf{X}} = \boldsymbol{\Delta}^{-1}\text{span}(\boldsymbol{\Gamma})$.
- $R(\mathbf{A}\mathbf{X}) \equiv R(\mathbf{X})$ and adapted to Y .
- \hat{R} computed as PFC's based on $\hat{\boldsymbol{\Sigma}}^{-1/2}\mathbf{X}$: Inner product of first d eigenvectors of $\hat{\boldsymbol{\Sigma}}^{-1/2}\hat{\boldsymbol{\Sigma}}_{\text{fit}}\hat{\boldsymbol{\Sigma}}^{-1/2}$ with $\hat{\boldsymbol{\Sigma}}^{-1/2}\mathbf{X}$.
- Robustness to non-normality and misspecification of \mathbf{f}_y still hold, $\text{rank}\{\text{cov}(\mathbf{f}_Y, \mathbf{g}_Y)\} = r$.
- SIR

Adding $\Delta \equiv \text{var}(\mathbf{X}_y) > 0$

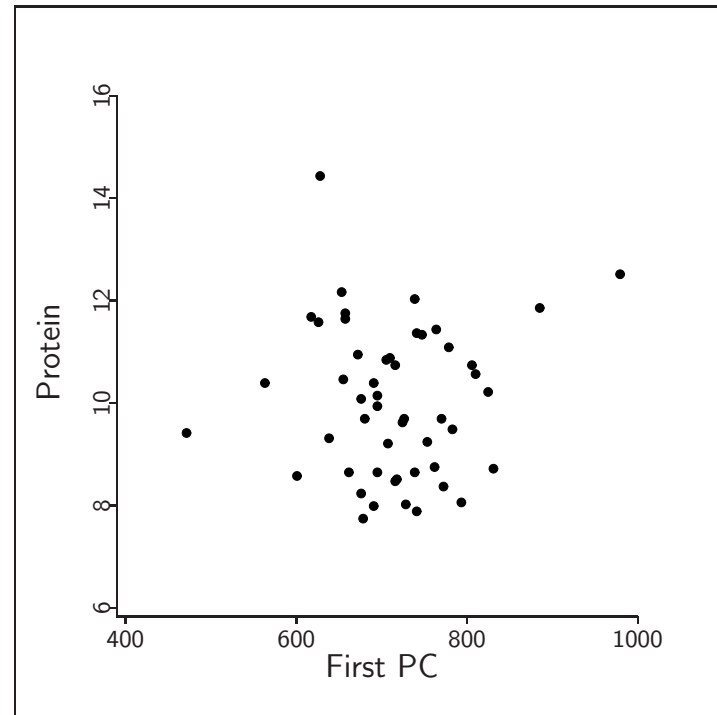
$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_y + \boldsymbol{\Delta}^{1/2}\boldsymbol{\varepsilon}, \quad \boldsymbol{\Delta} > 0$$

- $R(\mathbf{X}) = \boldsymbol{\alpha}^T \mathbf{X}$, $\boldsymbol{\alpha}$ any basis for $\mathcal{S}_{Y|\mathbf{X}} = \boldsymbol{\Delta}^{-1}\text{span}(\boldsymbol{\Gamma})$.
- $R(\mathbf{A}\mathbf{X}) \equiv R(\mathbf{X})$ and adapted to Y .
- \hat{R} computed as PFC's based on $\hat{\boldsymbol{\Sigma}}^{-1/2}\mathbf{X}$: Inner product of first d eigenvectors of $\hat{\boldsymbol{\Sigma}}^{-1/2}\hat{\boldsymbol{\Sigma}}_{\text{fit}}\hat{\boldsymbol{\Sigma}}^{-1/2}$ with $\hat{\boldsymbol{\Sigma}}^{-1/2}\mathbf{X}$.
- Robustness to non-normality and misspecification of \mathbf{f}_y still hold, $\text{rank}\{\text{cov}(\mathbf{f}_Y, \mathbf{g}_Y)\} = r$.
- SIR gives the MLE when Y is categorical.

Illustration: NIR for Protein Content of Wheat. $p = 6$,
 $n = 50$. $\mathbf{f}_y^T = (y, y^2, y^3)$; $.9 < \widehat{\text{corr}}(X_i, X_k) < .999$. $\hat{d} = 1$.



Sufficient reduction



First PC

Where can we go from here?

1. PC: $X_y = \mu + \Gamma\nu_y + \sigma\varepsilon$

2. PFC: $X_y = \mu + \Gamma\beta f_y + \sigma\varepsilon$

Envelope models for increased efficiency

3. PFC: $X_y = \mu + \Gamma\beta f_y + \Delta^{1/2}\varepsilon$

Models for increased scope

$$X_y = \mu + \Gamma\beta f_y + \Delta_y^{1/2}\varepsilon$$

Increasing Scope

$$\mathbf{X}_y = \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\beta}f_y + \boldsymbol{\Delta}_y^{1/2}\boldsymbol{\varepsilon}$$

$R(\mathbf{X}) = \boldsymbol{\alpha}^T \mathbf{X}$ is a sufficient reduction iff

1. $\text{span}(\boldsymbol{\Delta}^{-1}\boldsymbol{\Gamma}) \subseteq \text{span}(\boldsymbol{\alpha})$, where $\boldsymbol{\Delta} = \mathbf{E}(\boldsymbol{\Delta}_Y)$.
2. $\boldsymbol{\alpha}_0^T \boldsymbol{\Delta}_y^{-1}$ is constant in y , where $(\boldsymbol{\alpha}, \boldsymbol{\alpha}_0)$ is orthogonal.

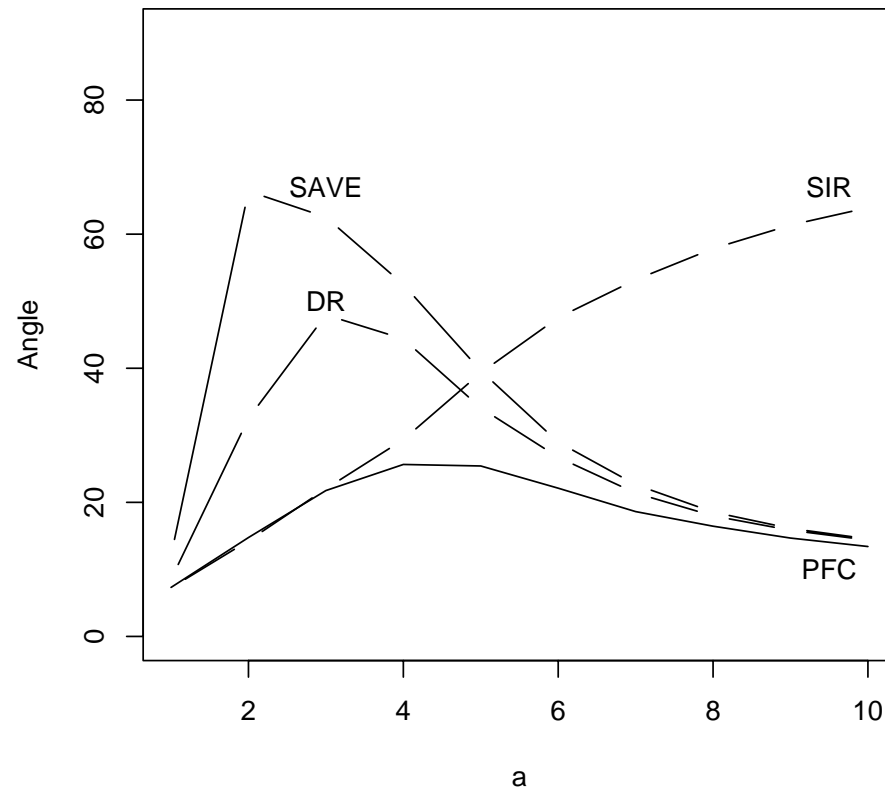
For y categorical, $\hat{R}(\mathbf{X}) = \hat{\boldsymbol{\alpha}}^T \mathbf{X}$, $\hat{\boldsymbol{\alpha}}$ is a basis for

$$\hat{\mathcal{S}}_{Y|\mathbf{X}} = \arg \max_{\mathcal{S} \in \mathcal{G}_{d,p}} n \log |P_{\mathcal{S}} \hat{\boldsymbol{\Sigma}} P_{\mathcal{S}}|_0 - \sum_y n_y \log |P_{\mathcal{S}} \hat{\boldsymbol{\Delta}}_y P_{\mathcal{S}}|_0$$

$|\mathbf{A}|_0 =$ product of non-zero eigenvalues of \mathbf{A} .

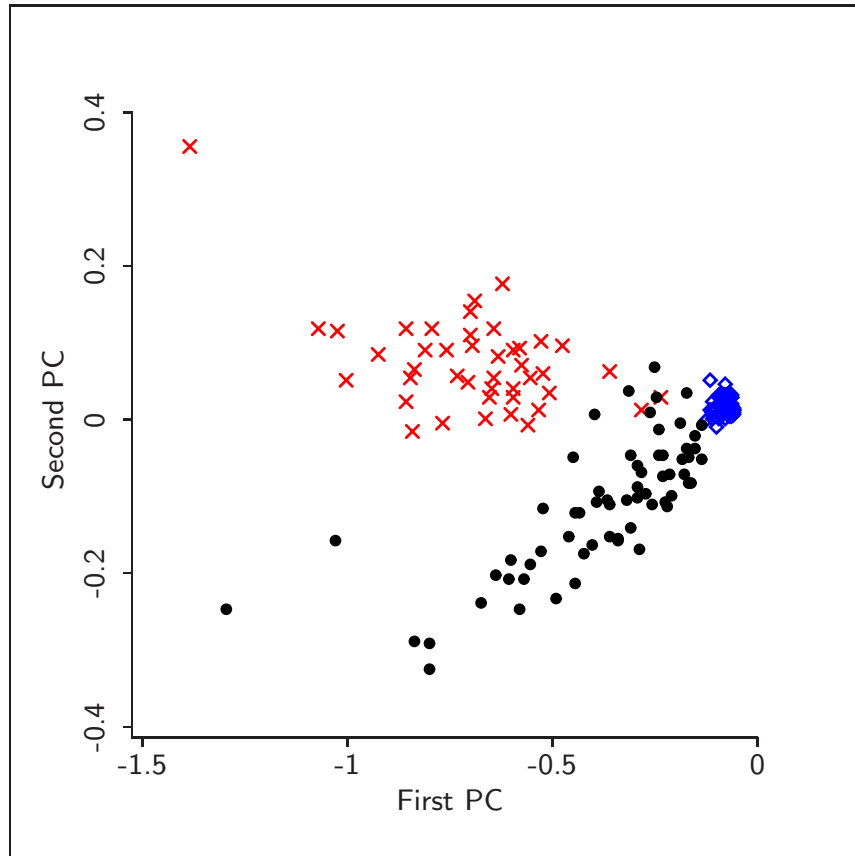
Simulation Model, $p = 8$

$$Y = X_1/a + aX_1^2/10 + \varepsilon, \mathcal{S}_{Y|\mathbf{X}} = \text{span}(1, 0 \dots, 0)^T$$

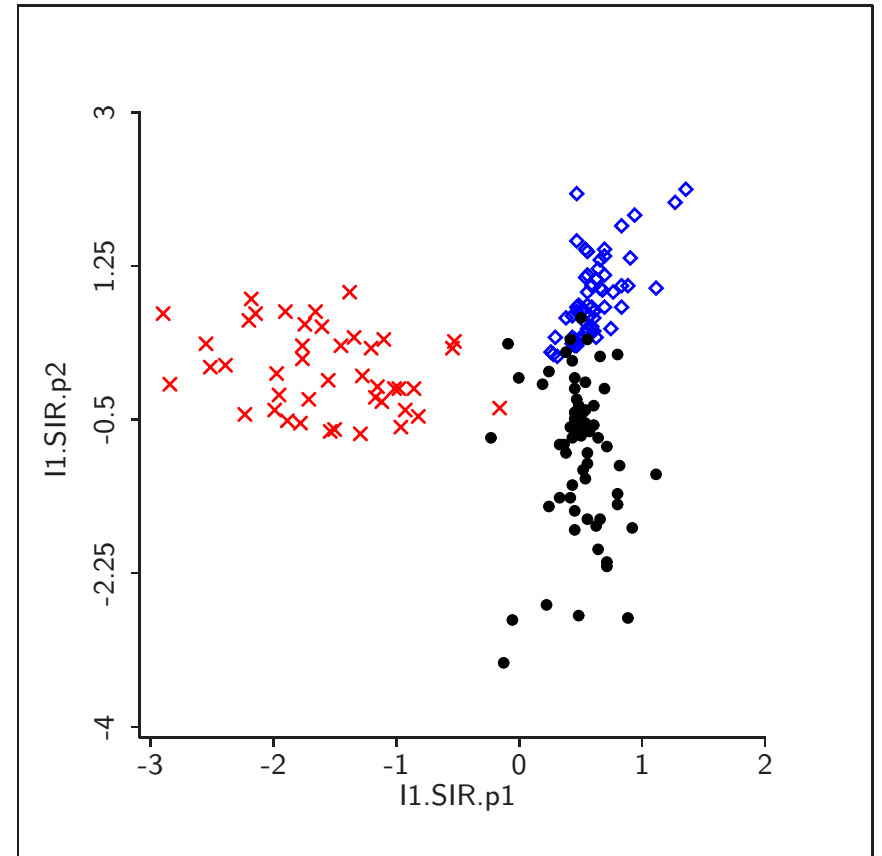


angle vs. a

Birds and Planes and Cars

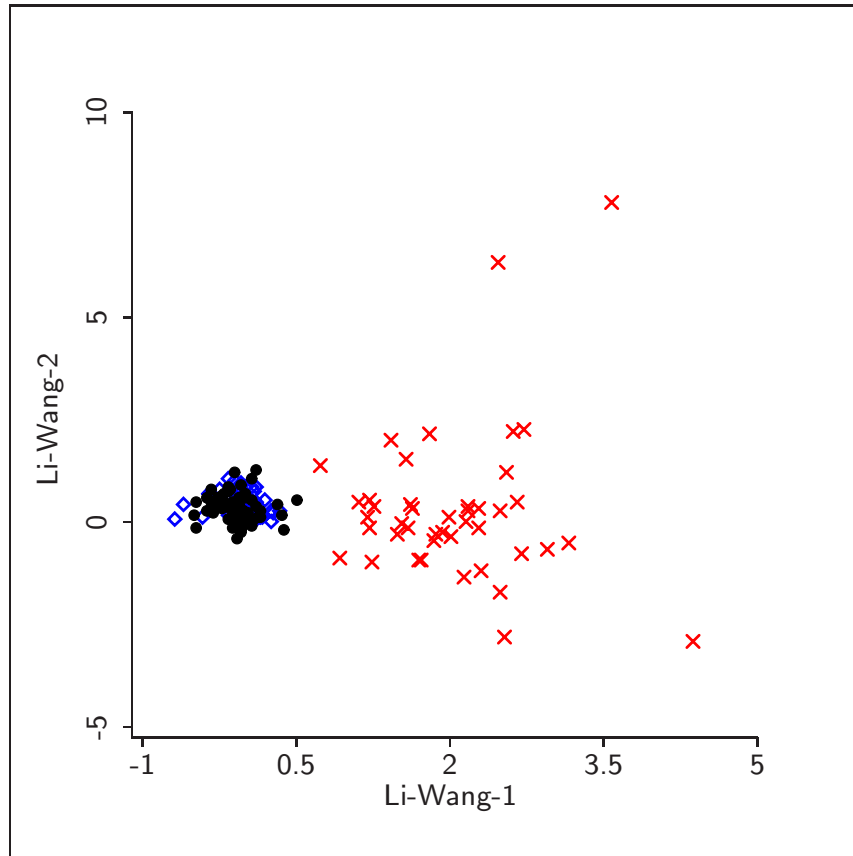


PC

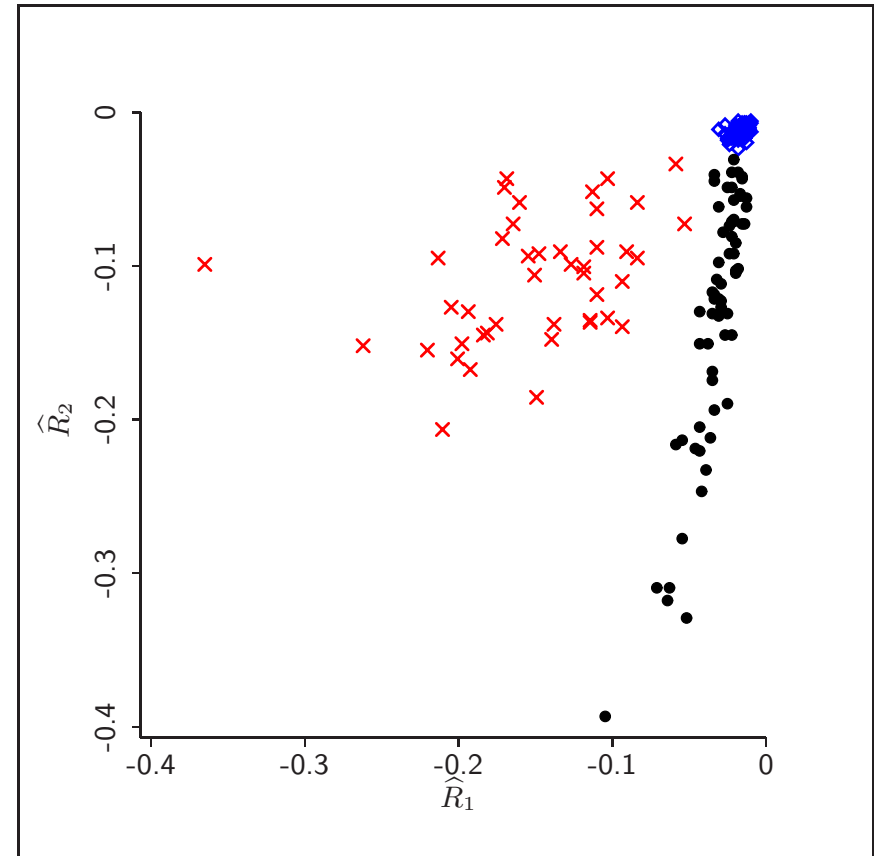


LDA

Birds and Planes and Cars



Li & Wang, JASA, 2007



Sufficient reductions: \hat{R}_1 & \hat{R}_2

Envelope Models for Increasing Efficiency

$$\begin{aligned}\mathbf{X}_y &= \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_y + \sigma\boldsymbol{\varepsilon} \\ \mathbf{X}_y &= \boldsymbol{\mu} + \boldsymbol{\Gamma}\boldsymbol{\beta}\mathbf{f}_y + \boldsymbol{\Delta}^{1/2}\boldsymbol{\varepsilon}\end{aligned}$$

$\mathcal{E}_\Delta(\mathcal{S}_\Gamma)$ = smallest reducing subspace of Δ that contains \mathcal{S}_Γ ,
the **Δ envelope of \mathcal{S}_Γ** . $u = \dim\{\mathcal{E}_\Delta(\mathcal{S}_\Gamma)\}$.

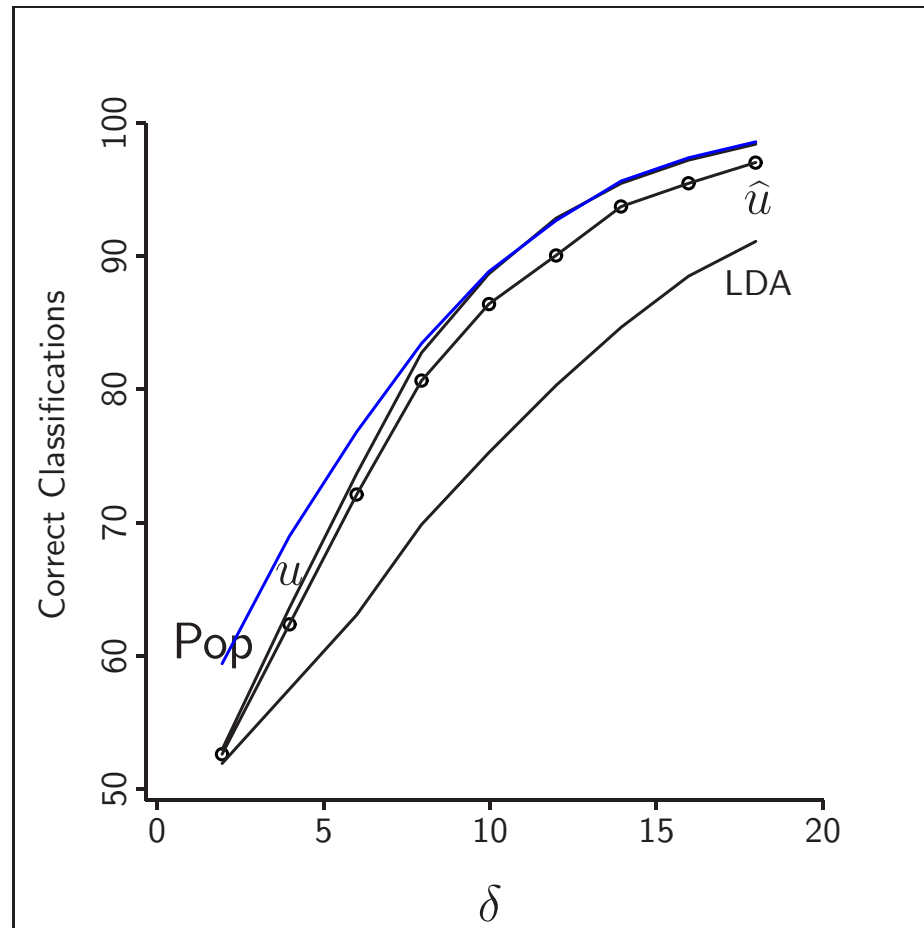
Let $\Phi : p \times u$ orthonormal basis for $\mathcal{E}_\Delta(\mathcal{S}_\Gamma)$. $d \leq u \leq p$.

$\Theta : u \times d$ semi-orthogonal matrix.

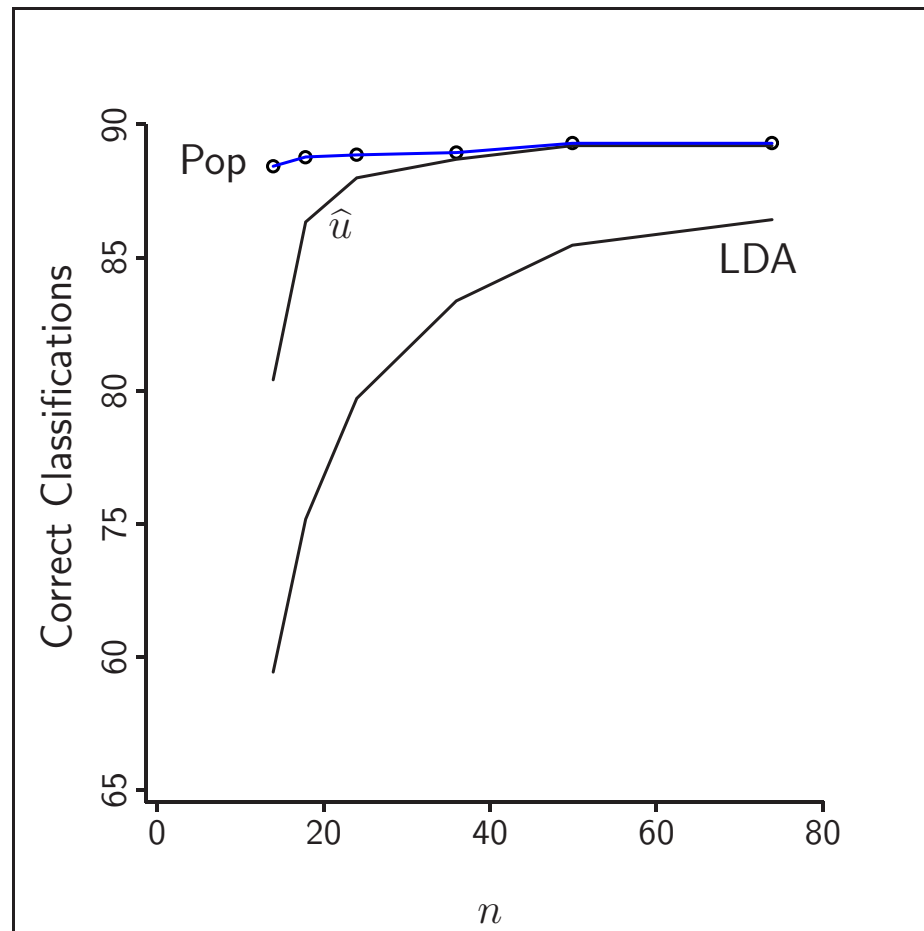
$$\begin{aligned}\mathbf{X}_y &= \boldsymbol{\mu} + \Phi\Theta\boldsymbol{\beta}\mathbf{f}_y + \boldsymbol{\Delta}^{1/2}\boldsymbol{\varepsilon} \\ \boldsymbol{\Delta} &= \Phi\mathbf{M}\Phi^T + \Phi_0\mathbf{M}_0\Phi_0^T \\ R(\mathbf{X}) &= \boldsymbol{\alpha}^T\mathbf{X}\end{aligned}$$

$\boldsymbol{\alpha}$ is a basis for $\mathcal{S}_{Y|\mathbf{X}} = \text{span}(\Phi\mathbf{M}^{-1}\Theta)$.

Discriminant analysis, two equally likely spherical normal populations, $p = 10$, $n = 18$, $\delta = \|\mu_1 - \mu_2\|$. 100 obs/pop for classification, 400 reps.



Same setting but varying n with $\delta = 10$.



Comments

- **Other distributions for X_y , like $X_y \sim QEF$.**
(Genotype data; SNPs)
- **Screening:** Regression of each elements of X_y on f_y ,
looking for any dependence.
- **Covariance reducing methods** with application in
evolutionary biology.
 - With y categorical, characterize $\{\Sigma_y\}$
 - Flury's spectral models
 - Instead, base methodology on finding $\alpha \in \mathbb{R}^{p \times d}$ so
that the distn of $\hat{\Sigma}_y | (\alpha^T \hat{\Sigma}_y \alpha, n_y)$ is constant in y .

- **Multivariate normal linear model**

$$\mathbf{Y}_i = \boldsymbol{\mu} + \boldsymbol{\beta}\mathbf{X}_i + \mathbf{V}^{1/2}N(0, I_r), \quad \mathbf{V} \geq 0, \quad i = 1, n$$

with inference on $\boldsymbol{\beta}$ and prediction as the goals.

- **Reparameterize in terms of $\mathcal{E}_{\mathbf{V}}(\text{span}(\boldsymbol{\beta}))$**

- **Use MLE and find the var in the limiting distn of $\hat{\boldsymbol{\beta}}$.**

- **Partial least squares**

- **$\mathbf{X}|Y$ vs $Y|\mathbf{X}$.**