

Absolute Penalty and B-spline-based Shrinkage Estimation in Partially Linear Models

Enayetur Raheem

University of Windsor / Windsor-Essex County Health Unit
Windsor, Ontario, Canada

International Workshop on High Dimensional Data Analysis
Fields Institute, June 9-11, 2011

With
K. A. Doksum (Wisconsin-Madison) and S Ejaz Ahmed (Windsor)

Outline

- Overview of shrinkage estimation
- Motivational Examples
- Asymptotic properties of shrinkage estimators
- Monte Carlo study
- Conclusion
- Software for shrinkage estimation

Brief introduction to shrinkage estimation

- Dates back to Stein (1956); then improved by James and Stein (1961).
- Suppose $\theta_{p \times 1}$ is an unknown parameter vector, and \mathbf{y} be a vector of observations of θ , s.t.

$$\mathbf{y} \sim N(\theta, \sigma^2 I)$$

We want an estimate $\hat{\theta}$ based on \mathbf{y}

- Since the noise has zero mean, we get using the concept of Least-squares

$$\hat{\theta}_{LS} = \bar{\mathbf{y}}$$

- Stein (1956) demonstrated that, in terms of a quadratic loss function, $E(\|\theta - \hat{\theta}\|^2)$, this approach is suboptimal. The result became known as **Stein's phenomenon**.

James-Stein Estimator

- If σ^2 is known, James-Stein showed that

$$\hat{\theta}_{JS} = \left\{ 1 - \frac{(p-2)\sigma^2}{\|\bar{\mathbf{y}}\|^2} \right\} \bar{\mathbf{y}}$$

dominates $\hat{\theta}_{LS}$ for $p \geq 3$

\Rightarrow JS estimator always achieves lower MSE than the LS estimator.

Shrinking towards ν

$$\hat{\theta}_{JS} = \left\{ 1 - \frac{(p-2)\sigma^2}{\|\bar{\mathbf{y}}\|^2} \right\} \bar{\mathbf{y}}$$

- Notice that if $(p-2)\sigma^2 < \|\bar{\mathbf{y}}\|^2$, this estimator shrinks the natural estimator $\bar{\mathbf{y}}$ towards zero.
- For any fixed vector ν of length p , there exists a JS estimator that shrinks towards ν , which is

$$\hat{\theta} = \left\{ 1 - \frac{(p-2)\sigma^2}{\|(\bar{\mathbf{y}} - \nu)\|^2} \right\} (\bar{\mathbf{y}} - \nu) + \nu$$

for any ν

Optimal choice of ν

- A natural question is whether the improvement over the usual estimator is independent of the choice of ν
- Improvement is large if $\|\theta - \nu\|$ is small

Usefulness

- Shrinkage estimators are more efficient than the classical estimators for $p \geq 3$ (in terms of a quadratic risk function, e.g., MSE)
- Ensures improved prediction accuracy

Motivating Examples

[Mroz (Econometrica, 1987)] used a sample of 1975 PSID¹ labour supply data to systematically study several theoretic and statistical assumptions used in many empirical models of female labour supply.

This data set was used in numerous Econometric studies and books.

- [Wooldridge, 2003] used this data in his book to demonstrate many applications of regression models.
- [Long, 1997] fitted parametric logistic regression model to this data.
- [Fox, 2002] used this data for a semiparametric logistic regression.

¹ Panel Study on Income Dynamics (PSID), University of Michigan. The data is in public domain and freely available from <http://ideas.repec.org/s/boc/bocins.html>

Motivating Examples ...

- [Fox, 2005] commented that semiparametric model may be used wherever there is reason to believe that one or more covariates enter the regression linearly. This could be known from prior studies, or there are prior reasons to believe so (although rare), or examination of the data might suggest a linear relationship for some covariates. A more general scenario is when some of the covariates are categorical and they enter in the model as dummy variables.
- [Engle (JASA, 1986)] considered such a PLM where demand for electricity was the outcome variable and four regions were entered as dummy variable while temperature was modelled nonparametrically.

Objective

- Motivated by [Engle (JASA, 1986)] and [Fox, 2002], we fit a semiparametric regression model with a continuous outcome variable and demonstrate shrinkage estimation using PSID data.
- We first show that semiparametric modeling is appropriate in our case.

Partially Linear Model

- A partially linear regression model:

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + g(t_i) + \varepsilon_i, i = 1, \dots, n, \quad (1)$$

where y_i 's are responses,

$\mathbf{x}_i = (x_{i1}, \dots, x_{ip})'$ and $t_i \in [0, 1]$ are design points,

$\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ is an unknown parameter vector,

$g(\cdot)$ is an unknown real-valued function defined on $[0, 1]$, and ε_i 's are unobservable random errors.

Related earlier work

- Extension of Ahmed et al. (2007)

Related earlier work

- Extension of Ahmed et al. (2007)

Objective

- We consider estimating β based on $g(\cdot)$ approximated by a B-spline series.

Data and Variables

- The female labour supply data consist of 753 observations on 19 variables. Data were collected from married white women between the ages 30 and 60 in 1975. Of them, 428 were working at some time during the year 1975.
- Similar to [Mroz (Econometrica, 1987)], we consider wife's annual hours of work (`hours`) as the response variable.
- We only used the portion of the data when the women were in labour force. Thus, we had 428 cases (rows) in our working data.

Table: Description of Variables in the Model for Working Women.

Covariates	Description	Remarks
hours	Hrs worked in 1975	Min=12, max=4950, med=1303
age	Age of woman	Min=30, max=60, med=42
nwifeinc	Non-wife income	Income in thousands
k5	Kids five and under	0-1, a few 2's and 3's, factor
k618	Kids six to 18	0-4, few >4, factor variable
wc	W attended college?	1 (if educ > 12), else 0
hc	H attended college?	1 (if huseduc > 12), else 0
unem	Unemployment rate	Min=3, max=14, median=7.5
mtr	Marginal tax rate	Min=0.44, max=0.94, med=0.69
exper	Labour market exp	Min=0, max=38, median=12

Algorithm

- When we have prior information about certain covariates: shrinkage estimators are directly obtain by combining the full and sub-model estimates.

if *a priori* information is not available, it take a two step approach

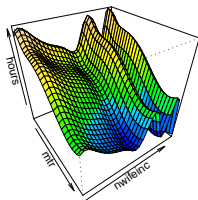
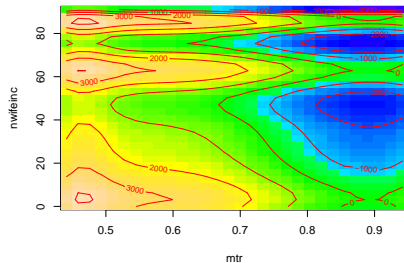
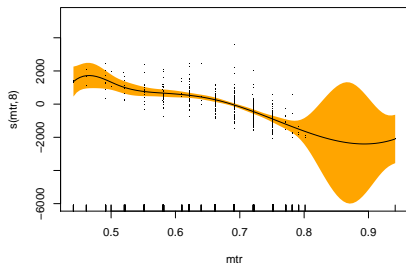
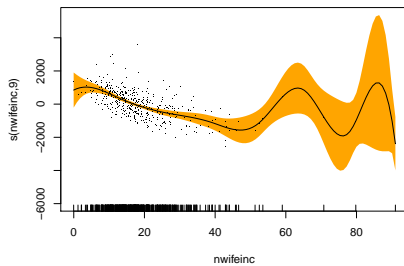
Step 1 a set of covariates are selected based on model selection criterion (AIC, BIC etc.) that form the main covariate-set. Consequently, the remaining covariates are taken as nuisance subset—forming the restriction on the full model.

Step 2 Full- and the sub-model estimates are combined in a way that minimizes the quadratic risk function.

Table: Selection of covariates by AIC, BIC and lasso.

Models	Selected Variables
Full Model	wc, nwifeinc, mtr, exper, unem, k5, age, k618, hc
AIC	wc, nwifeinc, mtr, exper, unem, k5, age
BIC	wc, nwifeinc, mtr, exper
Lasso	wc, nwifeinc, mtr, exper, unem, k5, age, hc

- Our analyses show that `nwifeinc` has significant nonlinear relationship ($p = 0.0065$) with hours of work.

(a)**(b)****(c)****(d)**

Full- and Sub-model

Finally, with the inclusion of a nonparametric part, our candidate full- and sub-models are listed below.

Full-Model:

$$\begin{aligned} \text{hours} = & \text{wc} + g(\text{nwifeinc}) + \text{mtr} + \text{exper} + \text{unem} \\ & + \text{k5} + \text{age} + \text{k618} + \text{hc} \end{aligned}$$

Sub-Model:

$$\text{hours} = \text{wc} + g(\text{nwifeinc}) + \text{mtr} + \text{exper}$$

Here $g(\cdot)$ denotes a covariate estimated by B-spline basis function.

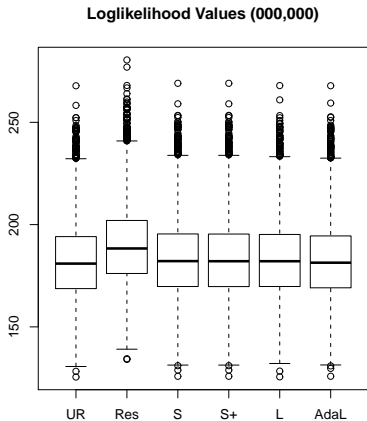
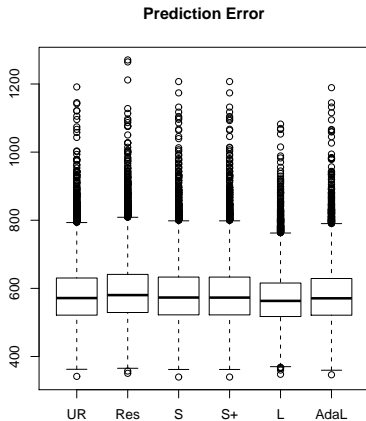


Figure: Comparison of the estimators through prediction errors and loglikelihood values.

THEORETICAL RESULTS:Parameter Estimation

Estimation of $g(\cdot)$ by B-spline

- Let k be an integer larger than or equal to ν where ν will be defined in the Assumption (2)
- Further, let $S_{m_n,k}$ be the class of functions $s(\cdot)$ on $[0, 1]$ with the following properties:
 - (i) $s(\cdot)$ is a polynomial of degree k on each of the sub-intervals $\left[\frac{(i-1)}{m_n}, \frac{i}{m_n}\right]$, $i = 1, \dots, m_n$, where m_n is a positive integer which depends on n .
 - (ii) $s(\cdot)$ is $(k - 1)$ times differentiable.
- Then, $S_{m_n,k}$ is called the class of all splines of degree k with m_n -equispaced knots.

Parameter Estimation...

Consequently,

- $S_{m_n, k}$ has a basis of $m_n + k$ normalized B-spline $\{B_{m_n j}(\cdot) : j = 1, \dots, m_n + k\}$, and
- $g(\cdot)$ can be approximated by a linear combination $\boldsymbol{\theta}' \mathbf{B}_{m_n}(\cdot)$ of the basis, where $\boldsymbol{\theta} \in \mathcal{R}^{m_n + k}$ and $\mathbf{B}_{m_n}(\cdot) = (B_{m_n 1}(\cdot), \dots, B_{m_n, m_n + k}(\cdot))'$. See de Boor (2001)
- Now, replacing $g(\cdot)$ by $\mathbf{B}_{m_n}(\cdot)\boldsymbol{\theta}$, in model (1) we get

$$\boxed{y = \mathbf{x}\boldsymbol{\beta} + \mathbf{B}_{m_n}(t)\boldsymbol{\theta} + \varepsilon} \quad (2)$$

Full Model Estimation

- $(\hat{\beta}, \hat{\theta})$ is obtained by minimizing

$$S_n(\beta, \theta) = n^{-1} \sum_{i=1}^n [y_i - \mathbf{x}_i' \beta - \theta' \mathbf{B}_{m_n}(t_i)]^2, \quad (3)$$

which gives

$$\hat{\beta} = (\mathbf{X}' \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X})^{-1} \mathbf{X}' \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{Y}$$

and

$$\hat{\theta} = (\mathbf{B}_{m_n}' \mathbf{B}_{m_n})^{-1} \mathbf{B}_{m_n}' (\mathbf{Y} - \mathbf{X} \hat{\beta}),$$

where, $\mathbf{Y} = (y_1, \dots, y_n)'$, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$,

$\mathbf{x}_s = (x_{1s}, \dots, x_{ns})'$, $s = 1, \dots, p$,

$\mathbf{M}_{\mathbf{B}_{m_n}} = \mathbf{I} - \mathbf{B}_{m_n} (\mathbf{B}_{m_n}' \mathbf{B}_{m_n})^{-1} \mathbf{B}_{m_n}'$ and

$\mathbf{B}_{m_n} = (B_{m_n}(t_1), \dots, B_{m_n}(t_n))'$.

- The estimator $\hat{\beta}$ is called a semiparametric least squares estimator (SLSE) of β .

Uncertain Prior Information

- β in the linear part can be partitioned as (β_1, β_2)
- β_1 is the coefficient vector for main effects (e.g., treatment effect, genetic effects) and β_2 is a vector for “nuisance” effects (e.g., age, laboratory).

Sub Model Estimation

Semi-Parametric Unrestricted Estimator (SURE)

$$\hat{\beta}_1^{UR} = (\mathbf{X}'_1 \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_2 \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_2 \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{Y},$$

where

\mathbf{X}_1 is composed of the first p_1 column vectors of \mathbf{X} ,

\mathbf{X}_2 is composed of the last p_2 column vectors of \mathbf{X} , and

$$\mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_2 = \mathbf{I} - \mathbf{B}_{m_n} \mathbf{X}_2 (\mathbf{X}'_2 \mathbf{B}'_{m_n} \mathbf{B}_{m_n} \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{B}'_{m_n}.$$

Sub Model Estimation

Semi-Parametric Unrestricted Estimator (SURE)

$$\hat{\beta}_1^{UR} = (\mathbf{X}_1' \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_2 \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_2 \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{Y},$$

where

\mathbf{X}_1 is composed of the first p_1 column vectors of \mathbf{X} ,

\mathbf{X}_2 is composed of the last p_2 column vectors of \mathbf{X} , and

$$\mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_2 = \mathbf{I} - \mathbf{B}_{m_n} \mathbf{X}_2 (\mathbf{X}_2' \mathbf{B}_{m_n}' \mathbf{B}_{m_n} \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{B}_{m_n}'.$$

Semi-Parametric Restricted Estimator (SRE)

$$\hat{\beta}_1^R = (\mathbf{X}_1' \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_{\mathbf{B}_{m_n}} \mathbf{Y}.$$

Shrinkage Estimation

Semiparametric Stein-Type Estimation

A semiparametric Stein-type estimator (SSTE) $\hat{\beta}_1^S$ of β_1 can be defined as

$$\hat{\beta}_1^S = \hat{\beta}_1^R + (\hat{\beta}_1^{UR} - \hat{\beta}_1^R) \{1 - (p_2 - 2)\psi_n^{-1}\}, \quad p_2 \geq 3.$$

where

$$\psi_n = \frac{n}{\hat{\sigma}_n^2} \hat{\beta}_2' \mathbf{X}_2' \mathbf{B}_{m_n}' \mathbf{M}_{\mathbf{B}_{m_n} \mathbf{X}_2} \mathbf{B}_{m_n} \mathbf{X}_2 \hat{\beta}_2,$$

with

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i' \hat{\beta} - \mathbf{B}_{m_n}'(t_i) \hat{\theta})^2.$$

Positive-part Semiparametric Stein-Type Estimator

A PSSTE has the form

$$\hat{\beta}_1^{S+} = \hat{\beta}_1^R + (\hat{\beta}_1^{UR} - \hat{\beta}_1^R) \{1 - (p_2 - 2)\psi_n^{-1}\}^+, \quad p_2 \geq 3.$$

where $z = \max(0, z)$.

Absolute Penalty Estimator

The lasso estimates are obtained as

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}.$$

Tibshirani (1996).

Asymptotics

Assumption 1

There exist bounded functions $h_s(\cdot)$ over $[0, 1]$, $s = 1, \dots, p$, such that

$$x_{is} = h_s(t_i) + u_{is}, \quad i = 1, \dots, n, s = 1, \dots, p, \quad (a)$$

where $u_i = (u_{i1}, \dots, u_{ip})'$ are real vectors satisfying

$$\lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n u_{ik} u_{ij}}{n} = b_{kj}, \quad \text{for } k = 1, \dots, p, j = 1, \dots, p, \quad (b)$$

and the matrix $\mathbf{B} = (b_{kj})$ is nonsingular. Moreover,

$$\max_{1 \leq k \leq p} \|\mathbf{A} \mathbf{u}_k^*\| = O\left(\left[\text{tr}(\mathbf{A}'\mathbf{A})\right]^{\frac{1}{2}}\right), \quad \text{for any matrix } \mathbf{A}, \quad (c)$$

where $\mathbf{u}_k^* = (u_{1k}, \dots, u_{nk})'$ and $\|\cdot\|$ denotes the Euclidean norm.

Asymptotics

Assumption 2

The functions $g(\cdot)$ and $h_j(\cdot)$ satisfy the Lipschitz condition of order ν , i.e., there exists a constant c such that

$$|f_j(s) - f_j(t)| \leq c|s - t|^{(\nu)}, \quad \text{for any } s, t \in [0, 1], \quad j = 0, 1, \dots, p,$$

where $f_0(\cdot) = g(\cdot)$ and $f_j(\cdot) = h_j(\cdot)$, $j = 1, \dots, p$.

Asymptotics

Lemma

If Assumptions 1 and 2 are satisfied, and ε_i are independent with mean zero and constant variance σ^2 and $\mu_{3i} = E\varepsilon_i^3$ being uniformly bounded, then

$$\sqrt{n}(\hat{\beta} - \beta) \rightarrow_D N(0, \sigma^2 \mathbf{B}^{-1}) \quad \text{and} \quad \mathbf{B}'_{m_n}(t)\hat{\theta} - g(t) = O_p(n^{-\frac{\nu}{2\nu+1}}),$$

where “ \rightarrow_D ” denotes convergence in distribution and \mathbf{B} is defined in Assumption 1. Proof is similar to Ahmed et al (2007).

Asymptotic Properties

- To study the asymptotic quadratic risks of $\hat{\beta}_1^{UR}$, $\hat{\beta}_1^R$, $\hat{\beta}_1^S$ and $\hat{\beta}_1^{S+}$, we consider a sequence of local alternatives
- Under fixed alternatives all the estimators are asymptotically equivalent to $\hat{\beta}_1$, while $\hat{\beta}_1^R$ has unbounded risk.
- A sequence $\{K_n\}$ of local alternatives defined by

$$K_n : \beta_{2(n)} = n^{-\frac{1}{2}}\omega, \omega \neq \mathbf{0} \text{ fixed} \quad (4)$$

Asymptotic Properties

Asymptotic Distributional Bias (ADB)

The asymptotic distributional bias (ADB) of an estimator δ is defined as

$$\text{ADB}(\delta) = \lim_{n \rightarrow \infty} E \left\{ n^{\frac{1}{2}} (\delta - \beta_1) \right\}.$$

Asymptotic Distirbutional Bias (ADB)

Suppose that Assumptions 1 and 2 hold. Under $\{K_n\}$, the ADB of the estimators are as follows:

$$\text{ADB}(\hat{\beta}_1^{UR}) = \mathbf{0},$$

$$\text{ADB}(\beta_1^R) = -\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\omega,$$

$$\text{ADB}(\hat{\beta}_1^S) = -(p_2 - 2)\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\omega E(\chi_{p_2, \alpha}^{-2}; \Delta),$$

$$\begin{aligned} \text{ADB}(\hat{\beta}_1^{S+}) = & -\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\omega H_{p_2+2}(p_2 - 2; \Delta) - (p_2 - 2)\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\omega \\ & \left\{ E \left[\chi_{p_2+2}^{-2}(\Delta) I(\chi_{p_2+2}^2(\Delta) > p_2 - 2) \right] \right\} \end{aligned}$$

where $\mathbf{B} = \begin{pmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{pmatrix}$ with \mathbf{B} defined in Assumption 1,

$\Delta = (\omega' \mathbf{B}_{22.1} \omega) \sigma^{-2}$, $\mathbf{B}_{22.1} = \mathbf{B}_{22} - \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{B}_{12}$, $H_\nu(x; \Delta)$ denotes the noncentral chi-square distribution function with noncentrality parameter Δ and ν degrees of freedom.

RISK PROPERTIES

Asymptotic Distirbutional Risk (ADR)

- Define a quadratic loss function using a positive definite matrix (p.d.m.) \mathbf{M} , namely,

$$\mathcal{L}(\delta, \beta_1) = n(\delta - \beta_1)' \mathbf{M}(\delta - \beta_1),$$

δ can be any one of estimator of $\hat{\beta}_1$,

- We assume that the asymptotic distribution function of δ under $\{K_n\}$ exists and is given by

$$F(\mathbf{x}) = \lim_{n \rightarrow \infty} P\{\sqrt{n}(\delta - \beta_1) \leq \mathbf{x} | K_n\}$$

- Then, the AQDR of δ is defined as

$$R(\delta, \mathbf{M}) = \text{tr} \left\{ \mathbf{M} \int_{\mathcal{R}^{p_1}} \int \mathbf{x} \mathbf{x}' dF(\mathbf{x}) \right\} = \text{tr}(\mathbf{M} \mathbf{V}),$$

where \mathbf{V} is the dispersion matrix for the asymptotic distribution $F(\mathbf{x})$.

RISK PROPERTIES

Asymptotic Quadratic Distributional Risk (AQDR)

Suppose that assumptions 1 and 2 hold, then under $\{K_n\}$ the AQDR of the estimators are:

$$\begin{aligned}R(\hat{\beta}_1^{UR}; \mathbf{M}) &= \sigma^2 \text{tr}(\mathbf{M}\mathbf{B}_{11.2}^{-1}), \\R(\beta_1^R; \mathbf{M}) &= \sigma^2 \text{tr}(\mathbf{M}\mathbf{B}_{11}^{-1}) + \omega' \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{M}\mathbf{B}_{11}^{-1} \mathbf{B}_{12} \omega, \\R(\hat{\beta}_1^S; \mathbf{M}) &= \sigma^2 \left[\text{tr}(\mathbf{M}\mathbf{B}_{11.2}^{-1}) - (p_2 - 2) \text{tr}(\mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{M}\mathbf{B}_{11}^{-1} \mathbf{B}_{12} \mathbf{B}_{22.1}^{-1}) \right. \\&\quad \cdot \left. \left\{ 2E(\chi_{p_2, \alpha}^2(\Delta)) - (p_2 - 2)E(\chi_{p_2+2}^{-4}(\Delta)) \right\} \right] \\&\quad + (p_2^2 - 4) \omega' \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{M}\mathbf{B}_{11}^{-1} \mathbf{B}_{12} \omega E(\chi_{p_2+4}^{-4}(\Delta)),\end{aligned}$$

RISK PROPERTIES

$$\begin{aligned} & R(\hat{\beta}^{S+}; \mathbf{M}) \\ = & R(\hat{\beta}^S; \mathbf{M}) + (p_2 - 2)\text{tr}(\mathbf{B}_{21}\mathbf{B}_{11}^{-1}\mathbf{M}\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\mathbf{B}_{22.1}^{-1}) \left\{ 2E \left[\chi_{p_2+2}^{-2}(\Delta) I(\chi_{p_2+2}^2(\Delta) \leq p_2 - 2) \right] \right. \\ & \left. - (p_2 - 2)E \left[\chi_{p_2+2}^{-4}(\Delta) I(\chi_{p_2+2}^2(\Delta) \leq p_2 - 2) \right] \right\} - \sigma^2\text{tr}(\mathbf{B}_{21}\mathbf{B}_{11}^{-1}\mathbf{M}\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\mathbf{B}_{22.1}^{-1}) \\ & \cdot H_{p_2+2}(p_2 - 2; \Delta) + \omega' \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{M} \mathbf{B}_{11}^{-1} \mathbf{B}_{12} \omega [2H_{p_2+2}(p_2 - 2; \Delta) - H_{p_2+4}(p_2 - 2; \Delta)] \\ - & (p_2 - 2)\omega \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{M} \mathbf{B}_{11}^{-1} \mathbf{B}_{12} \omega' \left\{ 2E \left[\chi_{p_2+2}^{-2}(\Delta) I(\chi_{p_2+2}^2(\Delta) \leq p_2 - 2) \right] \right. \\ - & \left. 2E \left[\chi_{p_2+2}^{-4}(\Delta) I(\chi_{p_2+2}^2(\Delta) \leq p_2 - 2) \right] + (p_2 - 2)E \left[\chi_{p_2+2}^{-4}(\Delta) I(\chi_{p_2+2}^2(\Delta) \leq p_2 - 2) \right] \right\}, \end{aligned}$$

THEORETICAL SUMMARY

Dominance of Shrinkage Estimator

By comparing $R(\hat{\beta}_1^S)$ and $R(\hat{\beta}_1^R)$ following dominance condition holds. If $\mathbf{M} \in \mathbf{M}^D$, $\hat{\beta}_1^S$ dominates $\hat{\beta}_1$ for any ω in the sense of AQDR, where

$$\mathbf{M}^D = \left\{ \mathbf{M} : \frac{\text{tr}(\mathbf{B}_{21}\mathbf{B}_{11}^{-1}\mathbf{M}\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\mathbf{B}_{22.1}^{-1})}{\text{ch}_{\max}(\mathbf{B}_{21}\mathbf{B}_{11}^{-1}\mathbf{M}\mathbf{B}_{11}^{-1}\mathbf{B}_{12}\mathbf{B}_{22.1}^{-1})} \geq \frac{p_2 + 2}{2} \right\}.$$

THEORETICAL SUMMARY

Dominance of Shrinkage Estimator

- $\hat{\beta}^{S+}$ dominates $\hat{\beta}^S$ for all the values of ω , with strict inequality holds for some ω .
- Risk of $\hat{\beta}^{S+}$ is also smaller than the risk of $\hat{\beta}_1$ in the entire parameter space and the upper limit is attained when Δ approaches ∞ .
- This implies that

$$R(\hat{\beta}_1^{S+}) \leq R(\hat{\beta}_1^S) \leq R(\hat{\beta}_1), \text{ for any } \mathbf{M} \in \mathbf{M}^D \text{ and } \omega,$$

with strict inequality holds for some ω .

Thus, we conclude that $\hat{\beta}_1^S$ and $\hat{\beta}_1^{S+}$ perform better than $\hat{\beta}_1$ in the entire parameter space induced by Δ . The gain in risk over $\hat{\beta}_1$ is substantial when $\Delta = 0$ or near.

MONTE CARLO RESULTS

Setup

We simulate the response from the following model:

$$y_i = x_{1i}\beta_1 + x_{2i}\beta_2 + \dots + x_{pi}\beta_p + g(t_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

where $t_i = (i - 0.5)/n$, $x_{1i} = (\zeta_{1i}^{(1)})^2 + \zeta_i^{(1)} + \xi_{1i}$,
 $x_{2i} = (\zeta_{2i}^{(1)})^2 + \zeta_i^{(1)} + 2\xi_{2i}$, $x_{si} = (\zeta_{si}^{(1)})^2 + \zeta_i^{(1)}$ with
 $\zeta_{si}^{(1)}$ i.i.d. $\sim N(0, 1)$, $\zeta_i^{(1)}$ i.i.d. $\sim N(0, 1)$,
 $\xi_{1i} \sim \text{Bernoulli}(0.45)$ and $\xi_{2i} \sim \text{Bernoulli}(0.45)$ for all
 $s = 3, \dots, p$ and $i = 1, \dots, n$.
Moreover, ε_i are i.i.d. $N(0, 1)$, and $g(t) = \sin(4\pi t)$.

MONTE CARLO RESULTS

Hypotheses that we are testing

- $H_0 : \beta_j = \mathbf{0}$, for $j = p_1 + 1, p_1 + 2, \dots, p_1 + p_2$, with $p = p_1 + p_2$.
- We are interested to estimating $\beta_1, \beta_2, \beta_3$ and β_4 when the remaining regression parameters may not be useful.
- We partition the regression coefficients as $\beta = (\beta_1, \beta_2) = (\beta_1, \mathbf{0})$ with $\beta_1 = (2, 1.5, 1, 0.6)$,
- We defined the parameter $\Delta^* = \|\beta - \beta^{(0)}\|$, where $\beta^{(0)} = (\beta_1, \mathbf{0})$ and $\|\cdot\|$ is the Euclidean norm.
- To determine the behavior of the estimators for $\Delta^* > 0$, further datasets were generated from those distributions under local alternative hypotheses. We considered $\Delta^* = 0, .1, .2, .3, .4, .5, .8, 1, 2$, and 4.

MONTE CARLO RESULTS

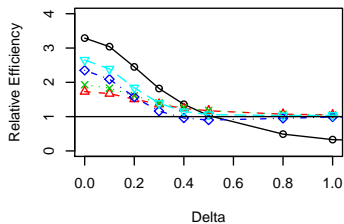
Relative MSE

Numerically calculated the relative MSEs of the proposed estimators $\tilde{\beta}_1$, $\hat{\beta}_1^S$, $\hat{\beta}_1^{S+}$, $\hat{\beta}_1^{PT}$, $\hat{\beta}_1^{IPT}$, relative to the unrestricted estimator $\hat{\beta}_1$, using:

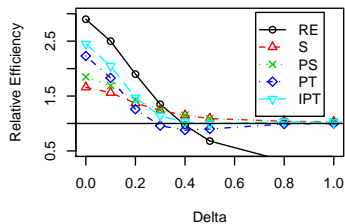
$$RMSE(\hat{\beta}_1 : \hat{\beta}_1^\diamond) = \frac{MSE(\hat{\beta}_1)}{MSE(\hat{\beta}_1^\diamond)}.$$

The amount by which an RMSE is larger than unity indicates the degree of superiority of the estimator $\hat{\beta}_1^\diamond$ over $\hat{\beta}_1$.

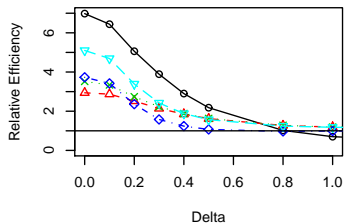
n=30, p1=4, p2=5



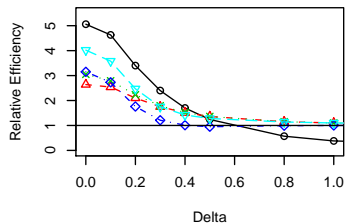
n=40, p1=4, p2=5



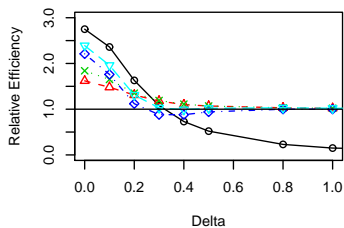
n=30, p1=4, p2=9



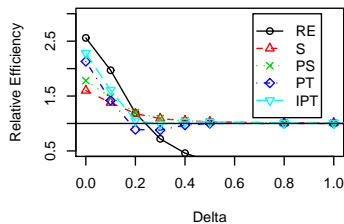
n=40, p1=4, p2=9



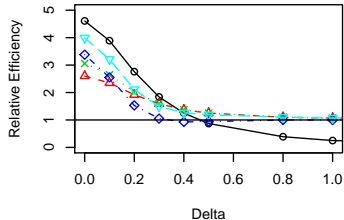
$n=50, p_1=4, p_2=5$



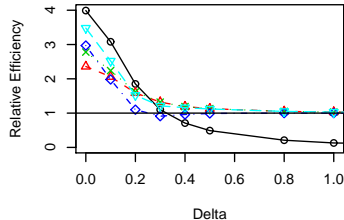
$n=80, p_1=4, p_2=5$



$n=50, p_1=4, p_2=9$



$n=80, p_1=4, p_2=9$



Shrinkage Vs. APE

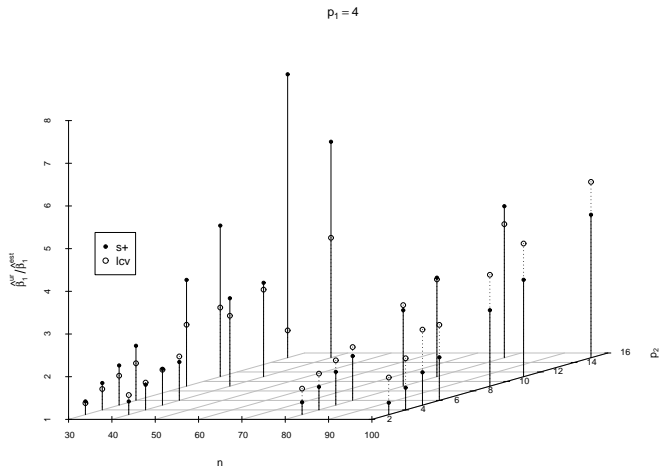


Figure: Three-dimensional plot of simulated relative MSE against n and p_2 to compare a positive shrinkage estimator and APE (CV) when $p_1 = 4$.

Summary

- Application to a real data shows the proposed method performs equally, in terms of loglikelihood values, with the absolute penalty estimator (our model has smaller number covariates)
- In Monte Carlo simulation, we found the positive-shrinkage estimator have smaller relative risk when the number of nuisance parameters are large.
- APE performs the best when p_1 and n are large

Future/Current work

- Shrinkage and APE estimation in multiple linear regression models
- Software for shrinkage estimation

shrink: R-package for Shrinkage Estimation

- Provides James-Stein-type shrinkage estimates of the regression parameters in a linear regression model.
- Computes unrestricted (ur), restricted (res), shrinkage (s), and positive shrinkage (ps) estimates of slope parameters.

Technically:

- `shrink()` produces an object of class “shrink” for which `coef()`, `predict()`, `residuals()` and `plot()` methods are available.
- Based on S3 methods of writing R extensions.

Usage: Options & Arguments

```
shrink(formula, formula2 = NULL, data = list(),  
method = "AIC", est = "ps", test = "F", ...)
```

<code>formula</code>	A full model formula in the form of $y \sim x_1 + x_2$;
<code>formula2</code>	Optional sub-model in the form of $y \sim x_1$; automatically selected if not supplied
<code>data</code>	Optional <code>data.frame</code> , or a matrix whose first column is the response
<code>method</code>	Sub-model selection method, either "AIC" or "BIC"
<code>est</code>	Estimators: one of "ur", "res", "s", "ps". Default is "ps"
<code>test</code>	Test statistic; either "Chisq" or "F"; default is "F"
<code>...</code>	Other arguments that can be passed to <code>shrink</code>

Example :

```
shrink(bodyfat ~., data=bodyfat, est='ps',  
method='BIC', test='Chisq')
```

Data source: <http://lib.stat.cmu.edu/datasets/bodyfat>

Thank you!

References



Ahmed, S.E., Doksum, K. A., Hossain, S. and You, J. (2007)
Shrinkage, Pretest and Absolute Penalty Estimators in Partially linear models
Australian & New Zealand Journal of Statistics, 49, 435-454.



Engle, R.F. and W.J. Granger and J. Rice and A. Weiss (1986)
Semiparametric estimates of the relation between weather and electricity sales
Journal of the American Statistical Association, 80, 310 – 319



Fox, John (2002)
An R and S-PLUS Companion to Applied Regression
Sage Publications, Thousand Oaks



Fox, John (2005)
Introduction to Nonparametric Regression
<http://socserv.socsci.mcmaster.ca/jfox/Courses/Oxford-2005/index.html>

References



Long, J. S. (1997)

Regression models for categorical and limited dependent variables,
SAGE Publications, Thousand Oaks.



Mroz, T. A. (1987)

The sensitivity of an empirical model of married womens hours of work to
economic and statistical assumptions,
Econometrica 55(4), 765799.



Wooldridge, J. (2003)

Introductory econometrics: a modern approach
South Western College Publishing.



Tibshirani, R. (1996).

Regression shrinkage and selection via the lasso.
J. Royal. Statist. Soc B., Vol. 58, No. 1, pages 267-288.

References



A.K.Md. E. Saleh (2006).

Theore of Preliminary Testand Stein-Type Estimation with Applications
Wiley.



Ahmed, S. E. , Saleh, E., Volodin , A. I. and I. N. Volodin (2007).

Asymptotic expansion of the coverage probability of James-Stein
Estimators.

Journal of Theory of Probability and its Applications, 51, 683-695.