# $L^\gamma$ Penalty Models

# Computation And Applications

# Part I

## Wenjiang Fu

**wfu@stat.tamu.edu**

**http://stat.tamu.edu/∼wfu.**

**Department of Statistics, Texas A&M University**

# OUTLINE

**Motivation**

**Lecture 1.** $L^{\gamma}$ **penalty: variable selection and computation for linear models**

**Lecture 2. Selection of tuning parameter and asymptotics**

**Lecture 3. Extension to non-Gaussian response and longitudinal studies**

**Lecture 4. Recent development in $L^{\gamma}$ penalty models and related topics**

# Linear regression model

$$Y = X\beta + \varepsilon,$$

where $Y$ is $n$-vector of responses, $\beta$ is $p$-vector of parameters, $X = (x_1 \ldots x_p)$ is $n \times p$ matrix with column vectors $x_1, \ldots, x_p$, and $\varepsilon$ is $n$-vector of random errors with $\mathrm{E}(\varepsilon) = 0$ and $\mathrm{var}(\varepsilon) = \sigma^2 I$.

Least-squares (LS) estimator $\widehat{\beta}_{\mathrm{ols}} = (X^T X)^{-1} X^T y$, if $X$ is of full rank, is BLUE (best linear unbiased estimator) . $\mathrm{var}(\widehat{\beta}_{\mathrm{ols}}) = (X^T X)^{-1} \sigma^2$.
If column vectors $x_1, \ldots, x_p$ are close to (but not exactly) linearly dependent, the vectors are said to be collinear. The determinant $\det(X^T X)$ is close to 0. Then $\mathrm{var}(\widehat{\beta}) \uparrow$.

# Problems of LS estimator $\widehat{\beta}$ with collinearity

- Large variance and mean squared error.

$$\mathrm{MSE} = \mathrm{bias}^2 + \mathrm{var.}$$

- Poor estimation and prediction.

- Three major phenomena (Land *et al.*1990, AJS):

◇ Large changes in parameter estimate when adding or deleting variables;

◇ Wide confidence interval, nonsignificant test statistics, and opposite signs to expected values of important independent variables;

◇ Unstable regression parameters from sample to sample.

## Diagnosis: condition number

Let $\lambda_1 \leq \ldots \leq \lambda_p$ be ordered eigenvalues of matrix $X^T X$. The condition number is defined as $\sqrt{\lambda_p/\lambda_1}$. Cutoff: 30.

## Q: How to improve performance?

James – Stein estimator.

If $\widehat{\theta} = x$ is an unbiased estimator for $\theta$ and $p \geq 3$, then $J_x = \left(1 - \frac{p-2}{\|x\|_2^2}\right) x$ is called James – Stein estimator.

Shrinkage estimators.
  Idea: Shrink parameters towards the origin
       to reduce variance (bias-variance trade-off).
  Recall: $\mathrm{MSE} = \mathrm{bias}^2 + \mathrm{var}.$

# $L^\gamma$ Penalty

## Ridge estimator (Hoerl and Kennard 1971)

$$\widehat{\beta}_{\mathbf{rdg}} = (X^T X + \lambda I)^{-1} X^T y,$$

where $I$ is identity matrix, $\lambda \geq 0$ is tuning parameter.

$$\widehat{\beta}_{\mathbf{rdg}} = \arg \min_{\beta} \{ (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \}.$$

Equivalently,
$$\widehat{\beta}_{\mathbf{rdg}} = \arg \min_{\beta} \{ (y - X\beta)^T (y - X\beta) \} \text{ subject to } \beta^T \beta \leq t,$$
with $t \geq 0$.

$$\mathrm{var}(\widehat{\beta}_{\mathbf{rdg}}) \leq \mathrm{var}(\widehat{\beta}_{\mathbf{ols}})$$

# $L^\gamma$ PENALTY

**Bridge estimator (Frank and Friedman 1993)**

$$\widehat{\beta}_{\mathrm{brdg}} = \arg\min_{\beta}\{(y - X\beta)^T(y - X\beta) + \lambda \sum_{j=1}^{p} |\beta_j|^\gamma\} .$$

Equivalently,

$$\widehat{\beta}_{\mathrm{brdg}} = \arg\min_{\beta}\{(y - X\beta)^T(y - X\beta)\} \text{ subject to}$$

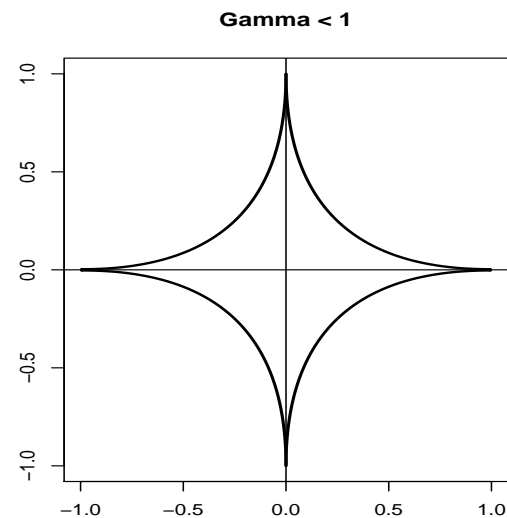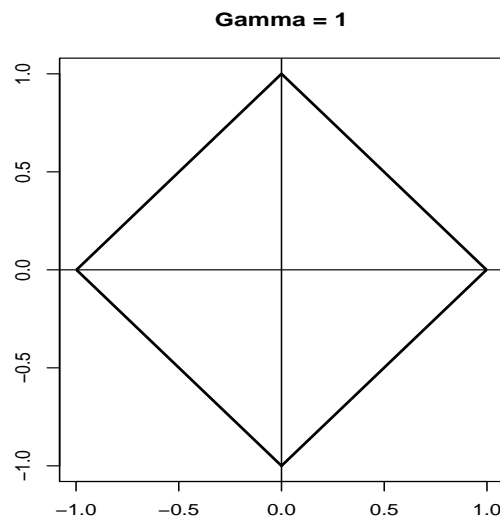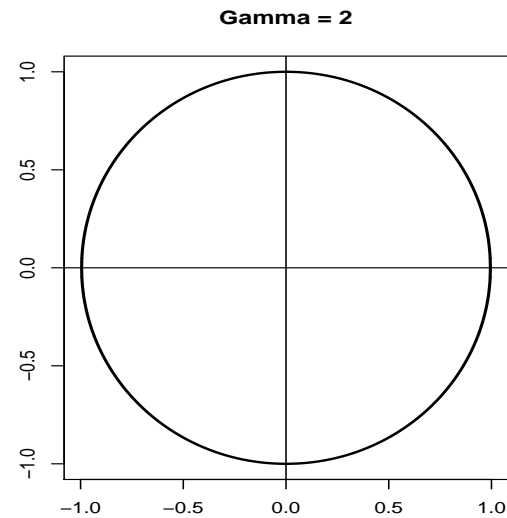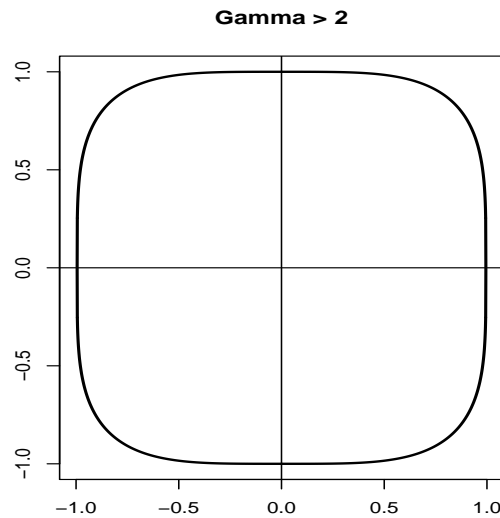$\sum_{j=1}^{p} |\beta_j|^\gamma \leq t$ , with $t \geq 0$.

$$\mathrm{var}(\widehat{\beta}_{\mathrm{brdg}}) \leq \mathrm{var}(\widehat{\beta}_{\mathrm{ols}})$$

**Bridge – generalization of ridge**

$\diamond\, \gamma = 2$, ridge;

$\diamond\, \gamma = 1$, lasso (Tibshirani 1996).

STATISTICS
TEXAS A&M UNIVERSITY

# $L^\gamma$ PENALTY

## Constraint area for different values of $\gamma > 0$.



Gamma > 2    Gamma = 2

Gamma = 1    Gamma < 1

# $L^\gamma$ PENALTY

Variable selection property of lasso $\widehat{\beta}_j = 0$ .

# $L^\gamma$ PENALTY

## Computation for bridge $\gamma > 1$

- $\gamma = 2$: closed form.
- $\gamma > 1$: modified Newton-Raphson (Fu 1998), complex!

## Notations:

$RSS = (y - X\beta)^T(y - X\beta),\ S_j = \partial RSS/\partial\beta_j,$

$d(\beta_j, \lambda, \gamma) = \lambda\gamma|\beta_j|^{\gamma-1}\text{sign}(\beta_j),\ l_j = S_j + d(\beta_j, \lambda, \gamma).$

Solve system of equations:

$$
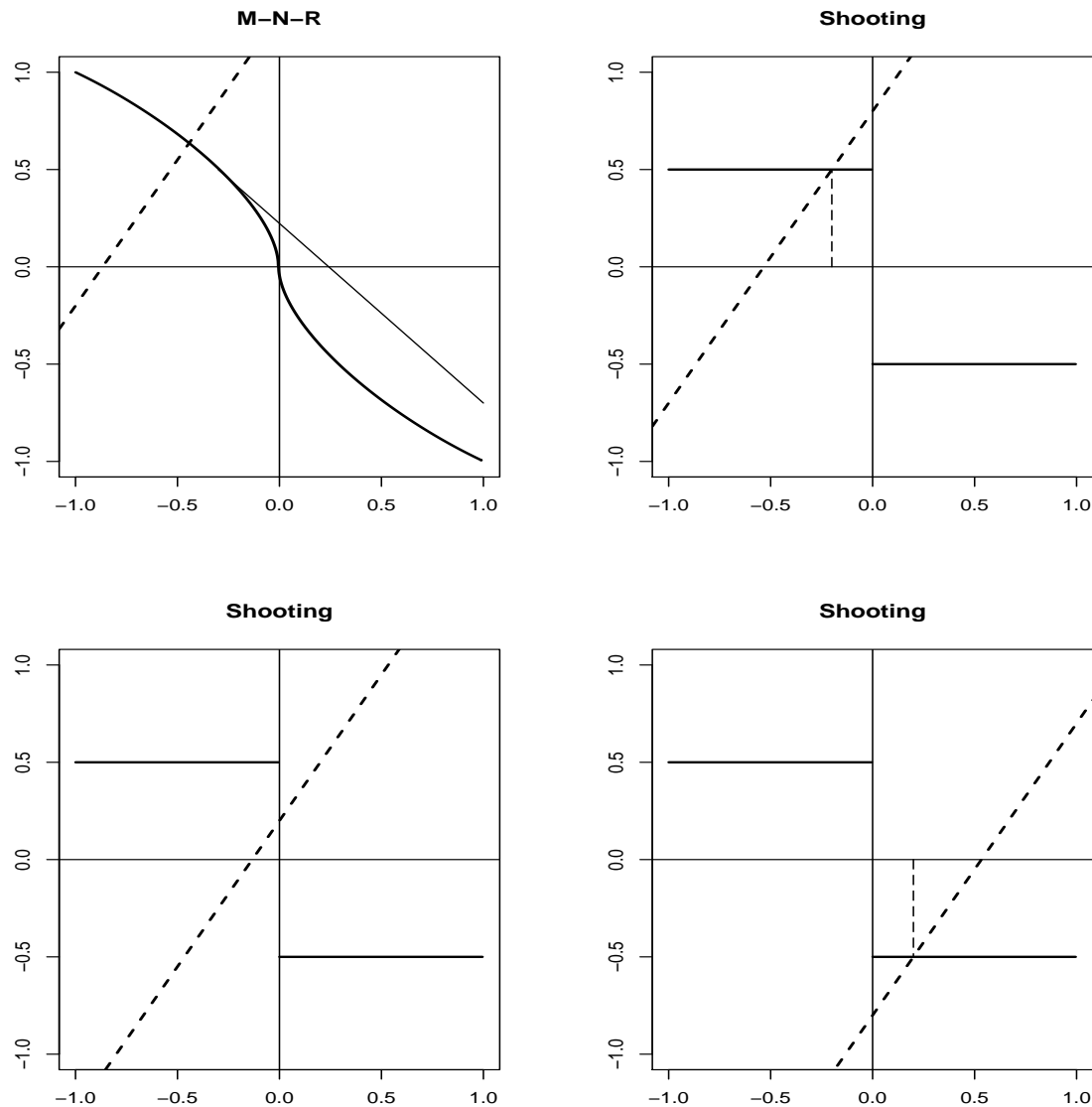\begin{cases}
l_1(\beta, X, y, \lambda, \gamma) = 0, \\
\qquad\qquad \ldots \\
l_p(\beta, X, y, \lambda, \gamma) = 0.
\end{cases}
\qquad (1)
$$

No closed form. Use N–R. $\beta_j^{new} = \beta_j^{old} - [\partial l_j/\partial\beta_j]^{-1}l_j$

Modify N-R since convexity changes at $\beta_j = 0$ for $1 < \gamma < 2$.

# $L^\gamma$ PENALTY

## M-N-R ($\gamma > 1$) and shooting algorithm ($\gamma = 1$).

# $L^\gamma$ PENALTY

## Computation for lasso $\gamma = 1$

- Combined quadratic programming (Tibshirani 1996).
  Quadratic programming:
  $$\min(y - X\beta)^T(y - X\beta) \text{ subject to } v^T\beta \geq 0.$$
  Constraint $\sum_{j=1}^{p} |\beta_j| \leq t$ is equivalent to
  $\sum_{j=1}^{p} w_j\beta_j \leq t$ with $w_j = \pm 1$.
  Total combinations of $2^p$ weights $w_j$. Complicated!

- Shooting algorithm (Fu 1998).
  Take limit $\gamma \to 1+$:

  not computationally – more complicated;
  but theoretically – iteration with simple closed form.

STATISTICS
TEXAS A&M UNIVERSITY

# $L^\gamma$ PENALTY

## Theorem 1

If $S_j$ is contin. diff., Jacobian $\partial S/\partial \beta$ pos-semi-def., then

1. $\widehat{\beta}(\lambda, \gamma)$ is unique and contin. in $(\lambda, \gamma)$.
2. $\lim_{\gamma \to 1+} \widehat{\beta}(\lambda, \gamma)$ exists for fixed $\lambda > 0$.
3. $\lim_{\gamma \to 1+} \widehat{\beta}(\lambda, \gamma) = \widehat{\beta}(\lambda, 1)$, the lasso estimator for L–S.

## Implication

1. Penalty (shrinkage) models do not need joint likelihood. Only Jacobian $\partial S/\partial \beta$ condition (p.s.d.). Potential extension!
2. If joint likelihood exists, the extension works perfectly.

STATISTICS
TEXAS A&M UNIVERSITY

# $L^\gamma$ Penalty

## Shooting algorithm for lasso

1). Start with $\widehat{\beta}^{(0)} = (\widehat{\beta}_1, \ldots, \widehat{\beta}_p)$.

2). At step $m$, for $j = 1, \ldots, p$, let $s_0 = S_j(0, \widehat{\beta}^{(-j)}, X, y)$ and $x_j$ be the $j$–th column vector of $X$. Set

$$
\widehat{\beta}_j = \begin{cases} \frac{\lambda - s_0}{2x_j^T x_j} & \text{if} \quad s_0 > \lambda \\ 0 & \text{if} \quad |s_0| \le \lambda \\ \frac{-\lambda - s_0}{2x_j^T x_j} & \text{if} \quad s_0 < -\lambda \end{cases}
$$

Form a new estimator $\widehat{\beta}^{(m)} = (\widehat{\beta}_1, \ldots, \widehat{\beta}_p)$ after updating all $\widehat{\beta}_j$.

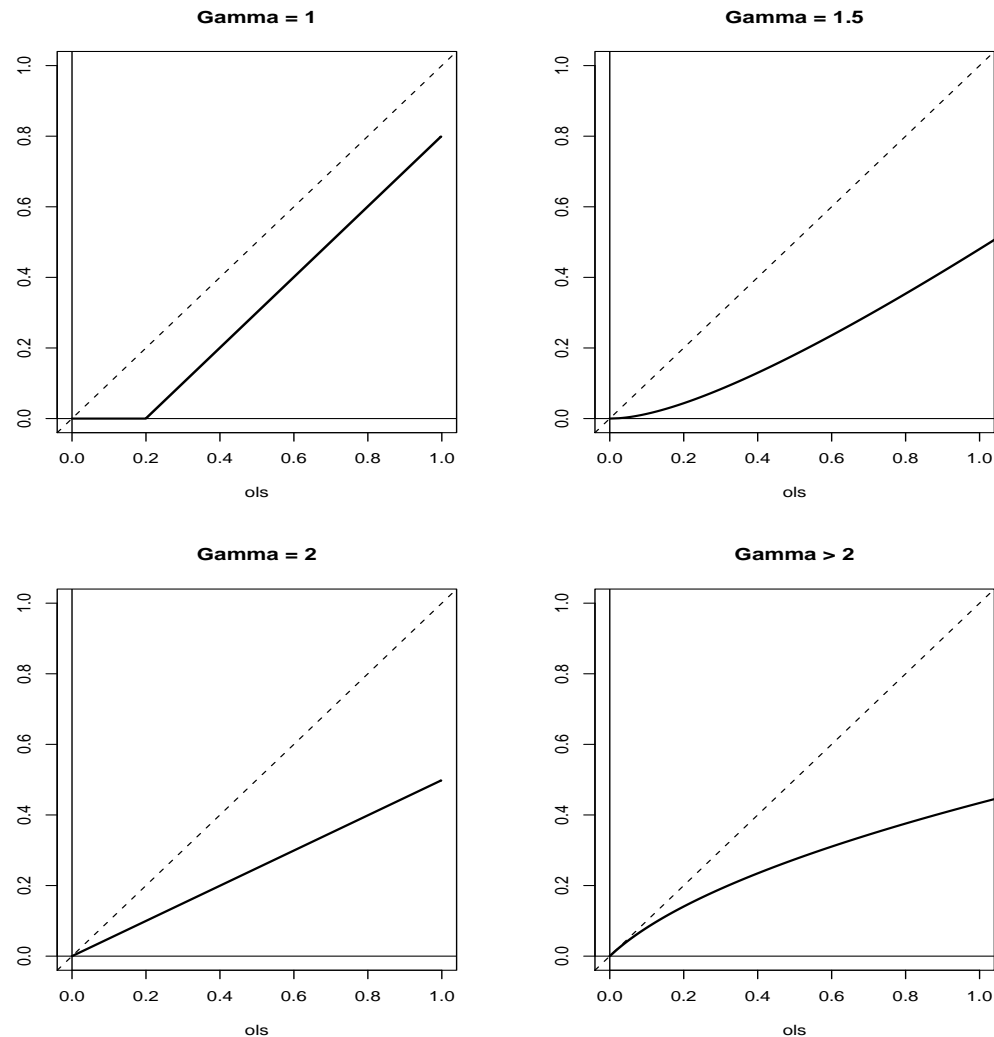3). Repeat step 2) until convergence of $\widehat{\beta}^{(m)}$.

## Convergence of algorithms

Let $G(\beta; \lambda, \gamma) = (y - X\beta)^T(y - X\beta) + \lambda \sum_j |\beta_j|^\gamma$ for given $\lambda > 0$ and $\gamma \geq 1$. $G(\beta; \lambda, \gamma)$ is convex and is minimized at finite $\beta = \beta_0$. Each step of updating $\widehat{\beta}_j$ through either M-N-R algorithm or the shooting algorithm decreases the function $G(\beta; \lambda, \gamma)$. Thus the estimator $\widehat{\beta}_m$

converges.

# $L^{\gamma}$ PENALTY

## Orthonormal matrix $X$: $X^T X = I$

Coordinate: $\widehat{\beta}_{\mathrm{brdg}} = \widehat{\beta}_{\mathrm{ols}} - \lambda\gamma/2|\widehat{\beta}_{\mathrm{brdg}}|^{\gamma-1}\mathrm{sign}(\widehat{\beta}_{\mathrm{brdg}})$

# $L^\gamma$ Penalty

## Variance of bridge estimator

● $\gamma > 1$, complex closed form: no zero-valued coordinates. (Fu 1998).

$$\mathrm{var}(\widehat{\beta}) =$$
$$(X^T X + D(\widehat{\beta})|_{y_0})^{-1} X^T X (X^T X + D(\widehat{\beta})|_{y_0})^{-1}\sigma^2, \quad (2)$$
$$D(\widehat{\beta}) = \lambda\gamma(\gamma - 1)/2 \operatorname{diag}(|\beta_j|^{\gamma-2}) .$$

$y_0$ is some point in the sample space.

Difficult to use for lasso due to $\widehat{\beta}_j = 0$.

## Variance of bridge estimator

$\bullet$ $\gamma = 1$, difficulty: zero-valued coordinates.

Method in Tibshirani (1996):

$$\mathrm{var}(\widehat{\beta}) =$$
$$(X^T X + \lambda W^-)^{-1} X^T X (X^T X + \lambda W^-)^{-1} \sigma^2, \qquad (3)$$
where $W = \mathrm{diag}(|\beta_1|, \ldots, |\beta_p|)$

Method in Osborne (2000);

Let

$$W = \frac{X^T (y - X\widehat{\beta})(y - X\widehat{\beta})^T X}{\|\widehat{\beta}\|_1 \|X^T(y - X\widehat{\beta})\|_\infty}$$

$$\mathrm{var}(\widehat{\beta}) = (X^T X + W)^{-1} X^T X (X^T X + W)^{-1} \sigma^2. \qquad (4)$$

# Comparison between two methods

- (3) is zero for $\widehat{\beta}_j = 0$, while (4) is non-zero.

- However, if set tuning parameter $\lambda > 0$ large, all $\beta_j = 0$. Then no variability. Hence variance should be zero. (3) is acceptable, but (4) still non-zero.

- All the above methods are approximations. No exact results except for $\gamma = 2$.

- Bootstrap method usually yields good estimation.

# $L^\gamma$ PENALTY

## Table 1. Analysis of prostate cancer data.

| Predictor | $\widehat{\beta}^a$ | SE by (1) | SE by (2) | $\widehat{\beta}^b$ | SE by bootstrap |
|-----------|------|-----------|-----------|------|-----------------|
| Intercept | 2.478 | 0.072 | 0.072 | 2.478 | 0.072 |
| lcavol | 0.559 | 0.079 | 0.101 | 0.618 | 0.103 |
| lweight | 0.097 | 0.060 | 0.081 | 0.190 | 0.076 |
| age | 0 | 0 | 0.079 | -0.048 | 0.046 |
| lbph | 0 | 0 | 0.080 | 0.103 | 0.066 |
| svi | 0.156 | 0.071 | 0.097 | 0.245 | 0.087 |
| lcp | 0 | 0 | 0.125 | 0 | 0.068 |
| gleason | 0 | 0 | 0.114 | 0 | 0.047 |
| pgg45 | 0 | 0 | 0.123 | 0.063 | 0.056 |

[a] Osborne (2000) ($t = 0.8114$). [b] Fu (1998) ($\lambda = 7.2$).

**STATISTICS**
TEXAS A&M UNIVERSITY

## Shrinkage trace

— Parameter estimates change with a special tuning parameter: standard shrinkage rate, $0 \leq s \leq 1$.
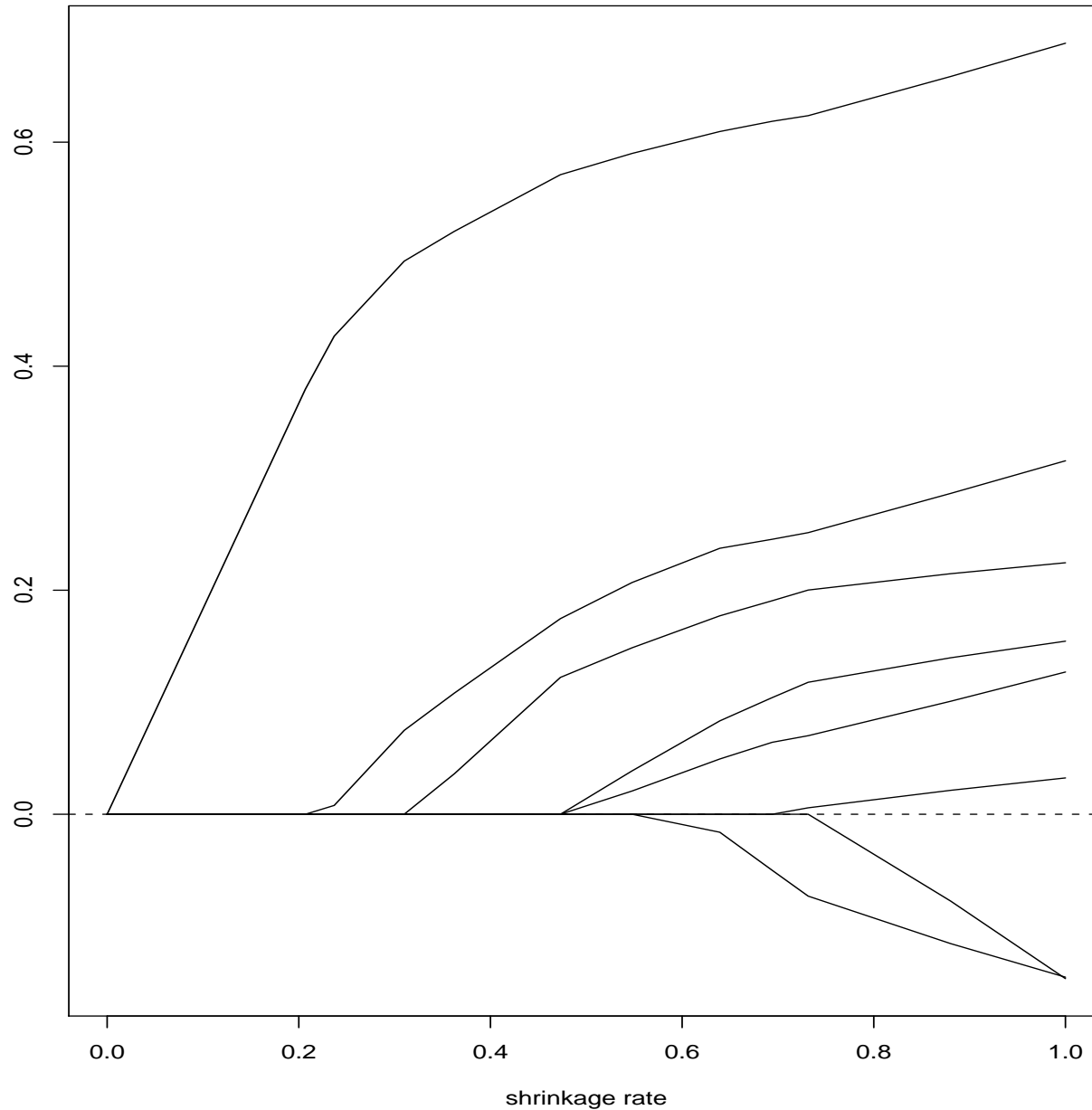
$$s = \frac{\|\beta_j(\lambda, \gamma)\|_\gamma}{\|\beta_j(\lambda = 0, \gamma)\|_\gamma},$$

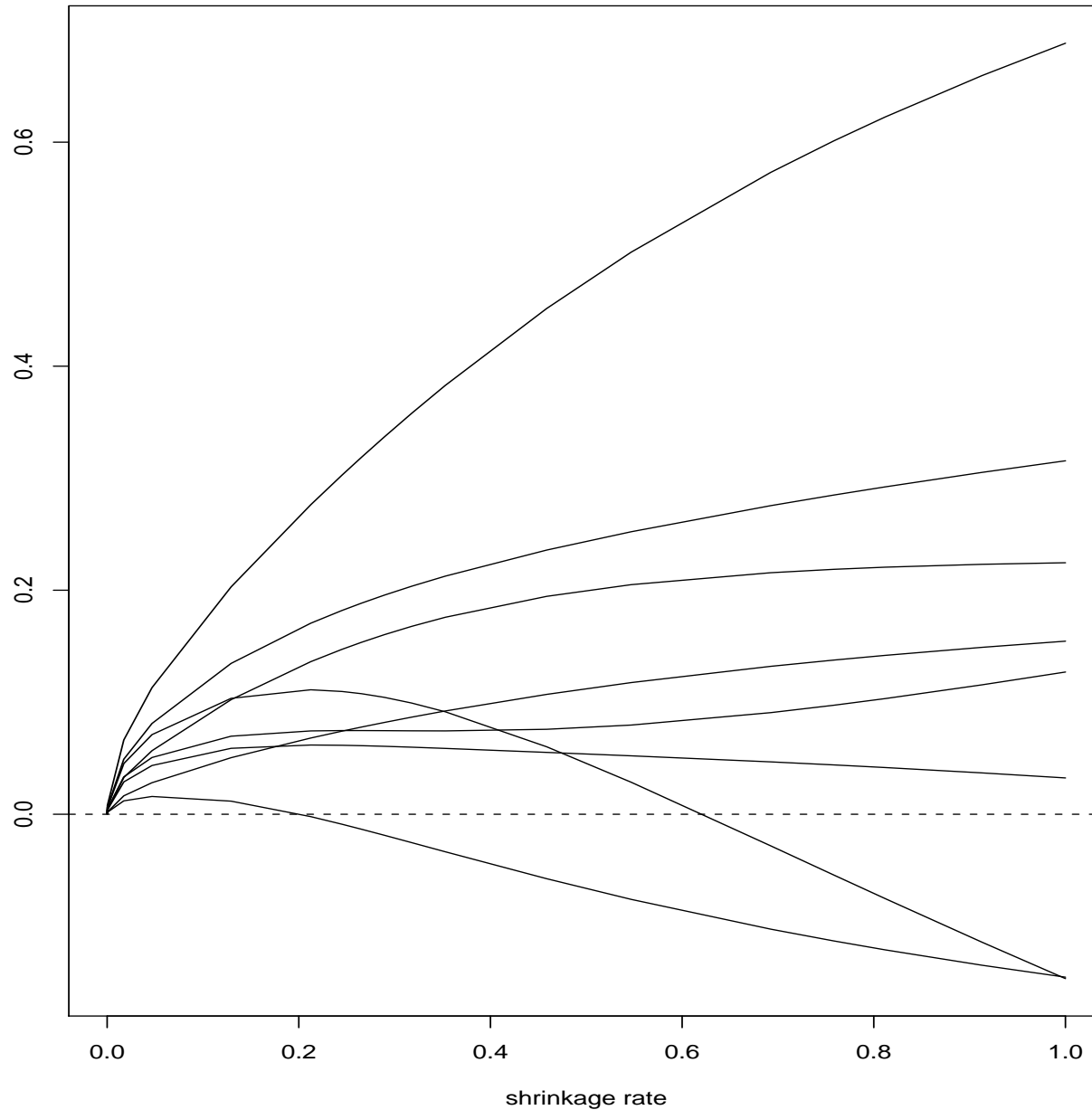where $\|\cdot\|_\gamma$ is the $L^\gamma$ norm of a $p$-vector.

$s = 0$, full shrinkage.

$s = 1$, no shrinkage.

shrinkage rate

STATISTICS
TEXAS A&M UNIVERSITY

shrinkage rate

STATISTICS
TEXAS A&M UNIVERSITY

## Seemingly contradictory results:

Given data $(X, y)$: $\quad Y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$.

Let $(x, y)$ be an arbitrary point in sample space.

$\diamond$ Prediction error increases with collinearity.
$$\text{PSE}(x) = \sigma^2 + \text{MSE}(x) = \sigma^2[1 + x^T(X^TX)^{-1}x]$$
$$= \sigma^2 + x^T\text{var}(\widehat{\beta})x.$$

$\diamond$ Prediction error at given data points is constant.
$$\frac{1}{n}\sum_1^n \text{PSE}(x_i) = \sigma^2[1 + \frac{1}{n}\sum_1^n x_i^T(X^TX)^{-1}x_i]$$
$$= \sigma^2[1 + \frac{1}{n}\sum_1^n \text{tr}\{(X^TX)^{-1}x_ix_i^T\} = \sigma^2(1 + \frac{p}{n}),$$
where $x_i$ are row vectors of matrix $X$.

$\diamond$ In fact, $\text{E}[\text{PSE}(x)] = \sigma^2(1 + \frac{p}{n})$.

## Collinearity increases variability of PSE:

**Proposition** (Fu 2005)

Assume existence of two moments $\mathrm{E}(xx^T)$ and $\mathrm{E}(xx^Txx^T)$. The expectation $\mathrm{E}\{\mathrm{PSE}(x)\}$ is independent of the collinearity for large samples with $\mathrm{E}\{\mathrm{PSE}(x)\} \sim \sigma^2(1 + p/n)$. The variance $\mathrm{var}\{\mathrm{PSE}(x)\}$ increases with the collinearity as the smallest eigenvalue of matrix $X^TX$ decreases to 0.

$$\mathrm{var}\{\mathrm{PSE}(x)\} = \mathrm{E}[\{\mathrm{PSE}(x)\}^2] - [\mathrm{E}\{\mathrm{PSE}(x)\}]^2$$
$$\sim \frac{\sigma^4}{n^2}[\mathrm{tr}\{V^{-1}\mathrm{E}(xx^TV^{-1}xx^T)\} - p^2],$$
where $V = X^TX/n$.

## Collinearity increases variability of PSE:

Assume $V = \mathrm{diag}(\lambda_1, \ldots, \lambda_p)$ with $\lambda_1 \geq \ldots \geq \lambda_p > 0$, without loss of generality.

Let $U = \mathrm{E}(xx^T V^{-1} xx^T) - \mathrm{E}(\lambda_1^{-1} xx^T xx^T)$, psd.

$U$ and $V$ can be diagonalized simultaneously. $V^{-1}U$ is psd.

$$
\begin{aligned}
\mathrm{tr}\{V^{-1}\mathrm{E}(xx^T V^{-1} xx^T)\} &\geq \mathrm{tr}\{V^{-1}\lambda_1^{-1}\mathrm{E}(xx^T xx^T)\} \\
&= \lambda_1^{-1}[\lambda_1^{-1}c_1 + \cdots + \lambda_p^{-1}c_p] \\
&> \lambda_1^{-1}\lambda_p^{-1}c_p \to \infty \text{ as } \lambda_p \to 0,
\end{aligned}
$$

where $c_1, \cdots, c_p > 0$ are the elements on the main diagonal of matrix $\mathrm{E}(xx^T xx^T)$ and are independent of matrix $V$.

# REFERENCES

Frank, I.E. and Friedman, J.H. (1993). A statistical view of some chemometrics regression tools, *Technometrics* 35:109-148.

Fu, W.J. (1998). Penalized regressions: the Bridge versus the Lasso, *J. Comp. Grap. Statist.* 7: 397-416.

Fu, W.J. (2005). Prediction error with collinearity, *Comm. Statist. - Theor. Meth.*, in press.

Gruber, M.H.J. (1990). *Regression Estimators: A Comparative Study*, Academic Press, Boston.

Hoerl, A.E. and Kennard, R.W. (1970a). Ridge regression: biased estimation for nonorthogonal problems, *Technometrics*, 12:55-67.

Hoerl, A.E. and Kennard, R.W. (1970b). Ridge regression: applications to nonorthogonal problems, *Technometrics*, 12:69-82.

James, W. and Stein, C (1961). Estimation with quadratic loss, *Proc. Fourth Berkeley Symp. Math. Statist. Prob.* 1, 311-319.

Land, K.C. McCall, P.L. and Cohen, L.E. (1990). Structural covariates of homicide rates: are there any invariances across time and social space? *American Journal of Sociology*, 95, 922-963.

Miller, A.J. (1990). *Subset Selection in Regression,* Chapman and Hall, New York.

Osborne, M. R., Presnell, B. and Turlach, B. A. (2000). On the LASSO and its dual. *Journal of Computational and Graphical Statistics* **9**, 319-337.

Sen, A. and Srivastava, M. (1990). *Regression Analysis: Theory, Methods, and Applications,* Springer, New York.

Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso, *J. Roy. Statist. Soc. B* 58:267-288.

**STATISTICS**
TEXAS A&M UNIVERSITY