Shrinkage Priors and Selection Readings Chapter 15 Christensen

STA721 Linear Models Duke University

Merlise Clyde

October 24, 2017

duke.eps

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

$\mathbf{Y} \mid \alpha, \boldsymbol{\beta}^{s}, \phi \quad \sim \quad \mathsf{N}(\mathbf{1}_{n}\alpha + \mathbf{X}^{s}\boldsymbol{\beta}^{s}, \mathbf{I}_{n}/\phi)$



$\mathbf{Y} \mid \alpha, \beta^{s}, \phi \sim \mathsf{N}(\mathbf{1}_{n}\alpha + \mathbf{X}^{s}\beta^{s}, \mathbf{I}_{n}/\phi)$ $\beta^{s} \mid \alpha, \phi, \tau, \lambda \sim \mathsf{N}(\mathbf{0}, \mathsf{diag}(\tau^{2})/\phi)$



$$\begin{aligned} \mathbf{Y} \mid \alpha, \boldsymbol{\beta}^{s}, \phi &\sim & \mathsf{N}(\mathbf{1}_{n}\alpha + \mathbf{X}^{s}\boldsymbol{\beta}^{s}, \mathbf{I}_{n}/\phi) \\ \boldsymbol{\beta}^{s} \mid \alpha, \phi, \boldsymbol{\tau}, \lambda &\sim & \mathsf{N}(\mathbf{0}, \mathsf{diag}(\boldsymbol{\tau}^{2})/\phi) \\ & p(\alpha, \phi) &\propto & 1/\phi \end{aligned}$$



$$\begin{array}{lll} \mathbf{Y} \mid \alpha, \beta^{s}, \phi & \sim & \mathsf{N}(\mathbf{1}_{n}\alpha + \mathbf{X}^{s}\beta^{s}, \mathbf{I}_{n}/\phi) \\ \beta^{s} \mid \alpha, \phi, \tau, \lambda & \sim & \mathsf{N}(\mathbf{0}, \mathsf{diag}(\tau^{2})/\phi) \\ & \rho(\alpha, \phi) & \propto & 1/\phi \end{array}$$

prior on τ_j



$$\begin{array}{rcl} \mathbf{Y} \mid \alpha, \beta^{s}, \phi & \sim & \mathsf{N}(\mathbf{1}_{n}\alpha + \mathbf{X}^{s}\beta^{s}, \mathbf{I}_{n}/\phi) \\ \beta^{s} \mid \alpha, \phi, \tau, \lambda & \sim & \mathsf{N}(\mathbf{0}, \mathsf{diag}(\tau^{2})/\phi) \\ & \rho(\alpha, \phi) & \propto & 1/\phi \end{array}$$

prior on τ_j Scale Mixture of Normals (Andrews and Mallows 1974)



Carvalho, Polson & Scott propose

Prior Distribution on

$$eta^{s} \mid \phi, au \sim \mathsf{N}(\mathbf{0}_{p}, rac{\mathsf{diag}(m{ au}^{2})}{\phi})$$



Carvalho, Polson & Scott propose

Prior Distribution on

$$oldsymbol{eta}^{s} \mid \phi, oldsymbol{ au} \sim \mathsf{N}(oldsymbol{0}_{oldsymbol{
ho}}, rac{\mathsf{diag}(oldsymbol{ au}^2)}{\phi})$$

• $\tau_j \mid \lambda \stackrel{\text{iid}}{\sim} C^+(0, \lambda^2)$ (difference in CPS notation)



Carvalho, Polson & Scott propose

Prior Distribution on

$$oldsymbol{eta}^{s} \mid \phi, oldsymbol{ au} \sim \mathsf{N}(oldsymbol{0}_{oldsymbol{
ho}}, rac{\mathsf{diag}(oldsymbol{ au}^2)}{\phi})$$

► $\tau_j \mid \lambda \stackrel{\text{iid}}{\sim} C^+(0, \lambda^2)$ (difference in CPS notation) ► $\lambda \sim C^+(0, 1)$



Carvalho, Polson & Scott propose

Prior Distribution on

$$oldsymbol{eta}^{s} \mid \phi, oldsymbol{ au} \sim \mathsf{N}(oldsymbol{0}_{oldsymbol{
ho}}, rac{\mathsf{diag}(oldsymbol{ au}^2)}{\phi})$$

duke ens

τ_j | $\lambda \stackrel{\text{iid}}{\sim} C^+(0, \lambda^2)$ (difference in CPS notation) $\lambda \sim C^+(0, 1)$ $p(\alpha, \phi) \propto 1/\phi$)

Carvalho, Polson & Scott propose

Prior Distribution on

$$oldsymbol{eta}^{s} \mid \phi, oldsymbol{ au} \sim \mathsf{N}(oldsymbol{0}_{p}, rac{\mathsf{diag}(oldsymbol{ au}^{2})}{\phi})$$

duke ens

► $\tau_j \mid \lambda \stackrel{\text{iid}}{\sim} C^+(0, \lambda^2)$ (difference in CPS notation) ► $\lambda \sim C^+(0, 1)$

•
$$p(\alpha, \phi) \propto 1/\phi)$$

In the case $\lambda=\phi=1$ and with canonical representation ${\bf Y}={\bf I}{\beta}+\epsilon$

Carvalho, Polson & Scott propose

Prior Distribution on

$$oldsymbol{eta}^{s} \mid \phi, oldsymbol{ au} \sim \mathsf{N}(oldsymbol{0}_{p}, rac{\mathsf{diag}(oldsymbol{ au}^{2})}{\phi})$$

► $\tau_j \mid \lambda \stackrel{\text{iid}}{\sim} C^+(0, \lambda^2)$ (difference in CPS notation) ► $\lambda \sim C^+(0, 1)$

•
$$p(\alpha, \phi) \propto 1/\phi)$$

In the case $\lambda=\phi=1$ and with canonical representation $\mathbf{Y}=\mathbf{I}\boldsymbol{\beta}+\boldsymbol{\epsilon}$

$$E[\beta_i \mid \mathbf{Y}] = \int_0^1 (1 - \kappa_i) y_i^* p(\kappa_i \mid \mathbf{Y}) \ d\kappa_i = (1 - \mathsf{E}[\kappa \mid y_i^*]) y_i^*$$

where $\kappa_i = 1/(1 + \tau_i^2)$ shrinkage factor

Carvalho, Polson & Scott propose

Prior Distribution on

$$oldsymbol{eta}^{s} \mid \phi, oldsymbol{ au} \sim \mathsf{N}(oldsymbol{0}_{p}, rac{\mathsf{diag}(oldsymbol{ au}^{2})}{\phi})$$

•
$$\tau_j \mid \lambda \stackrel{\text{iid}}{\sim} C^+(0, \lambda^2)$$
 (difference in CPS notation)
• $\lambda \sim C^+(0, 1)$

•
$$p(\alpha, \phi) \propto 1/\phi)$$

In the case $\lambda=\phi=1$ and with canonical representation $\mathbf{Y}=\mathbf{I}\boldsymbol{\beta}+\boldsymbol{\epsilon}$

$$E[\beta_i \mid \mathbf{Y}] = \int_0^1 (1 - \kappa_i) y_i^* p(\kappa_i \mid \mathbf{Y}) \ d\kappa_i = (1 - \mathsf{E}[\kappa \mid y_i^*]) y_i^*$$

where $\kappa_i = 1/(1 + \tau_i^2)$ shrinkage factor

Half-Cauchy prior induces a Beta(1/2, 1/2) distribution on κ_i a priori

duke ens

10 ω 9 density 4 N − 0.0 0.2 0.4 0.6 0.8 1.0

Beta(1/2, 1/2)

duke.eps

æ

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

х

Prior Comparison (from PSC)

Comparison of different priors



Tails of different priors



duke.eps

э

Normal means case $Y_i \stackrel{\text{iid}}{\sim} N(\beta_i, 1)$ (Equivalent to Canonical case)

duke ens

Posterior mean
 E[β | y] = y + d/dy log m(y)
 where m(y) is the
 predictive density under
 the prior (known λ)

Normal means case $Y_i \stackrel{\mathrm{iid}}{\sim} \mathcal{N}(\beta_i, 1)$ (Equivalent to Canonical case)

< ロ > < (型 > < 注 > < 注 > < 注 > 注

HS has Bounded Influence:

$$\lim_{|y|\to\infty}\frac{d}{dy}\log m(y)=0$$

Normal means case $Y_i \stackrel{\mathrm{iid}}{\sim} \mathcal{N}(\beta_i, 1)$ (Equivalent to Canonical case)

< ロ > < (型 > < 注 > < 注 > < 注 > 注

HS has Bounded Influence:

$$\lim_{|y|\to\infty}\frac{d}{dy}\log m(y)=0$$

•
$$\lim_{|y|\to\infty} E[\beta \mid y) \to y$$

(MLE)

Normal means case $Y_i \stackrel{\text{iid}}{\sim} N(\beta_i, 1)$ (Equivalent to Canonical case)

- Posterior mean
 E[β | y] = y + d/dy log m(y)
 where m(y) is the
 predictive density under
 the prior (known λ)
- HS has Bounded Influence:

$$\lim_{|y|\to\infty}\frac{d}{dy}\log m(y)=0$$

- $\lim_{|y| \to \infty} E[\beta \mid y) \to y$ (MLE)
- DE is also bounded influence, but bound does not decay to zero in tails



duke eps

R packages

The monomvn package in R includes

duke.eps

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- blasso
- bhs

See Diabetes.R code

Range of other scale mixtures used



Range of other scale mixtures used

Generalized Double Pareto (Armagan, Dunson & Lee)

Range of other scale mixtures used

► Generalized Double Pareto (Armagan, Dunson & Lee)

$$egin{aligned} & au_j^2 \mid \lambda \sim \mathsf{Exp}(\lambda^2/2) \ & \lambda \sim \mathsf{Gamma}(lpha,\eta) \ & eta_j^s \sim \mathsf{GDP}(\xi = \eta/lpha, lpha) \end{aligned}$$

duke.eps

Range of other scale mixtures used

► Generalized Double Pareto (Armagan, Dunson & Lee)

$$\begin{split} \tau_j^2 \mid \lambda &\sim \mathsf{Exp}(\lambda^2/2) \\ \lambda &\sim \mathsf{Gamma}(\alpha, \eta) \\ \beta_j^s &\sim \mathsf{GDP}(\xi = \eta/\alpha, \alpha) \end{split}$$

$$f(\beta_j^s) = \frac{1}{2\xi} (1 + \frac{|\beta_j^s|}{\xi\alpha})^{-(1+\alpha)}$$

duke ens

see http://arxiv.org/pdf/1104.0861.pdf

Range of other scale mixtures used

► Generalized Double Pareto (Armagan, Dunson & Lee)

$$\begin{split} \tau_j^2 \mid \lambda &\sim \mathsf{Exp}(\lambda^2/2) \\ \lambda &\sim \mathsf{Gamma}(\alpha, \eta) \\ \beta_j^s &\sim \mathsf{GDP}(\xi = \eta/\alpha, \alpha) \end{split}$$

$$f(\beta_j^s) = \frac{1}{2\xi} \left(1 + \frac{|\beta_j^s|}{\xi\alpha}\right)^{-(1+\alpha)}$$

duke ens

see http://arxiv.org/pdf/1104.0861.pdf

► Normal-Exponenetial-Gamma (Griffen & Brown 2005) λ² ~ Gamma(α, η)

Range of other scale mixtures used

► Generalized Double Pareto (Armagan, Dunson & Lee)

$$\begin{split} \tau_j^2 \mid \lambda &\sim \mathsf{Exp}(\lambda^2/2) \\ \lambda &\sim \mathsf{Gamma}(\alpha, \eta) \\ \beta_j^s &\sim \mathsf{GDP}(\xi = \eta/\alpha, \alpha) \end{split}$$

$$f(\beta_j^s) = \frac{1}{2\xi} \left(1 + \frac{|\beta_j^s|}{\xi\alpha}\right)^{-(1+\alpha)}$$

see http://arxiv.org/pdf/1104.0861.pdf

- ► Normal-Exponenetial-Gamma (Griffen & Brown 2005) λ² ~ Gamma(α, η)
- Bridge Power Exponential Priors (Stable mixing density)

Range of other scale mixtures used

► Generalized Double Pareto (Armagan, Dunson & Lee)

$$\begin{aligned} \tau_j^2 \mid \lambda \sim \mathsf{Exp}(\lambda^2/2) \\ \lambda \sim \mathsf{Gamma}(\alpha, \eta) \\ \beta_j^s \sim \mathsf{GDP}(\xi = \eta/\alpha, \alpha) \end{aligned}$$

$$f(\beta_j^s) = \frac{1}{2\xi} \left(1 + \frac{|\beta_j^s|}{\xi\alpha}\right)^{-(1+\alpha)}$$

see http://arxiv.org/pdf/1104.0861.pdf

► Normal-Exponenetial-Gamma (Griffen & Brown 2005) λ² ~ Gamma(α, η)

Bridge - Power Exponential Priors (Stable mixing density)
 See the monomvn package on CRAN

Range of other scale mixtures used

► Generalized Double Pareto (Armagan, Dunson & Lee)

$$\begin{aligned} \tau_j^2 \mid \lambda \sim \mathsf{Exp}(\lambda^2/2) \\ \lambda \sim \mathsf{Gamma}(\alpha, \eta) \\ \beta_j^s \sim \mathsf{GDP}(\xi = \eta/\alpha, \alpha) \end{aligned}$$

$$f(\beta_j^s) = \frac{1}{2\xi} \left(1 + \frac{|\beta_j^s|}{\xi\alpha}\right)^{-(1+\alpha)}$$

see http://arxiv.org/pdf/1104.0861.pdf

► Normal-Exponential-Gamma (Griffen & Brown 2005) λ² ~ Gamma(α, η)

Bridge - Power Exponential Priors (Stable mixing density)
 See the monomvn package on CRAN
 Choice of prior? Properties?

Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties



Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties

• Model
$$Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$



Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties

• Model
$$Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

• Assume $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ (orthonormal) and $\epsilon \sim N(0, \mathbf{I}_n)$



Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties

- Model $Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$
- Assume $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ (orthonormal) and $\boldsymbol{\epsilon} \sim N(0, \mathbf{I}_n)$
- Penalized Likelihood

$$\frac{1}{2} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \frac{1}{2} \sum_j (\beta_j - \hat{\beta}_j)^2 + \sum_j p_\lambda(|\beta_j|)$$

とうかん かんかく ふく きょう ふしゃ

Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties

- Model $Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$
- Assume $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ (orthonormal) and $\boldsymbol{\epsilon} \sim N(0, \mathbf{I}_n)$
- Penalized Likelihood

$$\frac{1}{2} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \frac{1}{2} \sum_j (\beta_j - \hat{\beta}_j)^2 + \sum_j p_\lambda(|\beta_j|)$$

duality $p_{\lambda}(|\beta|)$ is negative log prior

Requirements on penality

Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties

- Model $Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$
- Assume $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ (orthonormal) and $\boldsymbol{\epsilon} \sim N(0, \mathbf{I}_n)$
- Penalized Likelihood

$$\frac{1}{2} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \frac{1}{2} \sum_j (\beta_j - \hat{\beta}_j)^2 + \sum_j p_\lambda(|\beta_j|)$$

duality $p_{\lambda}(|\beta|)$ is negative log prior

- Requirements on penality
 - Unbiasedness: for large $|\beta_j|$

Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties

- Model $Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$
- Assume $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ (orthonormal) and $\boldsymbol{\epsilon} \sim N(0, \mathbf{I}_n)$
- Penalized Likelihood

$$\frac{1}{2} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \frac{1}{2} \sum_j (\beta_j - \hat{\beta}_j)^2 + \sum_j p_\lambda(|\beta_j|)$$

duality $p_{\lambda}(|\beta|)$ is negative log prior

- Requirements on penality
 - Unbiasedness: for large $|\beta_j|$
 - Sparsity: thresholding rule sets small coefficients to 0

Fan & Li (JASA 2001) discuss Variable selection via nonconcave penalties and oracle properties

- Model $Y = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$
- Assume $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ (orthonormal) and $\boldsymbol{\epsilon} \sim N(0, \mathbf{I}_n)$
- Penalized Likelihood

$$\frac{1}{2} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \frac{1}{2} \sum_j (\beta_j - \hat{\beta}_j)^2 + \sum_j p_\lambda(|\beta_j|)$$

duality $p_{\lambda}(|\beta|)$ is negative log prior

- Requirements on penality
 - Unbiasedness: for large $|\beta_j|$
 - Sparsity: thresholding rule sets small coefficients to 0

とうかん かんかく ふく きょう ふしゃ

• Continuity: continuous in $\hat{\beta}_j$

Conditions

Derivative of
$$\frac{1}{2} \sum_{j} (\beta_j - \hat{\beta}_j)^2 + \sum_{j} p_{\lambda}(|\beta_j|)$$
 is
 $\operatorname{sgn}(\beta_j) \{ |\beta_j| + p'_{\lambda}(|\beta_j|) \} - \hat{\beta}_j$

Conditions:

- unbiased: if $p'_{\lambda}(|\beta|) = 0$ for large $|\beta|$; estimator is $\hat{\beta}_j$
- ► thresholding: min $\{|\beta_j| + p'_{\lambda}(|\beta_j|)\} > 0$ then estimator is 0 if $|\hat{\beta}_j| < \min\{|\beta_j| + p'_{\lambda}(|\beta_j|)\}$

duke ens

• continuity: minimum of $|\beta_j| + p'_{\lambda}(|\beta_j|)$ is at zero

Choice?

- Lasso does not satisfy conditions
- GDP does



Posterior Mode (may set some coefficients to zero)



- Posterior Mode (may set some coefficients to zero)
- Posterior Mean (no selection, just shrinkage) (Squared error loss)

duke ens

► Minimize L₁ posterior loss E[|β_j - a|] (Shrinkage and Selection)

- Posterior Mode (may set some coefficients to zero)
- Posterior Mean (no selection, just shrinkage) (Squared error loss)

► Minimize L₁ posterior loss E[|β_j - a|] (Shrinkage and Selection)

Bayesian Posterior does not assign any probability to $\beta_i^s = 0$

- Posterior Mode (may set some coefficients to zero)
- Posterior Mean (no selection, just shrinkage) (Squared error loss)

► Minimize L₁ posterior loss E[|β_j - a|] (Shrinkage and Selection)

Bayesian Posterior does not assign any probability to $\beta_i^s = 0$

Selection solved as a post-analysis decision problem

- Posterior Mode (may set some coefficients to zero)
- Posterior Mean (no selection, just shrinkage) (Squared error loss)

► Minimize L₁ posterior loss E[|β_j - a|] (Shrinkage and Selection)

Bayesian Posterior does not assign any probability to $\beta_i^s = 0$

- Selection solved as a post-analysis decision problem
- Selection part of model uncertainty \Rightarrow add prior

- Posterior Mode (may set some coefficients to zero)
- Posterior Mean (no selection, just shrinkage) (Squared error loss)
- ► Minimize L₁ posterior loss E[|β_j a|] (Shrinkage and Selection)

Bayesian Posterior does not assign any probability to $\beta_i^s = 0$

- Selection solved as a post-analysis decision problem
- Selection part of model uncertainty ⇒ add prior probability that β^s_i = 0 and combine with decision problem

Remember all models are wrong, but some may be useful!