# STAT 525 Lecture 17

October 26, 2017

## 1   Last time

$$\boldsymbol{y} \sim N\left(\boldsymbol{X}\boldsymbol{\beta}, \sigma^2 \mathbb{I}\right)$$
$$\boldsymbol{\beta} \sim N\left(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta\right)$$

$$\boldsymbol{\mu}_\beta = 0 \Rightarrow \begin{cases} \text{g-prior} & \boldsymbol{\Sigma}_\beta = g\sigma^2 \left[\boldsymbol{X}^T \boldsymbol{X}\right]^{-1} \\ \text{ridge} & \boldsymbol{\Sigma}_\beta = \sigma_\beta^2 \mathbb{I} \end{cases}$$

Let $\tau = \sigma^{-2}$ then

$$\boldsymbol{y} \sim N\left(\boldsymbol{X}\boldsymbol{\beta}, \tau^{-1}\mathbb{I}_n\right)$$
$$\boldsymbol{\beta} \sim N\left(\boldsymbol{\mu}_\beta, \tau^{-1}\boldsymbol{\Sigma}_\beta\right)$$
$$\tau \sim \text{Gamma}\left(\alpha_\tau, \beta_\tau\right)$$

Then the posterior would be

$$
\begin{aligned}
p\left(\tau \mid \boldsymbol{y}, \boldsymbol{X}, \boldsymbol{\beta}\right) &\propto p\left(\tau, \boldsymbol{\beta}\right) p\left(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \tau\right) \\
&= p\left(\tau\right) p\left(\boldsymbol{\beta} \mid \tau\right) p\left(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \tau\right) \\
&\propto \tau^{\alpha_\tau - 1} \exp\left(-\beta_\tau \tau\right) \left|2\pi\tau^{-1}\mathbb{I}_n\right|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\tau \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2\right) \times \\
&\quad \left|2\pi\tau^{-1}\boldsymbol{\Sigma}_\beta\right|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\left(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta\right)^T \left[\tau^{-1}\boldsymbol{\Sigma}_\beta\right]^{-1} \left(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta\right)\right) \\
&\propto \tau^{\alpha_\tau + \frac{n}{2} + \frac{p}{2} - 1} \exp\left(-\tau\left[\beta_\tau + \frac{1}{2}\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \frac{1}{2}\left(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta\right)^T \left[\boldsymbol{\Sigma}_\beta\right]^{-1} \left(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta\right)\right]\right)
\end{aligned}
$$

Therefore,

$$\tilde{\alpha}_\tau = \alpha_\tau + \frac{n}{2} + \frac{p}{2}$$
$$\tilde{\beta}_\tau = \beta_\tau + \frac{1}{2}\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \frac{1}{2}\left(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta\right)^T \left[\boldsymbol{\Sigma}_\beta\right]^{-1} \left(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta\right)$$

Instead if we don't have $\tau$ in the prior on beta then

$$\boldsymbol{\beta} \sim N\left(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta\right)$$

$$\tilde{\alpha}_\tau = \alpha_\tau + \frac{n}{2}$$

$$\tilde{\beta}_\tau = \beta_\tau + \frac{1}{2}\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2$$

If we have the following $g$ prior

$$\boldsymbol{\mu}_\beta = 0$$

$$\boldsymbol{\Sigma}_\beta = g\sigma^2 \left[\boldsymbol{X}^T\boldsymbol{X}\right]^{-1}$$

After marginizing $\boldsymbol{\beta}$ out, we have

$$[\tau \mid \boldsymbol{y}, \boldsymbol{X}] \sim \text{Gamma}\left(\alpha_\tau + \frac{n}{2}, \beta_\tau + SSR_g\right)$$

where

$$SSR_g = \boldsymbol{y}^T \left[\mathbb{I}_n - \frac{g}{g+1}\boldsymbol{X}\left(\boldsymbol{X}^T\boldsymbol{X}\right)^{-1}\boldsymbol{X}\right]\boldsymbol{y}$$

$$= \|y\|^2 - \frac{g}{g+1}\boldsymbol{y}^T\left[\boldsymbol{X}\hat{\boldsymbol{\beta}}_{\text{ols}}\right]$$

# 2    Prediction

Given new $\tilde{\boldsymbol{X}}$ to predict $\tilde{\boldsymbol{y}}$ or $p\left(\tilde{\boldsymbol{y}} \mid \boldsymbol{y}\right)$

Sources of uncertainty

1. Model variability $\sigma^2$ (not accounted for by $\boldsymbol{X}\boldsymbol{\beta}$)

2. Posterior uncertainty in $p\left(\boldsymbol{\beta}, \sigma^2 \mid \boldsymbol{y}\right)$ (due to finite sample size)

$$\left[\tilde{\boldsymbol{y}} \mid \tilde{\boldsymbol{X}}, \beta, \sigma^2\right] \sim N\left(\tilde{\boldsymbol{X}}\boldsymbol{\beta}, \sigma^2\mathbb{I}\right)$$

where we simulate $\boldsymbol{\beta}$, $\sigma$ and $\tilde{\boldsymbol{y}}$.

# 3    Bayesian Robustness

Instead of $\epsilon_i \stackrel{\text{iid}}{\sim} N\left(0, \sigma^2\right)$ try $\epsilon_i \stackrel{\text{iid}}{\sim} t_v\left(0, \sigma^2\right)$ where $\nu$ is the degree of freedom. Can be augmented as

$$\epsilon_i \mid \xi_i \stackrel{\text{iid}}{\sim} N\left(0, \sigma^2/\xi_i\right)$$

$$\xi_i \stackrel{\text{iid}}{\sim} \text{Gamma}\left(\frac{\nu}{2}, \frac{\nu}{2}\right)$$

(Scaled Mixture of Gaussian)

1. (Asymmertic) Laplace

2. Skew normal

3. Discrete mixtures of Gaussian

# 4 Shrinkage and penalized regression

- Many predictors $(p \gg n)$ but may be unrelated to $\boldsymbol{y}$.

- Including unnecessary predictors. Can cause poor performance.

- Nice to represent a small sets of predictors.

**Key trade off for several methods**

- Discrete (or two groups model)
$$p\,(\beta_j = 0) > 0$$

  - Can select $\{j : \beta_j \neq 0\}$
  - Problem: computationally feasible for moderate $p$

- Continuous (One group)

  - no true zero but small $|\beta_j| \approx 0$
  - Scalable

**Penalized Regression**

Setting:
$$\hat{\boldsymbol{\beta}} = \arg\min_{\boldsymbol{\beta}} \mathcal{L}\,(\boldsymbol{y}, \boldsymbol{X}\boldsymbol{\beta}) + \lambda P\,(\boldsymbol{\beta})$$

where $\mathcal{L}\,(\boldsymbol{y}, \boldsymbol{X}\boldsymbol{\beta})$ is the loss function (corresponding to (-log) likelihood) and $P\,(\boldsymbol{\beta})$ is penalty (corresponding to (-log) prior). $\lambda$control the trade off (corresponding to prior precision)

$$\lambda \to 0\mathcal{L} \text{ dominates}$$
$$\lambda \to \infty P \text{ dominates}$$

## 4.1 Lasso Regression
$$\hat{\boldsymbol{\beta}}_L = \arg\min_{\boldsymbol{\beta}} \mathcal{L}\,(\boldsymbol{y}, \boldsymbol{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^{p} |\beta_j|$$

1. unpenalized intercept ($\bar{\boldsymbol{y}}$ centered)

2. variable $X$ should be on the same scale

3. penalized MLE (post mode) produces sparse solution.

## 4.2 Bayesian Lasso (Park & Casella 2008)

$$\left[\beta_j \mid \sigma^2, \lambda\right] \overset{\text{iid}}{\sim} DE\,(\lambda/\sigma)$$

where
$$p\,(\beta_j) = \frac{\lambda}{2\sigma} \exp\left(-\lambda\,|\beta_j|\,/\sigma\right)$$

notice that

$$- \log P(\boldsymbol{\beta}) = -\sum_{j=1}^{p} \log\left(p\left(\beta_j\right)\right)$$

$$= A^{-1} \sum_{j=1}^{p} |\beta_j|$$

which is equivalent to

$$[\beta_j \mid \sigma, \eta_j] \overset{\text{indep}}{\sim} N\left(0, \sigma^2/\eta_j\right)$$

$$[\eta_j^{-1} \mid \lambda] \overset{\text{indep}}{\sim} Exp\left(\lambda^2/2\right)$$

$$\boldsymbol{\beta} \sim N\left(0, \boldsymbol{\Sigma}_\beta\right)$$

where

$$\boldsymbol{\Sigma}_\beta = \begin{bmatrix} \sigma^2/\eta_1 & & & \\ & \sigma^2/\eta_2 & & \\ & & & \\ & & & \sigma^2/\eta_p \end{bmatrix}$$

- $\eta_j$ is inverse-Gaussian

- do not get true zeros in $\boldsymbol{\beta}$

- but do get SEs (even when $|\beta_j| \approx 0$)

## 4.3   Horseshoe prior

$$\left[\beta_j \mid \sigma^2, \lambda_j^2\right] \overset{\text{indep}}{\sim} N\left(-, \sigma^2 \lambda_j^2\right)$$

$$\lambda_j \overset{\text{indep}}{\sim} C^+(0, A)$$

Priors for $A$:

- $A \sim C^+(0, 1)$

- $A \sim \text{Uniform}(0, 1)$

$$\boldsymbol{\beta} \sim N\left(0, \boldsymbol{\Sigma}_\beta\right)$$

$$\boldsymbol{\Sigma}_\beta = \text{diag}\left(\sigma^2 \lambda_1, \sigma^2 \lambda_2, \cdots, \sigma^2 \lambda_p\right)$$

Comments

- Many priors in this family: global local priors

$$\beta_j \sim N\left(0, \sigma^2 A^2 \lambda_j^2\right)$$

where $p\left(A, \lambda_1, \cdots, \lambda_p\right)$ determines behavior of $\boldsymbol{\beta}$.

- Do we need $\beta_j = 0$ or is $|\beta_j| \approx 0$ good enough?

- Threshold procedures:

– Recall

$$\kappa_j = \frac{1}{1 + \lambda_j^2} \sim \text{Beta}\left(\frac{1}{2}, \frac{1}{2}\right)$$

is the amount of shrinkage to zero

– Reasonable:

$$\hat{b}_j = \mathbb{I}\left(\kappa_j < \frac{1}{2}\right)\hat{\beta}_j$$

where $\hat{\beta}_j$ is the posterior mean.