

---

# $L^\gamma$ Penalty Models

## Computation And Applications

### Part IV

**Wenjiang Fu**

**wfu@stat.tamu.edu**

**<http://stat.tamu.edu/~wfu>**

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

## Bayesian interpretation

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

where  $Y$ :  $n$ -vector of responses,  $X$ : regression matrix,  $\beta = (\beta_1, \dots, \beta_p)^T$   $p$ -vector of parameters,  $\varepsilon$ :  $n$ -vector of random errors with mean  $E(\varepsilon) = 0$  and  $\text{var}(\varepsilon) = \sigma_0^2 I$ .

$L^\gamma$  penalty has a Bayesian prior interpretation (Lectures 1–2).

Assume  $\varepsilon \sim N(0, \sigma_0^2 I)$ . Let  $\beta_j \sim L^\gamma$  prior,  $j = 1, \dots, p$ .

Study posterior  $\pi(\beta|y)$  for given data  $y$ .

However, it is difficult to compute the posterior due to lack of conjugate property in general.

Notice that two members of the  $L^\gamma$  family are special and play a major role: lasso ( $\gamma = 1$ ) and ridge ( $\gamma = 2$ ), which correspond to Laplacian and Gaussian priors, respectively.

## Laplacian prior

We study a novel family of priors including Laplacian and Gaussian as special cases.

## Why Laplacian prior?

- ◇ Achieve variable selection. Same idea for Lasso.
- ◇ Representation by simple distributions.

$\text{Lap}(1) \stackrel{d}{=} N(0, 2\Lambda)$  with  $\Lambda \sim \text{Exp}(1)$  (Kotz et al. 2000).

- ◇ Studied for variable selection with applications to microarray studies (Bae and Mallick 2004).

# BAYESIAN APPROACH

---

## Extension to a Bayesian prior family

Special properties of  $\text{Gamma}(\lambda, k)$ :

Mean  $\mu = \lambda/k$  and variance  $\sigma^2 = \lambda/k^2$ .

How to achieve  $N(0, 2C)$ ?

Consider  $\text{Gamma}(1 + Ct, 1 + t)$  with  $t \geq 0$  and constant  $C > 0$ .

Two special cases:

◇  $t = 0$ .  $\beta_j \sim \text{Lap}(1) \stackrel{d}{=} \text{Gamma}(1, 1)$ .

◇  $t \rightarrow \infty$ ,  $\text{Gamma}(1 + Ct, 1 + t) \xrightarrow{p} C$ .  $\beta_j \sim N(0, 2C)$ .

For  $\text{Gamma}(\lambda, k)$ ,  $\mu = \lambda/k$ ,  $\sigma^2 = \lambda/k^2$ .

# BAYESIAN APPROACH

---

## Laplacian – Gaussian mixture (LGM) prior

For given  $t \geq 0$  and constant  $C > 0$ .

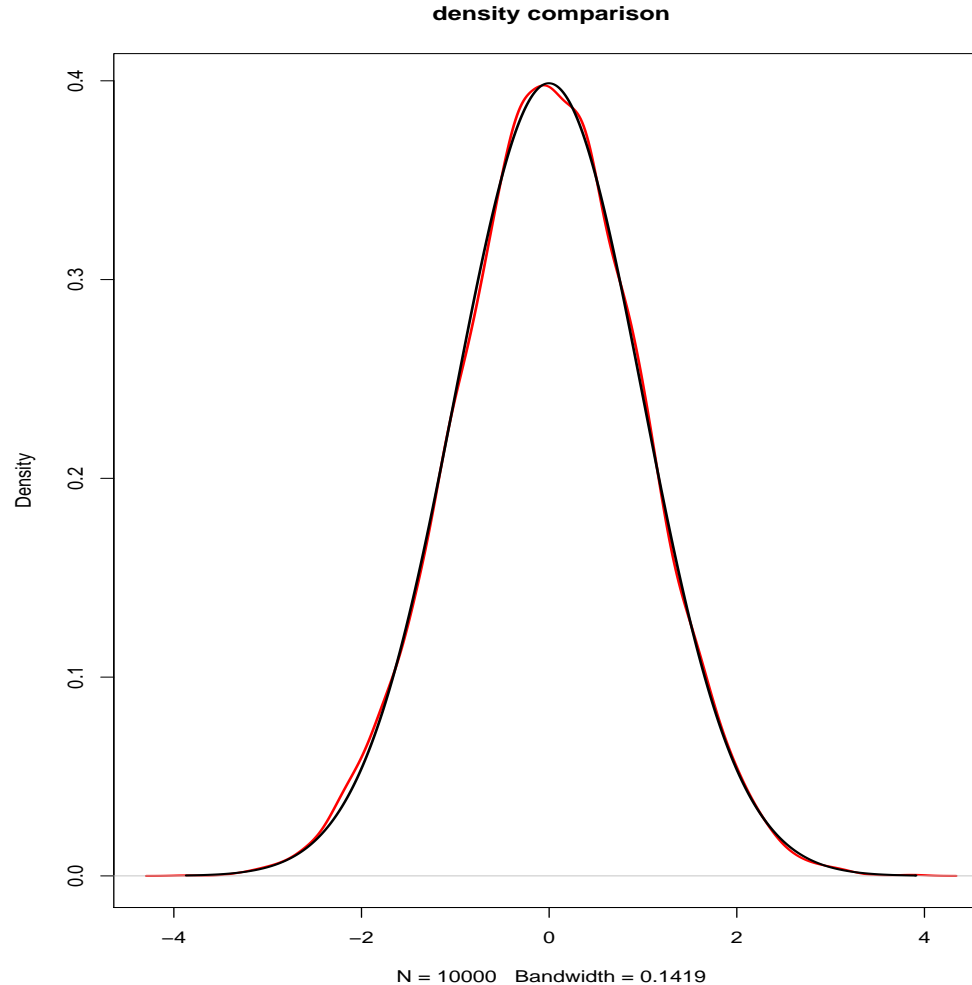
Consider  $\beta_j \sim N(0, 2\Lambda) \cdot \text{Gamma}(1 + Ct, 1 + t)$ ,  
the Laplacian–Gaussian mixture prior.

- ◇ LGM is a natural extension of Laplacian to a family including Gaussian.
- ◇ Posterior computation will take the advantage of  $\beta_j \sim N(0, 2\Lambda)$ .
- ◇  $L^\gamma$  is another one, but the posterior is difficult to handle. If use  $\beta_j \sim \exp(-\lambda|\beta|^\gamma)$ , it will involve stable distribution, complicated and difficult.

# BAYESIAN APPROACH

## Laplacian – Gaussian mixture (LGM) prior

LGM behaves like a Gaussian for  $t \geq 300$ .



Black curve:  $N(0, 1)$ ; Red curve: LGM  $t = 300$

## Advantages of LGM prior

- **Prior on hyperparameter  $t \geq 0$ :  $\pi(t)$ .**
  - ◇  $\pi(t)$ : point mass at  $t = 0$ , Laplacian prior.
  - ◇  $\pi(t)$ : point mass at  $t = t_0$  large ( $t_0 \geq 300$ ), Gaussian prior.
  - ◇  $\pi(t)$ : point mass at  $t = 0$  and  $t = t_0 > 0$  large.  
Combines Lasso and ridge, Elastic Net (Zou and Hastie 2004).
  - ◇  $\pi(t)$ : continuous  $t \geq 0$ , Bayesian model averaging.

## Posterior with LGM prior

Denote  $\Lambda^{-1} = \text{diag}(\Lambda_1^{-1}, \dots, \Lambda_p^{-1})$ , the inverse of the diagonal matrix of elements  $(\Lambda_1, \dots, \Lambda_p)$ .

$$\begin{aligned}\pi(\boldsymbol{\beta}, \Lambda | \mathbf{y}, t) &\propto \exp\left[-\frac{1}{2\sigma_0^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right] \\ &\times \exp\left(-\frac{1}{4}\boldsymbol{\beta}^T \Lambda^{-1} \boldsymbol{\beta}\right) \frac{(1+t)^{p(1+ct)}}{[\Gamma(1+ct)]^p} \\ &\times \Lambda_1^t \dots \Lambda_p^t \exp[-(1+t)(\Lambda_1 + \dots + \Lambda_p)].\end{aligned}$$

$$\pi(\boldsymbol{\beta} | \mathbf{y}, t) = \int_{\Lambda} \pi(\boldsymbol{\beta}, \Lambda | \mathbf{y}, t) d\Lambda$$



## Bayesian variable selection (SSVS)

**SSVS** (Stochastic search variable selection) (George and McCulloch 1993)

For given  $\gamma_j$  Bernoulli (0 or 1) – **index for variable selection.**

$\beta_j | \gamma_j \sim (1 - \gamma_j)N(0, \tau_j^2) + \gamma_j N(0, c_j^2 \tau_j^2)$ ,  $c_j > 0$  large,  
 $P(\gamma_j = 1) = 1 - P(\gamma_j = 0) = p_j$ .

**Priors:**  $\beta | \gamma \sim N_p(\mathbf{0}, D_\gamma R D_\gamma)$ ,

Variance component:  $\sigma^2 | \gamma \sim \text{IG}(\nu_\gamma/2, \nu_\gamma \lambda_\gamma/2)$ ,

$\gamma \sim f(\gamma) = \prod p_j^{\gamma_j} (1 - p_j)^{(1-\gamma_j)}$ .

## Computational Methods for SSVS

### Gibbs sampling for best subset:

$$\beta^0, \sigma^0, \gamma^0, \beta^1, \sigma^1, \gamma^1, \dots, \beta^m, \sigma^m, \gamma^m, \dots,$$

Variable selection by determining posterior distribution of  $\gamma$ .

Computationally intensive !

### Metropolis – Hastings search.

Brown, Vannucci and Fearn (1998, 2002),

Brown, Fearn and Vannucci (1999),

Vannucci, Brown and Fearn (2001),

Lee, Sha, Dougherty, Vannucci and Mallick (2003).

## Representation of Laplace Distribution

Theorem 7 (Kotz 2000)

A standard classical Laplace r.v.  $X$  has the representation  $X \stackrel{d}{=} \sqrt{2W}Z$ , where the r.v.s  $W$  and  $Z$  have the standard exponential and normal distributions, respectively.

The moment generating function of exponential  $W$  is  $M_w(t) = (1 - t)^{-1}$ ,  $t < 1$ . Characteristic function of normal  $Z$  is  $\exp(-t^2/2)$ . The characteristic function of  $\sqrt{2W}Z$

$$\begin{aligned} \mathbb{E}[\exp(it\sqrt{2W}Z)] &= \mathbb{E}\{\mathbb{E}[\exp(it\sqrt{2W}Z)|W]\} \\ &= \mathbb{E}[\phi_z(t\sqrt{2W})] = \mathbb{E}[\exp(-t^2W)] \\ &= M_W(-t^2) = (1 + t^2)^{-1}, \end{aligned}$$

where  $\phi_z(t) = \exp(-t^2/2)$  is the characteristic function of standard normal.

# BAYESIAN APPROACH

---

## Representation of Laplace Distribution

The density of  $X = \sqrt{2W}Z$  is given by

$$\int_0^{\infty} \frac{1}{2\sqrt{\pi w}} \exp\left[-\frac{1}{2}\left(\frac{x^2}{2w} + 2w\right)\right] dw.$$

Consider transformation  $Y_1 = W, Y_2 = \sqrt{2W}Z$ . Calculate joint density and the marginal of  $Y_2$  by integrating out  $Y_1$ .

# FUSED LASSO

**Fused Lasso** (Tibshirani, Saunders, Rosset and Zhu 2005)

Consider linear model

$$y_i = \sum X_{ij}\beta_j + \varepsilon_i,$$

where  $x_j = (x_{1j}, \dots, x_{nj})$  are standardized and ordered variables (Protein mass spectroscopy data: time of flight with mass/charge  $m/z$ ).

**Idea** Penalize both the parameters  $|\beta_j|$  and their differences  $|\beta_j - \beta_{j-1}|$ .

Min RSS subject to  $\sum |\beta_j| \leq s_1$  and  $\sum |\beta_j - \beta_{j-1}| \leq s_2$ .

**Goal** Achieve sparsity and smoothness.

Two special cases:

1). Lasso:  $s_2$  is large; 2). Fusion:  $s_1$  is large (Land and Friedman 1996).

# FUSED LASSO

**Performance** using prostate cancer protein mass

spectroscopy data with random split of training 216 + test 108 samples (total 157 healthy + 167 cancer) (Tibshirani 2005)

Model	test err	df	sites	$s_1$	$s_2$
Lasso	6/108	116	116	144	262
Fusion	19/108	168	171	175	200
Fused Lasso	6/108	122	344	184	222

# FUSED LASSO

**Comparison** with leukemia classification using microarrays

(training: 27+11; test: 34) (Tibshirani 2005)

Method	$s_1$	$s_2$	10-FdCV	Test err	genes
Golub (50 genes)			3/38	4/34	50
Lasso 37 df	0.65	1.32	1/38	1/34	37
Fused Lasso 38 df	1.08	0.71	1/38	2/34	135
Fused Lasso 20 df	1.35	1.01	1/38	4/34	737

# ELASTIC NET MODEL

**Elastic Net (ENet)** (Zou and Hastie 2004)

$$L(\beta, \lambda_1, \lambda_2) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda_1 |\beta|_1 + \lambda_2 |\beta|_2^2,$$

where  $|\beta|_1 = \sum |\beta_j|$ ,  $|\beta|_2^2 = \sum \beta_j^2$ .

Naive elastic net estimator is defined as

$$\hat{\beta} = \operatorname{argmin}_{\beta} L(\beta, \lambda_1, \lambda_2).$$

It's equivalent to

$$\hat{\beta} = \operatorname{argmin}_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \text{ subject to}$$

$$(1 - \alpha) |\beta|_1 + \alpha |\beta|_2^2 \leq t \text{ for some } t.$$

ENet combines Lasso penalty and ridge penalty.



# ELASTIC NET MODEL

## Algorithm for Naive ENet

Given data  $(\mathbf{y}, \mathbf{X})$  and fixed  $(\lambda_1, \lambda_2)$ . Define artificial data set  $(\mathbf{y}^*, \mathbf{X}^*)$  by

$$\mathbf{X}_{(n+p) \times p}^* = (1 + \lambda_2)^{-1/2} \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda_2} \mathbf{I} \end{pmatrix}, \quad \mathbf{y}_{(n+p)}^* = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}.$$

Let  $\gamma = \lambda_1 / \sqrt{1 + \lambda_2}$  and  $\beta^* = \sqrt{1 + \lambda_2} \beta$ .

$$L^*(\beta^*, \gamma) = (\mathbf{y}^* - \mathbf{X}^* \beta^*)^T (\mathbf{y}^* - \mathbf{X}^* \beta^*) + \gamma |\beta^*|_1.$$

Let  $\hat{\beta}^* = \operatorname{argmin}_{\beta^*} L(\beta^*, \gamma)$ , then the ENet estimator

$$\hat{\beta} = \frac{1}{\sqrt{1 + \lambda_2}} \hat{\beta}^*.$$

# ELASTIC NET MODEL

## Relationship with Lasso estimator

For orthonormal design matrix  $X$ ,

$$\hat{\beta}_j(\text{NENet}) = \frac{(|\hat{\beta}_j(\text{ols})| - \lambda_1/2)_+}{1 + \lambda_2} \text{sign}(\hat{\beta}_j(\text{ols})).$$

Two special cases:

- 1)  $\lambda_1 = 0$ , ridge estimator  $\hat{\beta}(\text{ridge}) = 1/(1 + \lambda_2)\hat{\beta}(\text{ols})$ ;
- 2)  $\lambda_2 = 0$ , lasso estimator

$$\hat{\beta}_j(\text{lasso}) = (|\hat{\beta}_j(\text{ols})| - \lambda_1/2)_+ \text{sign}(\hat{\beta}_j(\text{ols})).$$

# ELASTIC NET MODEL

## Grouping effect of ENet

Given data  $(y, X)$  and  $(\lambda_1, \lambda_2)$  with centered  $y$  and standardized  $X$ . Let  $\hat{\beta}(\lambda_1, \lambda_2)$  be the NENet estimator. Suppose  $\hat{\beta}_i(\lambda_1, \lambda_2)\hat{\beta}_j(\lambda_1, \lambda_2) > 0$ . Define

$$D_{\lambda_1, \lambda_2}(i, j) = \frac{1}{|y|_1} |\hat{\beta}_i(\lambda_1, \lambda_2) - \hat{\beta}_j(\lambda_1, \lambda_2)|,$$

then  $D_{\lambda_1, \lambda_2}(i, j) \leq [\sqrt{2(1 - \rho)}] / \lambda_2$ ,  
where  $\rho = x_i^T x_j$ , the sample correlation.  
Highly correlated covariates tend to be selected together.

# ELASTIC NET MODEL

## ENet estimator

Given data  $(y, X)$  and  $(\lambda_1, \lambda_2)$ , and augmented data  $(y^*, X^*)$ .  
Naive ENet estimator

$$\hat{\beta}^* = \operatorname{argmin}_{\beta^*} (y^* - X^* \beta^*)^T (y^* - X^* \beta^*) + \frac{\lambda_1}{\sqrt{1 + \lambda_2}} |\beta^*|_1.$$

$$\hat{\beta}(\text{ENet}) = \sqrt{1 + \lambda_2} \hat{\beta}^*.$$

So

$$\hat{\beta}(\text{ENet}) = (1 + \lambda_2) \hat{\beta}(\text{NENet})$$

possesses all properties of the Lasso.

# ELASTIC NET MODEL

---

## ENet estimator

$$\hat{\beta}(\text{ENet}) = \operatorname{argmin}_{\beta} \beta^T \left( \frac{X^T X + \lambda_2 I}{1 + \lambda_2} \right) \beta - 2y^T X \beta + \lambda_1 |\beta|_1.$$

$$\hat{\beta}(\text{lasso}) = \operatorname{argmin}_{\beta} \beta^T (X^T X) \beta - 2y^T X \beta + \lambda_1 |\beta|_1.$$

# ELASTIC NET MODEL

**Comparison** with leukemia classification using microarrays  
(training: 27+11; test: 34) (Zou and Hastie 2004)

Method	10-FdCV	Test err	genes
Golub	3/38	4/34	50
SVM	1/38	1/34	31
PenLogitReg	1/38	2/34	26
NSC(PAM)	2/38	2/34	21
ENet	3/38	0/34	45

# REFERENCES

---

- Bae, and Mallick, B. (2004) Gene selection using a two-level hierarchical Bayesian model. *Bioinfo.* **20**, 18: 3423-3430.
- Brown, P. Vannucci, M. and Fearn (1998) Multivariate Bayesian variable selection and prediction, *JRSSB*, **60**, 627-641.
- Brown, P. Vannucci, M. and Fearn (2002) Bayesian model averaging with selection of regressors, *JRSSB*, **64**, 519-536.
- Fu, WJ. (1998) Penalized regressions: the Bridge versus the Lasso. *J. Comp. Grap. Statist.* **7**, 3: 397-416.
- George, El. and McCulloch, R.E. (1993) Variable selection via Gibbs sampling. *J. Amer. Statist. Assoc.*, **88**, 881-889.
- Golub, T. et al. (1999) Molecular classification of cancer: class discovery and class prediction by gene expression monitoring, *Science*, **286**, 531-536.
- Kotz, S. Kozubowski, T.J. and Podgorski, K. (2000) *The Laplace Distribution And Generalizations*, Birkhauser, Boston.

# REFERENCES

---

- Land, S. and Friedman, J. (1996) Variable Fusion: a new method of adaptive signal regression, Technical report, Department of Statistics, Stanford University.
- Lee, KE. Sha, N. Dougherty, E. Vannucci, M. and Mallick, B. (2003) Gene selection: a Bayesian variable selection approach. *Bioinfo.*, **19**, 90-97.
- Tibshirani R, Saunders M, Rosset S, Zhu, J. (2005) Sparsity and smoothness via the fused lasso, *J. Roy. Statist. Soc. B.* 67: 91-108.
- Vannucci, M. Brown, P. Fearn, T. and (2001) Predictor selection for model averaging. In *Bayesian Methods with Applications to Science, Policy and Official Statistics*, (eds E.I. George and P. Nanopoulos), pp.553-562.
- Zou, H. and Hastie, T. (2004) Regularization and variable selection via the elastic net, *Technical report, Stanford University*.