



# LASSO算法与广义贝叶斯信息准则

高明

2018年4月

——毕业论文中期汇报



上海交通大学

SHANGHAI JIAO TONG UNIVERSITY

# LASSO (Least absolute shrinkage and selection operator) + normal distribution

Given a data set  $\{y_i, x_{i1}, \dots, x_{ip}\}_{i=1}^n$ , suppose  $y_i \sim N(\mathbf{x}_i^\top \beta, \sigma^2)$ , then the model takes the form

$$y_i = \beta_0 1 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i = \mathbf{x}_i^\top \beta + \varepsilon_i, \quad i = 1, \dots, n, \quad (1)$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

$$X = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix},$$

Often these  $n$  equations are stacked together and written in vector form as  $\mathbf{y} = X\beta + \varepsilon$ .



We can calculate the likelihood function

$$\sum_{i=1}^n f(y_i|\beta, \sigma^2) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \sigma^{-n} e^{-\frac{(y-X\beta)^T(y-X\beta)}{2\sigma^2}} \quad (2)$$

$$\text{MLE: } \hat{\beta} = (X^T X)^{-1} X^T y, \hat{\sigma}^2 = \frac{1}{n} (y - X\hat{\beta})^T (y - X\hat{\beta}).$$

When we look it up in the matrix space, it is the same to minimize (when fix  $\sigma^2$ )

$$l(\beta) = (y - X\beta)^T (y - X\beta) \quad (3)$$

OLS: same.

LASSO:

$$\tilde{\beta} = \operatorname{argmin} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 \right\} \quad \text{subject to } |\beta|_1 \leq t \quad (4)$$



# Group Lasso and Sparse Group Lasso

$$\tilde{\beta} = \operatorname{argmin}\left\{\frac{1}{2n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda |\beta|_1\right\} \triangleq \operatorname{argmin}\{Q(\beta) + \phi(\beta)\} \quad (5)$$

GRPL:

we divide  $\{1 \dots m\}$  into  $J$  group  $\{G_1, \dots, G_J\}$ .  $\#\{G_j\} \triangleq p_j$ .  $\beta_{G_j} = (\beta_k)_{k \in G_j}$

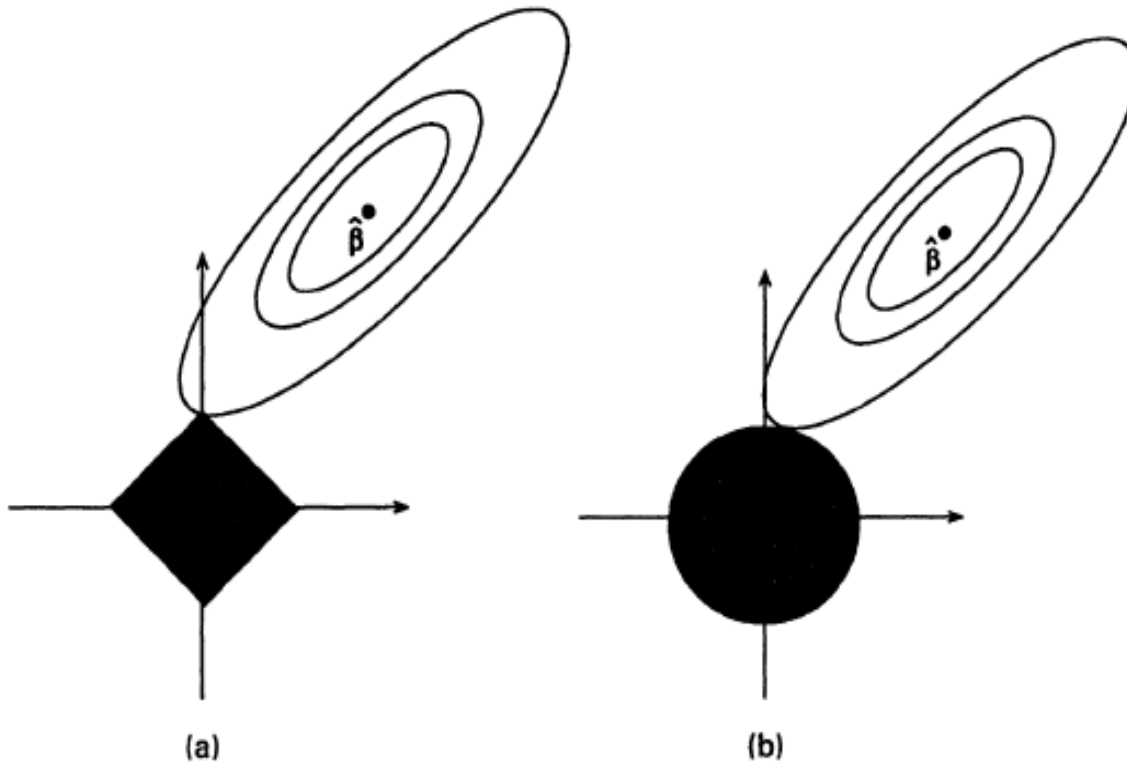
$$\phi(\beta) = \lambda \sum_{j=1}^J \sqrt{p_j} |\beta_{G_j}|_2 \quad (6)$$

SGRPL:

$$\phi(\beta) = \lambda \left\{ (1 - \alpha) \sum_{j=1}^J \sqrt{p_j} |\beta_{G_j}|_2 + \alpha |\beta|_1 \right\} \quad (7)$$



# Lasso vs Ridge



Estimation picture for (a) the lasso and (b) ridge regression

$$\tilde{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 \right\} \quad \text{subject to } |\beta|_1 \leq t$$



# Solve lasso with normal distribution

---

**Algorithm 1** Coordinate Descent

---

- 1: **repeat**
- 2:     Choose next index  $j$ .
- 3:     Update  $\beta_j$ .
- 4:

$$\beta_j = \arg \min_{\beta_j} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda |\beta|_1 \right\}$$

- 5: **until** stopping condition is met.
- 

If we standardize the design matrix  $X$ , we may get an explicit solution.

$$\beta_j = s[\hat{\beta}_j, \lambda] \triangleq \text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda)_+ \quad (8)$$

$s$  is called soft-thresholding operator.  $\hat{\beta}_j$  is the OLS of residual, fixing  $\beta_k$ ,  $k \neq j$ .

Also, check whether the gradient=0 fall in to the non-differentiable point.  
Subgradient.

It can be shown that  $\beta \rightarrow \tilde{\beta}$





# Solve grplasso with normal distribution

---

**Algorithm 2** Block Coordinate Descent

---

- 1: repeat
- 2:   Choose next block index  $G_j$ .
- 3:   Update  $\beta_{G_j}$ .
- 4:

$$\beta_{G_j} = \arg \min_{\beta_{G_j}} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^J \sqrt{p_j} |\beta_{G_j}|_2 \right\}$$

- 5: until stopping condition is met.
- 

$$\beta_{G_j} = s[\hat{\beta}_{G_j}, \lambda \sqrt{p_j}] \quad (9)$$



# Solve sgrplasso with normal distribution

---

## Algorithm 3 Block Coordinate Descent

---

- 1: repeat
- 2:     Choose next block index  $G_j$ .
- 3:     repeat
- 4:         Choose next index  $k$  in  $G_j$ .
- 5:         Update  $\beta_k$ .
- 6:

$$\beta_k = \arg \min_{\beta_k} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda(1 - \alpha) \sum_{j=1}^J \sqrt{p_j} |\beta_{G_j}|_2 + \lambda\alpha |\beta|_1 \right\}$$

- 7:     until stopping condition is met.
  - 8: until stopping condition is met.
- 

the coordinate-wise algorithms for the lasso, the grouped lasso and elastic net, etc. converge to their optimal solutions.





# Lasso with binomial distribution

Suppose  $y_i \sim \text{Bin}(1, p_i)$ , then  $\mu_i \triangleq E(y_i) = p_i$ . Let  $x_i^T \beta = g(\mu_i) \triangleq \ln \frac{p_i}{1-p_i}$ .

Hence,

$$p_i = P(y_i = 1|x_i) = \frac{1}{1 + e^{-x_i^T \beta}} \quad (10)$$

$$\begin{aligned} Q(\beta) &= \frac{1}{n} \sum_{i=1}^n y_i \ln p_i + (1 - y_i) \ln(1 - p_i) \\ &= \frac{1}{n} \sum_{i=1}^n y_i (x_i^T \beta) - \ln(1 + e^{x_i^T \beta}) \end{aligned}$$



---

#### Algorithm 4 Newton Method

---

- 1: repeat
  - 2:    Evaluate  $g = \nabla Q(\beta) = X^T(\mu - y)$
  - 3:    Evaluate  $H = \nabla g = X^T S X$ , here  $S = \text{diag}\{\mu_1(1 - \mu_1), \dots, \mu_n(1 - \mu_n)\}$
  - 4:    Solve  $d = -H^{-1}g$
  - 5:    Update  $\beta$ :  $\beta = \beta + d$
  - 6:     $\beta^{t+1} = (X^T S^t X)^{-1} X^T S^t z^t$ , here  $z^t = X\beta^t + S^{-1t}(y - \mu^t)$
  - 7: until stopping condition is met.
- 

---

#### Algorithm 5 Iteratively Reweighted Least Squares

---

- 1: repeat at  $t^{th}$  iterate
  - 2:  
$$Q(\beta^t) = -\frac{1}{2n} \sum_{i=1}^n s_i^t (z_i^t - x_i^T \beta^t)^2 + \text{CONSTANT}$$
  - 3:    Update  $\beta^{t+1} = \text{argmin} Q(\beta^t)$
  - 4:    From OLS,  $\beta^{t+1} = (X^T S^t X)^{-1} X^T S^t z^t$
  - 5: until stopping condition is met.
-



# Solve Lasso with binomial distribution

---

## Algorithm 6 IWLS + Coordinate Descent

---

- 1: **repeat** at  $t^{th}$  iterate
  - 2:     Update the quadratic approximation  $Q(\beta^t)$  at  $\beta^t$
  - 3:     **repeat**
  - 4:         Choose next index  $j$
  - 5:         Update  $\beta_j^t = \arg \min_{\beta_j^t} \{Q(\beta^t) + \lambda|\beta_j^t|\}$ .
  - 6:     **until** stopping condition is met.
  - 7: **until** stopping condition is met.
-



# Solve grpLasso with binomial distribution

---

## Algorithm 7 Block Coordinate Gradient Descent

---

- 1: repeat at  $t^{th}$  iterate
  - 2:     Update the quadratic approximation  $Q(\beta^t)$  at  $\beta^t$
  - 3:     repeat
  - 4:         Choose next block index  $G_j$
  - 5:         Approximate  $H_{G_j G_j}^t = h_{G_j}^t I$
  - 6:         Find direction  $d_{G_j}^t$
  - 7:         
$$d_{G_j}^t = \arg \min_{d_{G_j}^t} \{Q(\beta^t) + \lambda \sum_{j=1}^J \sqrt{p_{G_j}^t} |\beta_{G_j}^t|_2\}$$
  - 8:         Line search  $\tau$  using the Armijo rule.
  - 9:         Update  $\beta_{G_j}^t = \beta_{G_j}^t + \tau d_{G_j}^t$
  - 10:     until stopping condition is met.
  - 11: until stopping condition is met.
-



# Solve sgrpLasso with binomial distribution

---

## Algorithm 8 Block Coordinate Gradient Descent

---

- 1: repeat at  $t^{th}$  iterate
  - 2:     Update the quadratic approximation  $Q(\beta^t)$  at  $\beta^t$
  - 3:     repeat
  - 4:         Choose next block index  $G_j$
  - 5:         Approximate  $H_{G_j G_j}^t = h_{G_j}^t I$
  - 6:         repeat
  - 7:             Choose next index  $k$  in  $G_j$
  - 8:             Find direction  $d_k^t$
  - 9:             
$$d_k^t = \arg \min_{d_k^t} \{Q(\beta^t) + \lambda(1 - \alpha) \sum_{j=1}^J \sqrt{p_{G_j}^t} |\beta_{G_j}^t|_2 + \lambda\alpha |\beta^t|_1\}$$
  - 10:             Line search  $\tau$  using the Armijo rule.
  - 11:             Update  $\beta_k^t = \beta_k^t + \tau d_k^t$
  - 12:         until stopping condition is met.
  - 13:     until stopping condition is met.
  - 14: until stopping condition is met.
-



# Model selection: CV

In K-fold cross-validation, the original sample is randomly partitioned into K equal sized subsamples. Of the K subsamples, a single subsample is retained as the validation data for testing the model, and the remaining  $K - 1$  subsamples are used as training data.

for each  $k = 1, \dots, K$ , fit the model with parameter  $\lambda$  to the other  $K - 1$  parts, giving  $\tilde{\beta}^{-k}(\lambda)$  and compute its loss  $LOSS_k(\lambda)$  in predicting the  $k^{th}$  part. This gives the cross-validation error

$$CV(\lambda) = \frac{1}{K} \sum_{k=1}^K LOSS_k(\lambda) \quad (11)$$

$$\lambda^* = \operatorname{argmin} CV(\lambda) \quad (12)$$



# Model selection: IC

## 7 AIC

$$AIC(\lambda) = -2\ln l(\hat{\beta}(\lambda)) + 2\nu(\lambda) \quad (13)$$

## 8 BIC

$$BIC(\lambda) = -2\ln l(\hat{\beta}(\lambda)) + \nu(\lambda)\ln n \quad (14)$$

Here,  $\nu(\lambda) = df(\lambda)$

## 9 EBIC

$$EBIC_{\gamma}(\lambda) = -2\ln l(\hat{\beta}(\lambda)) + \nu(\lambda)\ln n + 2\gamma\nu(\lambda)\ln p \quad (15)$$

**Theorem 1** *Suppose  $\lambda_0$  is the true model. Under some mild conditions with  $n \rightarrow \infty$ , we have*

$$P\{\min EBIC_{\gamma}(\lambda) \leq EBIC_{\gamma}(\lambda_0)\} \rightarrow 0 \quad (16)$$



谢谢！

