

# Sparse Modeling: Best Subset and the Lasso

Statistical Learning

CLAMSES - University of Milano-Bicocca

Aldo Solari

# References

- Tibshirani, Wasserman (2017). Sparsity, the Lasso, and Friends.  
Lecture notes on Statistical Machine Learning

## Three norms: $\ell_0$ , $\ell_1$ and $\ell_2$

- Let's consider three canonical choices: the  $\ell_0$ ,  $\ell_1$  and  $\ell_2$  norms:

$$\|\beta\|_0 = \sum_{j=1}^p \mathbb{1}\{\beta_j \neq 0\}, \quad \|\beta\|_1 = \sum_{j=1}^p |\beta_j|, \quad \|\beta\|_2 = \sqrt{\sum_{j=1}^p \beta_j^2}$$

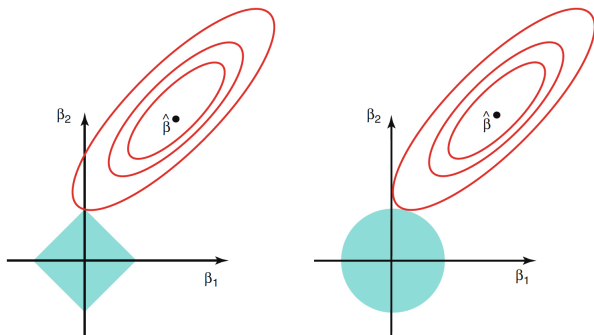
- $\ell_0$  is not a proper norm: it does not satisfy positive homogeneity, i.e.  $\|a\beta\|_0 \neq |a|\|\beta\|_0$  for  $a \in \mathbb{R}$

# Constrained form

$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 \text{ subject to } \|\beta\|_0 \leq c \quad \text{Best Subset Selection}$$

$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 \text{ subject to } \|\beta\|_1 \leq c \quad \text{Lasso Regression}$$

$$\min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 \text{ subject to } \|\beta\|_2^2 \leq c \quad \text{Ridge Regression}$$



The “classic” illustration comparing lasso and ridge constraints.  
From Chapter 3 of ESL

# Sparsity

- *Signal sparsity* is the assumption that only a small number of predictors have an effect, i.e. have  $\beta_j \neq 0$
- In this case we would like our estimator  $\hat{\beta}$  to be sparse, meaning that  $\hat{\beta}_j = 0$  for many components  $j \in \{1, \dots, p\}$
- Sparse estimators are desirable because perform variable selection and improve interpretability of the result
- The best subset selection and the lasso estimators are sparse, the ridge estimator is not sparse

## Penalized form

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_0 \quad \text{Best Subset Selection}$$

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \quad \text{Lasso Regression}$$

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \quad \text{Ridge Regression}$$

- Suppose that  $y \sim N(\mu, 1)$
- $\ell_0$  penalty

$$\min_{\mu} \frac{1}{2}(y - \mu)^2 + \lambda \mathbb{1}\{\mu \neq 0\}, \quad \hat{\mu} = H_{\sqrt{2\lambda}}(y)$$

where  $H_a(y) = y \mathbb{1}\{|y| > a\}$  is the hard-thresholding operator

- $\ell_1$  penalty

$$\min_{\mu} \frac{1}{2}(y - \mu)^2 + \lambda |\mu|, \quad \hat{\mu} = S_{\lambda}(y)$$

where

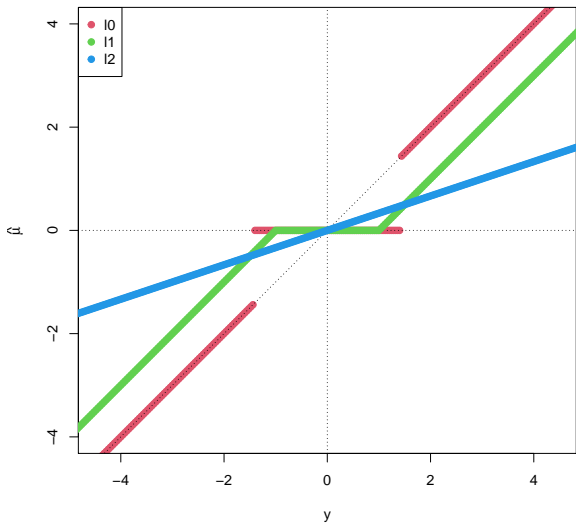
$$S_a(y) = \begin{cases} y - a & \text{if } y > a \\ 0 & \text{if } -a \leq y \leq a \\ y + a & \text{if } y < -a \end{cases}$$

is the soft-thresholding operator

- $\ell_2$  penalty

$$\min_{\mu} \frac{1}{2}(y - \mu)^2 + \lambda \mu^2, \quad \hat{\mu} = \left(\frac{1}{1 + 2\lambda}\right)y$$





$$\lambda = 1$$

# Hard and soft thresholding

- $\ell_0$  penalty creates a zone of sparsity but it is discontinuous (hard thresholding)
- $\ell_1$  penalty creates a zone of sparsity but it is continuous (soft thresholding)
- $\ell_2$  penalty creates a nice smooth estimator but it is never sparse

## Orthogonal case

- Suppose  $X^t X = I_p$

- OLS estimator

$$\hat{\beta} = X^t y$$

- BSS estimator

$$\hat{\beta} = H_{\sqrt{2\lambda}}(X^t y)$$

- Lasso estimator

$$\hat{\beta} = S_{\lambda}(X^t y)$$

- Ridge estimator

$$\hat{\beta} = \left(\frac{1}{1 + 2\lambda}\right) X^t y$$

where  $H_a(\cdot)$ ,  $S_a(\cdot)$  are the componentwise hard- and soft-thresholding operators

$$\begin{aligned}g(b) &= \frac{1}{2}\|y - Xb\|^2 + \lambda\|b\|_1 \\&= \frac{1}{2}(y^t y + b^t X^t X b - 2b^t X^t y) + \lambda\|b\|_1 \\&= \frac{1}{2}y^t y + \frac{1}{2}\sum_{j=1}^p \{b_j^2 - 2b_j X_j^t y + 2\lambda|b_j|\} \\&= \frac{1}{2}y^t y + \frac{1}{2}\sum_{j=1}^p f_j(b_j)\end{aligned}$$

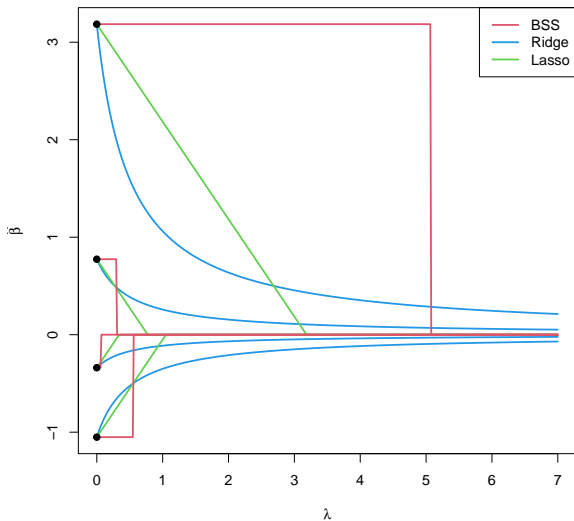
We can minimize each quantity  $f_j$  inside the sum independently.

Suppose  $b_j \geq 0$  and remove the absolute value:

$$f_j(b_j) = b_j^2 + 2b_j(\lambda - X_j^t y)$$

To minimize this, take the derivative and set it equal to zero:

$$b_j = X_j^t y - \lambda$$



Solution paths of  $\ell_0$ ,  $\ell_1$  and  $\ell_2$  penalties as a function of  $\lambda$

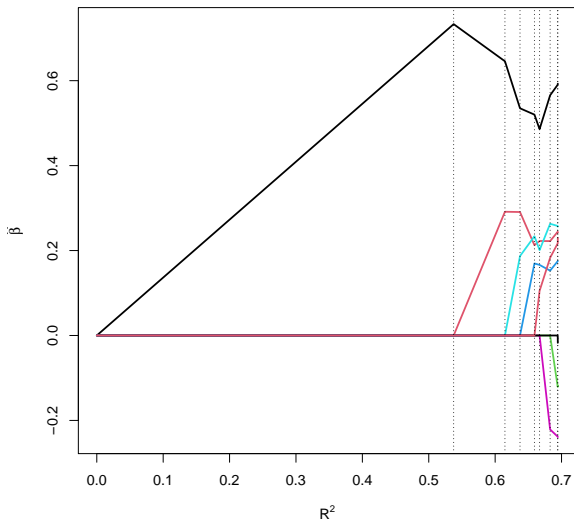
# Convexity

- Consider using the norm  $\|\beta\|_q = (\sum_{j=1}^q |\beta_j|^q)^{1/q}$  as a penalty. Sparsity requires  $q \leq 1$  and convexity requires  $q \geq 1$ . The only norm that gives sparsity and convexity is  $q = 1$
- The lasso and ridge regression are *convex optimization problems*, best subset selection is not
- The ridge regression optimization problem is always *strictly convex* for  $\lambda > 0$
- The best subset selection optimization problem is N-P-complete because of its combinatorial complexity (there are  $2^p$  subsets), the worst kind of non convex problem

# Forward Stepwise Selection

Greedy forward algorithm, sub-optimal but feasible alternative to BSS and applicable when  $p > n$

- Set  $S_0$  as the null model (intercept only)
- For  $k = 0, \dots, \min(n - 1, p - 1)$ :
  1. Consider all  $p - k$  models that augment the predictors in  $S_k$  with one additional predictor
  2. Choose the best among these  $p - k$  models and call it  $S_{k+1}$ , where best is defined having the smallest RSS
- Select a single best model from among  $S_0, S_1, S_2, \dots$  (e.g. using  $C_p$ , BIC, Cross-Validation, validation set, etc.)



Forward Stepwise solution path as a function of training  $R^2$



## The Lasso

The name “lasso” was also introduced as an acronym for *Least Absolute Selection and Shrinkage Operator* (Tibshirani, 1996)

The lasso finds the solution  $(\hat{\alpha}, \hat{\beta})$  to the optimization problem

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y - 1\alpha - X\beta\|_2^2 + \lambda \|\beta\|_1$$

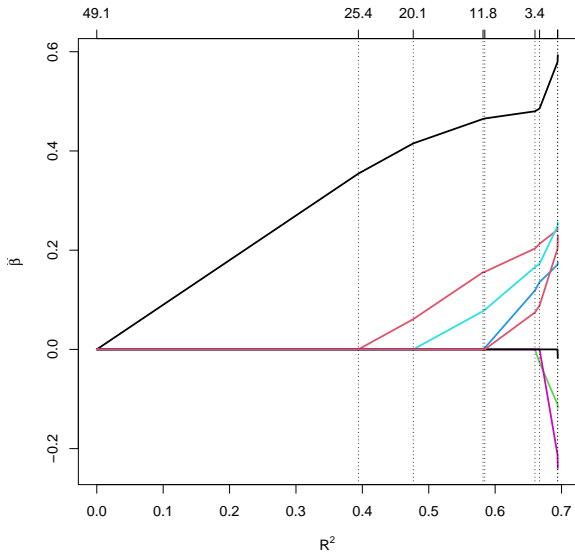
- Typically, we first standardize the predictors  $X$  so that each column is centered ( $(1/n) \sum_{i=1}^n x_{ij} = 0$ ) and has unit variance ( $(1/n) \sum_{i=1}^n x_{ij}^2 = 1$ )
- Without standardization, the lasso solutions would depend on the units (e.g., feet versus meters) used to measure the predictors. On the other hand, we typically would not standardize if the features were measured in the same units
- For convenience, we also assume that the outcome values  $y_i$  have been centered ( $(1/n) \sum_{i=1}^n y_i = 0$ ). Centering is convenient, since we can omit the intercept term  $\alpha$  in the lasso optimization, and given the solution  $\hat{\beta}$

$$\hat{\alpha} = \bar{y} - \sum_{j=1}^p \bar{x}_j \hat{\beta}_j$$

- Lagrange form

$$\frac{1}{2n} \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

- Intercept term omitted (center / scale the columns of  $X$  and  $y$ )
- The coefficient profiles for the lasso are continuous and piecewise linear over the range of  $\lambda$ , with knots occurring whenever the *active set* changes, or the sign of the coefficients changes



Lasso solution path as a function of training  $R^2$

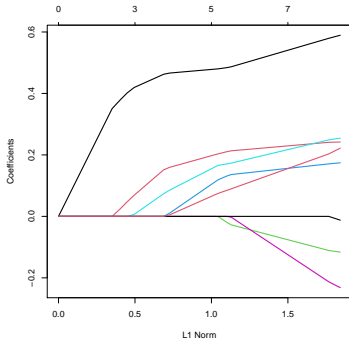
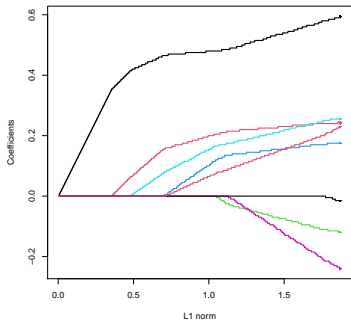
# Boosting with componentwise linear least squares

- Response and predictors are standardized to have mean zero and unit norm
- Initialize  $\hat{\beta}^{(0)} = 0$
- For  $b = 1, \dots, B$ 
  - compute the residuals  $r = y - X\hat{\beta}^{(b-1)}$
  - find the predictor  $X_j$  most correlated with the residuals  $r$
  - update  $\hat{\beta}^{(b-1)}$  to  $\hat{\beta}^{(b)}$  with

$$\hat{\beta}_j^{(b)} = \hat{\beta}_j^{(b-1)} + \epsilon \cdot s_j$$

where  $s_j$  is the sign of the correlation

- This is known as *forward stagewise regression* and converges to the least squares solution when  $n > p$
- Forward stagewise regression with infinitesimally small step-sizes, i.e.  $\epsilon \rightarrow 0$ , produces a set of solutions which is approximately equivalent to the set of Lasso solutions



Left: forward stagewise regression with  $\epsilon = 0.005$ ; Right: lasso

## Convex optimization and the elastic net



# Convex function

The objective function of the  $\ell_1$  penalty, unlike the  $\ell_0$  penalty, is a continuous function in the regression vector. Not only is the objective function continuous, it is also convex.

We say that a function  $f: \mathbb{R}^p \mapsto \mathbb{R}$  *convex* if for any values  $b_1$  and  $b_2$  and quantity  $t \in [0, 1]$  we have

$$f(tb_1 + (1 - t)b_2) \leq tf(b_1) + (1 - t)f(b_2)$$

Replacing the  $\leq$  with  $<$  for  $t \in (0, 1)$  yields a definition of *strict* convexity.

The  $\ell_1$ -penalized objective function is in fact convex (but not strictly convex, which makes the solutions non-linear in the  $y_i$ , and there is no closed form solution as in ridge).

## Convex optimization

A convex function does not have any local minima that are not also global minima. In other words, if the value  $b_0$  minimizes  $f$  over a neighborhood of  $b_0$ , it must also minimize  $f$  over its entire domain.

To see this, assume that  $b_1$  is any point that is not a global minima but set  $b_2$  equal to a global minima of  $f$ . Then, for any  $t \in [0, 1)$ , we have

$$tf(b_1) + (1 - t)f(b_2) < tf(b_1) + (1 - t)f(b_1) = f(b_1)$$

and by convexity this implies that

$$f(tb_1 + (1 - t)b_2) < f(b_1)$$

For any neighborhood around  $b_1$  we can find  $t$  close enough to 1 such that  $tb_1 + (1 - t)b_2$  is in that neighborhood, and therefore  $b_1$  cannot be a local minimum.

The lack of local optima makes it possible to optimize convex objective functions using e.g. the *coordinate descent algorithm*.

## Elastic net

Define the objective function  $f$  for some  $\lambda > 0$  and  $\alpha \in [0, 1]$  as

$$f(\beta; \lambda, \alpha) = \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \left( (1 - \alpha) \frac{1}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$

and the corresponding *elastic net* estimator as

$$\hat{\beta}_{\lambda, \alpha} = \arg \min_{\beta} f(\beta; \lambda, \alpha)$$

Setting  $\alpha$  to 1 yields the Lasso regression and setting it to 0 the ridge regression.

Adding a small  $\ell_2$ -penalty preserves the variable selection and convexity properties of the  $\ell_1$ -penalized regression, while reducing the variance of the model when  $X$  contains sets of highly correlated variables.

# Coordinate descent

Coordinate descent is a general purpose convex optimization algorithm particularly well-suited to solving the elastic net equation.

Coordinate descent successively minimizes the objective function along each variable. In every step all but one variable is held constant and a value for the variable of interest is chosen to minimize the constrained problem.

This process is applied iteratively over all of the variables until the algorithm converges.

Writing the problem in terms of the individual values of  $b$ :

$$f(b) = \frac{1}{2n} \sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} b_j)^2 + \lambda \sum_{j=1}^p \left\{ \frac{1}{2} (1 - \alpha) b_j^2 + \alpha |b_j| \right\}$$

Let  $\tilde{b}$  be a vector of candidate values of the regression vector  $b$  and assume  $\tilde{b}_l > 0$ . Then the derivative of this function with respect to  $b_l$  at  $b = \tilde{b}$  is

$$\begin{aligned} \frac{\partial f}{\partial b_l} \Big|_{b=\tilde{b}} &= -\frac{1}{n} \sum_{i=1}^n x_{il} (y_i - \sum_{j \neq l}^p x_{ij} \tilde{b}_j - x_{il} \tilde{b}_l) + \lambda ((1 - \alpha) \tilde{b}_l + \alpha) \\ &= -\frac{1}{n} \sum_{i=1}^n x_{il} (y_i - \tilde{y}_i^{(l)}) + \frac{1}{n} \sum_{i=1}^n x_{il}^2 \tilde{b}_l + \lambda (1 - \alpha) \tilde{b}_l + \lambda \alpha \end{aligned}$$

where  $\tilde{y}_i^{(l)} = \sum_{j \neq l}^p x_{ij} \tilde{b}_j$  is the contribution of all regressors in the model except the  $l$ th.

By setting the function to zero and solving for  $\tilde{b}_l$  resulting in

$$\tilde{b}_l = \frac{\frac{1}{n} \sum_{i=1}^n x_{il}(y_i - \tilde{y}_i^{(l)}) - \lambda\alpha}{\frac{1}{n} \sum_{i=1}^n x_{il}^2 + \lambda(1 - \alpha)}$$

We can then generalize this equation using the soft-thresholding function as

$$\tilde{b}_l \leftarrow \frac{S_{\lambda\alpha}(\frac{1}{n} \sum_{i=1}^n x_{il}(y_i - \tilde{y}_i^{(l)}))}{\frac{1}{n} \sum_{i=1}^n x_{il}^2 + \lambda(1 - \alpha)}$$

The algorithm updates each of the values of  $\tilde{b}$  based on the initial guess of a  $p \times 1$  vector of zeros.

## KKT conditions

The solution must satisfy the subgradient / Karush-Kuhn-Tucker conditions

$$-\frac{1}{n} \langle X_j, y - X\hat{b} \rangle + \lambda s_j = 0$$
$$\sum_{i=1}^n x_{ij} (y_i - \sum_{j=1}^p x_{ij} \hat{b}_j) = \lambda s_j$$

for  $j = 1, \dots, p$ , where

$$s_j \in \begin{cases} 1 & \text{if } \hat{b}_j > 0 \\ [-1, 1] & \text{if } \hat{b}_j = 0 \\ -1 & \text{if } \hat{b}_j < 0 \end{cases}$$

which means that if these conditions have not been met, then our  $\hat{b}$  vector cannot be optimal.

Cross-validation



# Maximum $\lambda$

What values of  $\lambda$  should we use when performing cross-validation?  
Consider the case that  $\alpha = 1$ . From the coordinate descent updates

$$\tilde{b}_j = \frac{1}{n} X_j^t y - \lambda$$

It follows that if  $|X_j^t y| \leq n\lambda$  then  $\tilde{b}_j = 0$ . The smallest value of  $\lambda$  for which all  $\tilde{b}_j$  are zero is therefore:

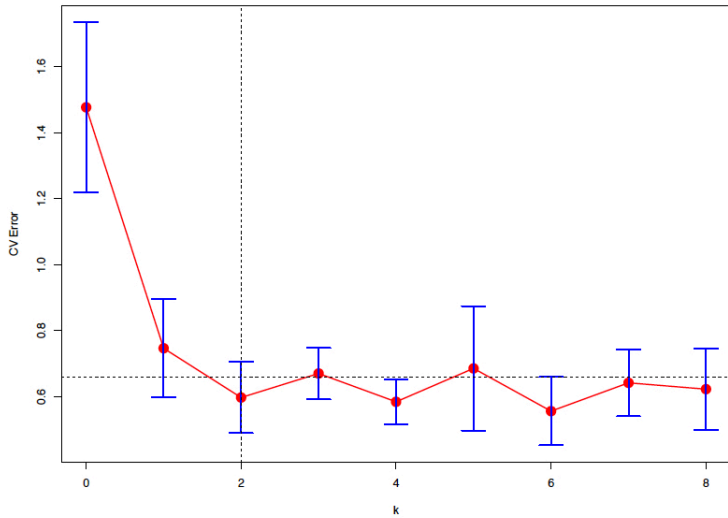
$$\lambda_{max} = \max_{j \in \{1, \dots, p\}} \left| \frac{X_j^t y}{n} \right|$$

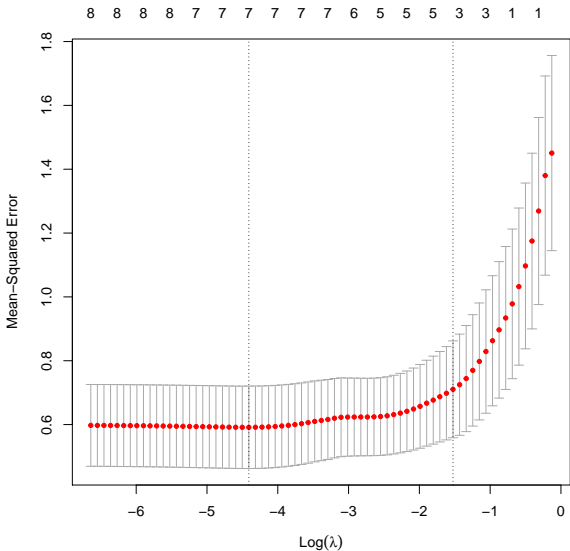
## Minimum CV error and the 1se rule

- `lambda.min`:  $\lambda$  that minimize the cross-validation error
- `lambda.1se`: largest value of `lambda` such that error is within 1 standard error of the minimum (*one standard error rule*). To compute cross-validation "standard errors"

$$se = \frac{1}{\sqrt{K}} \text{sd}(\text{Err}^{-1}, \dots, \text{Err}^{-K})$$

where  $\text{Err}^{-k}$  denotes the error incurred in predicting the observations in the  $k$  hold-out fold,  $k = 1, \dots, K$ .





$$\lambda_{\min} = 0.012 \text{ (7 nonzero)}, \lambda_{\text{lse}} = 0.21 \text{ (3 nonzero)}$$

# Degrees of freedom

- Let  $A(\lambda) = \{j \in \{1, \dots, p\} : \hat{\beta}_j(\lambda) \neq 0\}$  denotes the active set
- The degrees of freedom of the Lasso are the

$$\text{df}(\lambda) = |A(\lambda)|$$

i.e. the size of the active set

# Bayesian interpretation

- A Bayesian viewpoint assumes that  $\beta$  has a double-exponential (Laplace) prior distribution with mean zero and scale parameter a function of  $\lambda$

$$(1/2\tau) \exp(-\|\beta\|_1/\tau)$$

with  $\tau = 1/\lambda$

- It follows that the posterior mode for  $\beta$  is the lasso solution
- However, the lasso solution is not the posterior mean and, in fact, the posterior mean does not yield a sparse coefficient vector

Extensions of the lasso

# Group Lasso

- Suppose we have a partition  $G_1, \dots, G_q$  of  $\{1, \dots, p\}$
- The group Lasso penalty (Yuan and Lin, 2006) is given by

$$\lambda \sum_{k=1}^q m_k \|\beta_{G_k}\|_2$$

The multipliers  $m_k > 0$  serve to balance cases where the groups are of very different sizes; typically we choose  $m_k = \sqrt{|G_k|}$

- This penalty encourages either an entire group  $G$  to have  $\hat{\beta}_G = 0$  or  $\hat{\beta}_j \neq 0$  for all  $j \in G$
- Such a property is useful when groups occur through coding for categorical predictors or when expanding predictors using basis functions.

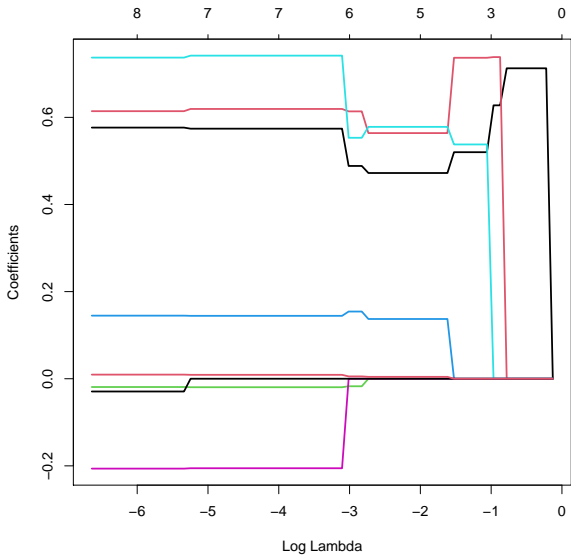


# Relaxed Lasso

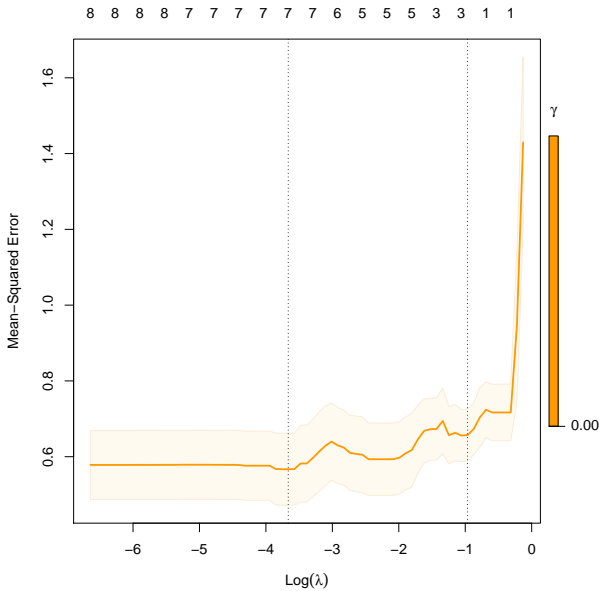
- Originally proposed by Meinshausen (2006). We present a simplified version.
- Suppose  $\hat{\beta}_\lambda$  is the lasso solution at  $\lambda$  and let  $\hat{A}$  be the active set of indices with nonzero coefficients in  $\hat{\beta}_\lambda$
- Let  $\hat{\beta}^{\text{LS}}$  be the coefficients in the least squares fit, using only the variables in  $\hat{A}$ . Let  $\hat{\beta}_\lambda^{\text{LS}}$  be the full-sized version of this coefficient vector, padded with zeros.  $\hat{\beta}_\lambda^{\text{LS}}$  debiases the lasso, while maintaining its sparsity.
- Define the Relaxed Lasso

$$\hat{\beta}_\lambda^{\text{RELAX}} = \gamma \hat{\beta}_\lambda + (1 - \gamma) \hat{\beta}_\lambda^{\text{LS}}$$

with  $\gamma \in [0, 1]$  is an additional tuning parameter which can be selected by cross-validation



$$\gamma = 0$$



$$\gamma = 0$$