

Regularization Technique for Linear Regression in High Dimension: Ridge Regression and LASSO



Xiang Zhou

School of Data Science
Department of Mathematics
City University of Hong Kong

Focus point of this part is still the bias-variance tradeoff by introducing the shrinkage methodology to allow the **biased** estimator. The linear regression problem here is used as a second example ¹of **model assessment and model selection**. (Chapter 7 [ESL])

Not intention to cover the followings: the optimization theory and numerical methods for the ridge regression and LASSO, which is very important in practice but more close to the optimization field.

¹the first is KNN

Revisit the Bias-Variance tradeoff

For a model fit $\hat{f}_D(x)$ based on the data $D = (\mathbf{X}, \mathbf{y})$, where \mathbf{X} is the design matrix and $\mathbf{y} = \mathbf{X}\beta + \epsilon$ is the response data, a good measure of the quality of this model at a new test input $x_0 \in \mathbb{R}^p$ is the mean square error (MSE). Let $f(x) = x^\top \beta$ be the true value of the output at the point x , then

$$\begin{aligned}\text{MSE}(\hat{f}_D(x_0)) &:= \mathbb{E}_D(\hat{f}_D(x_0) - f(x_0))^2 \\ &= \text{Var}_D(\hat{f}_D(x_0)) + \left(\mathbb{E}_D \hat{f}_D(x_0) - f(x_0)\right)^2.\end{aligned}$$

- Typically, when bias is low, variance will be high and vice-versa. Choosing estimators often involves a tradeoff between bias and variance.
- So far, OLS estimator $\hat{\beta}^{OLS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ is unbiased and also has the minimal MSE if restricted to be *unbiased* and linear in \mathbf{y} .

For the OLS fit, its MSE (written with dependency on x_0) is

$$\begin{aligned}\text{MSE}^{OLS}(x_0) &= \text{Var}_{\mathbf{y}}(x_0^\top \hat{\beta}^{OLS}) \\ &= \text{Var}_{\mathbf{y}}(x_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} x_0 \mathbb{V}(\mathbf{y}) x_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X} \\ &= \sigma_\varepsilon^2 \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} x_0 x_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ &= \sigma_\varepsilon^2 \left\| \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} x_0 \right\|_2^2\end{aligned}$$

Here \mathbf{X} is treated as deterministic or say that we work on the condition variance of given \mathbf{X} .

Exercise

In the above, $x_0 \in \mathcal{X}$ is arbitrary. Now let x_0 take to be each of training sample value x_i , $1 \leq i \leq N$, show that the in-sample error, which is defined by $\frac{1}{N} \sum_{i=1}^N \text{MSE}^{OLS}(x_i)$, equals to $\frac{p}{N} \sigma_\varepsilon^2$ where p is the dimension of \mathcal{X} . (Equation (7.29) in [ESL])

There can be biased estimators with smaller MSE. The following property is quite general.

Exercise

Define a biased estimate of the coefficient β in the following special form

$$\tilde{\beta} = (1 + \alpha)\hat{\beta}$$

with a scalar α where $\hat{\beta}$ is an unbiased estimate. Calculate the MSE of $\tilde{\beta}$ and find a condition that $MSE(\tilde{\beta}) < MSE(\hat{\beta})$.

- Generally, by regularizing the estimator in some way, its variance will be reduced; if the corresponding increase in bias is small, this will be worthwhile.
- Examples of regularization: subset selection (forward, backward, all subsets)¹; ridge regression, lasso

¹read [ISL][ESL] by yourself; not to cover in lecturing
Xiang Zhou CityU

Ridge Regression

Regression in high dimension

- High dimensional problem: the input dimension $d = p + 1$ is close or greater than the number of samples n .¹
- When the design matrix \mathbf{X} is high-dimensional, the covariates (the columns of \mathbf{X}) are super-collinear. *collinearity* in regression analysis refers to the event of two (or multiple) covariates being highly linearly related.
- In OLS estimate, the dimension of the subspace \mathbf{X} may be less than d and then the matrix inverse $(\mathbf{X}^T \mathbf{X})^{-1}$ does not exist or $\mathbf{X}^T \mathbf{X}$ is close to singular even if invertible².

¹we switch symbols both d and $p + 1$ from now. The notations for design matrix and response variable are changed to bold font.

²Numerical linear algebra uses the **condition number** (defined as the ratio of largest to the smallest eigenvalue) to represent the illness of the problem: the larger, the worse.

Solution of ridge regression

Ridge Regression solves

$$\min_{\beta} \quad \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2$$

The solution is

$$\hat{\beta}^\lambda = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Inclusion of λ makes problem non-singular even if $\mathbf{X}^\top \mathbf{X}$ is not invertible:
This was the original motivation for ridge regression (Hoerl and Kennard, 1970)

Note $\lambda = 0$ gives the ordinary least squares estimator, and if $\lambda \rightarrow \infty$ then $\hat{\beta}_\lambda \rightarrow 0$. In general, with a good choice of λ , $\hat{\beta}_\lambda$ is a biased estimator that may have smaller mean squared error than the least squares estimator

Eigenvalue shrinkage in ridge regression

Let $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ is a singular value decomposition of \mathbf{X} : \mathbf{D} is the $(p+1) \times (p+1)$ diagonal matrix consisting of singular values $d_0 \geq d_1 \geq \dots \geq d_p$. \mathbf{U} and \mathbf{V} are $n \times (p+1)$ and $(p+1) \times (p+1)$ matrices, respectively ¹. Then $\mathbf{X}^\top \mathbf{X} = \mathbf{V}\mathbf{D}^2\mathbf{V}^\top$ and the ridge regression solution

$$\begin{aligned}\hat{\beta}^\lambda &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}, \\ &= (\mathbf{V}\mathbf{D}^2\mathbf{V}^\top + \lambda \mathbf{V}\mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{D}\mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}(\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D}\mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V} \operatorname{diag} \left\{ \frac{d_j}{d_j^2 + \lambda} \right\} \mathbf{U}^\top \mathbf{y},\end{aligned}$$

which is well defined for any d_j when λ is strictly positive.

¹The column space of \mathbf{U} is the column space of \mathbf{X} in \mathbb{R}^n and the column space of \mathbf{V} is the row space of \mathbf{X} in \mathbb{R}^{p+1} . $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_n$ and $\mathbf{V}\mathbf{V}^\top = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_{p+1}$. The number of nonzero $\{d_j\}$ is the rank of \mathbf{X} .

Let $\tilde{\beta} = \mathbf{V}^T \hat{\beta}$ and $\tilde{\mathbf{y}} = \mathbf{U}^T \mathbf{y}$. Then

$$\tilde{\beta}^\lambda = \mathbf{V}^T \hat{\beta}^\lambda = \text{diag} \left\{ \frac{d_j}{d_j^2 + \lambda} \right\} \tilde{\mathbf{y}}$$

If \mathbf{D} is nonsingular, then $\tilde{\beta}^{OLS} = \tilde{\beta}^0 = \mathbf{D}^{-1} \tilde{\mathbf{y}}$ exists. So

$$\tilde{\beta}^\lambda = \mathbf{D}_\lambda \mathbf{D}^{-1} \tilde{\mathbf{y}} = \mathbf{D}_\lambda \tilde{\beta}^{OLS} \quad \text{where} \quad \mathbf{D}_\lambda := \text{diag} \left\{ \frac{d_j^2}{d_j^2 + \lambda} \right\} \leq \mathbf{I}$$

$$\tilde{\beta}_j^{ridge} = \frac{d_j^2}{d_j^2 + \lambda} \tilde{\beta}_j^{OLS}$$

Exercise

Find the bias, the variance and the MSE for the transformed ridge coefficient $\tilde{\beta}^\lambda$ in terms of $\beta, \mathbf{X}, \mathbf{y}, \sigma^2$. Find the optimal λ in theory that minimize the MSE.

Smoother matrix and effective degree of freedom

- A smoother matrix S is a linear operator satisfying

$$\mathbf{y} = \mathbf{S}\mathbf{y}$$

where S may depend on \mathbf{X} , but not \mathbf{y} .

- define the **effective degrees of freedom** (or effective number of parameters) for a smoother S :

$$\text{df}(\mathbf{S}) = \text{Trace}(\mathbf{S})$$

Exercise

Ex. 7.4 and 7.5 in [ESL]

Effective degree of freedom

The ridge solution can be rewritten as

$$\begin{aligned}\hat{\mathbf{y}}^{ridge} &= \mathbf{X}\hat{\boldsymbol{\beta}}^{ridge} = \mathbf{UDV}^T \mathbf{VD}_\lambda \mathbf{D}^{-1} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{UD}_\lambda \mathbf{U}^T \mathbf{y} \\ &= \sum_{j=0}^p \left(\frac{d_j^2}{d_j^2 + \lambda} \right) \mathbf{U}_j \mathbf{U}_j^T \mathbf{y},\end{aligned}$$

\mathbf{U}_j is the column vector of \mathbf{U} . The **effective degree of freedom** of the ridge regression is

$$\begin{aligned}df(\lambda) &:= \text{Trace}(\mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T) = \text{Trace}(\mathbf{VD}_\lambda \mathbf{D}^{-1} \mathbf{U}^T) \\ &= \text{Trace}(\mathbf{D}_\lambda \mathbf{D}^{-1}) = \sum_{j=0}^p \frac{d_j^2}{d_j^2 + \lambda}\end{aligned}$$

decreases from $p + 1$ to 0 as λ increases from 0 to ∞ .

Ridge regression equivalent form

$$\hat{\beta}^{\lambda} = \operatorname{argmin}_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2,$$

$$\hat{\beta}^t = \operatorname{argmin}_{\|\beta\|_2^2 \leq t} \|\mathbf{y} - \mathbf{X}\beta\|_2^2.$$

One can prove that there exists a bijection between λ and t .

Exercise

Find the bijection between two positive scalars λ and t if two vectors β^{λ} and β^t are the same.

The square error can be rewritten as the \mathbf{D} -weighted L_2 distance to $\hat{\beta}^{OLS}$ in the transformed coordinate system

$$\begin{aligned} & \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \\ &= \left\| \mathbf{U}\mathbf{U}^\top \mathbf{y} - \mathbf{U}\mathbf{D}\mathbf{V}^\top \beta \right\|_2^2 = \left\| \mathbf{U}^\top \mathbf{y} - \mathbf{D}\mathbf{V}^\top \beta \right\|_2^2 \\ &= \|\tilde{\mathbf{y}} - \mathbf{D}\tilde{\beta}\|_2^2 = \|\mathbf{D}(\tilde{\beta}^{OLS} - \tilde{\beta})\|_2^2 \end{aligned}$$

For the constraint $\|\beta\|_2 = \|\mathbf{V}^\top \beta\| = \|\tilde{\beta}\| \leq t$, one can show the inequality problem $\min_{\|\tilde{\beta}\|_2 \leq t} \|\tilde{\mathbf{y}} - \mathbf{D}\tilde{\beta}\|_2^2$ actually attains the equality constraint for this case. The Lagrangian function is $\|\tilde{\mathbf{y}} - \mathbf{D}\tilde{\beta}\|_2^2 + \mu(\|\tilde{\beta}\|_2^2 - t)$. So, KKT gives $\mathbf{D}^2\tilde{\beta} - \mathbf{D}\tilde{\mathbf{y}} + \mu\tilde{\beta} = 0$, i.e., $\tilde{\beta} = (\mathbf{D}^2 + \mu\mathbf{I})^{-1}\mathbf{D}\tilde{\mathbf{y}}$. So, μ is the same as λ . The equality constraint $\|\tilde{\beta}\| = t$ determines uniquely $\beta = \mu = \mu(t)$.

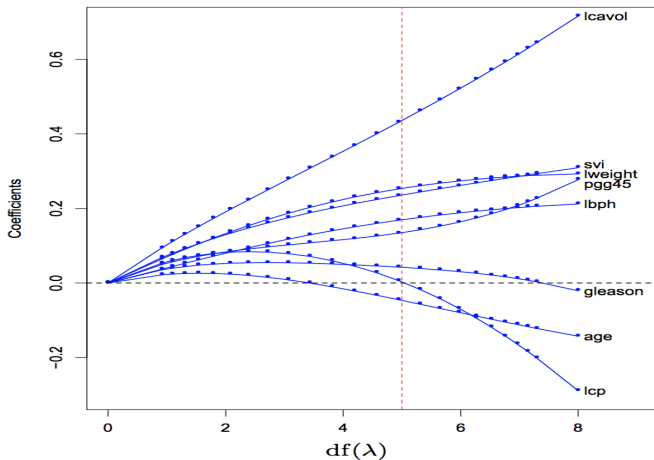


FIGURE 3.8. Profiles of ridge coefficients for the prostate cancer example, as the tuning parameter λ is varied. Coefficients are plotted versus $df(\lambda)$, the effective degrees of freedom. A vertical line is drawn at $df = 5.0$, the value chosen by cross-validation.

1

LASSO

(Frank and Friedman, 1993)

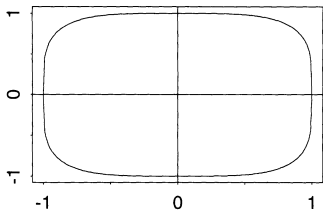
With $L_r(\beta) = \sum_{j=0}^p |\beta_j|^r$,

$$\hat{\beta}^{bridge} = \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda L_r(\beta)$$

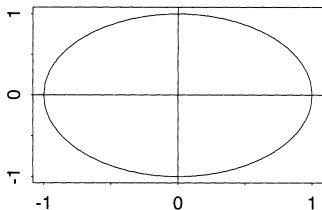
- $L_0(\beta) = \sum_{j=0}^p I(\beta_j \neq 0)$; (Hard thresholding) ¹
- $L_1(\beta) = \sum_{j=0}^p |\beta_j|$; (Lasso)
- $L_2(\beta) = \sum_{j=0}^p \beta_j^2$; (Ridge regression)
- $L_\infty(\beta) = \max_j |\beta_j|$.

¹This gives subset selection

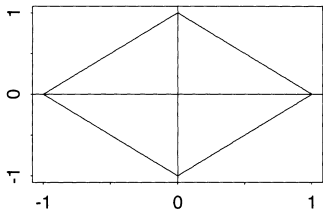
$\gamma > 2$



$\gamma = 2$



$\gamma = 1$



$\gamma < 1$

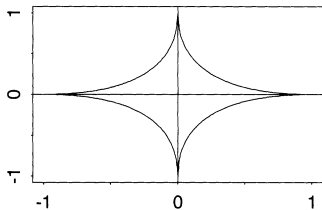


Figure 1. Constrained Areas of Bridge Regressions with $t = 1$.

Least Absolute Shrinkage and Selection Operator (Lasso)

Tibshirani (Journal of the Royal Statistical Society 1996) introduced the LASSO.

Lasso estimator: let $r = 1$ in bridge estimator

$$\hat{\beta}^{\lambda} = \underset{\beta}{\operatorname{argmin}} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right\},$$

$$\hat{\beta}^s = \underset{\|\beta\|_1 \leq s}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2,$$

where $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$.

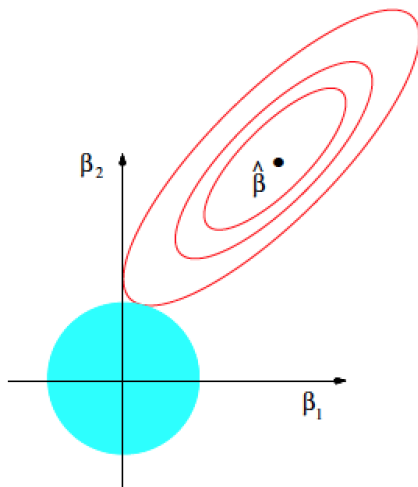
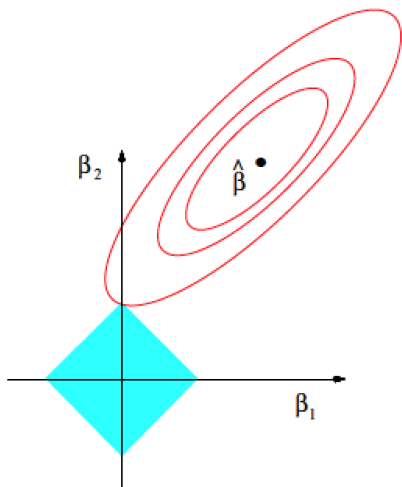
Again there exists a bijection between λ and s .

Sparse solution

- Due to the nature of the l_1 -norm constraint, if t is small enough some coefficients of the lasso solution become *exactly* zero.
- The elliptical contour is likely to hit the corner of the polytope, corresponding to sparse $\hat{\beta}$. Variable selection: drop the features with $\hat{\beta}_j = 0$.
- The l_r regularization results in sparsity when $0 \leq r \leq 1$, and is convex when $1 \leq r < \infty$.
- **Lasso is sparse and convex:**¹ the original implementation involves quadratic programming techniques from convex optimization
- Efron et al. (Annals of Statistics 2004) proposed **LARS (least angle regression)**, which computes the LASSO path efficiently
 - ▶ Interesting modification called is called forward stagewise
 - ▶ In many cases it is the same as the LASSO solution
 - ▶ Forward stagewise is easy to implement:

<http://www-stat.stanford.edu/~hastie/TALKS/nips2005.pdf>

¹not strictly convex. At a fixed λ , the coefficient $\hat{\beta}$ exists but may not be unique but the prediction $\hat{y} = \mathbf{X}\hat{\beta}$ is unique (Tibshirani, R. J. (2013). The lasso problem and uniqueness. Electronic Journal of Statistics, 7, 1456– 1490.).



$\hat{\beta}$ in the center is the $\hat{\beta}^{OLS}$

Consider a simple but illuminating example: Show the solutions

$$\hat{\beta}^{ridge} = \operatorname{argmin}_{\beta} \left\{ \|\beta - \hat{\beta}\|_2^2/2 + \lambda \|\beta\|_2^2 \right\},$$

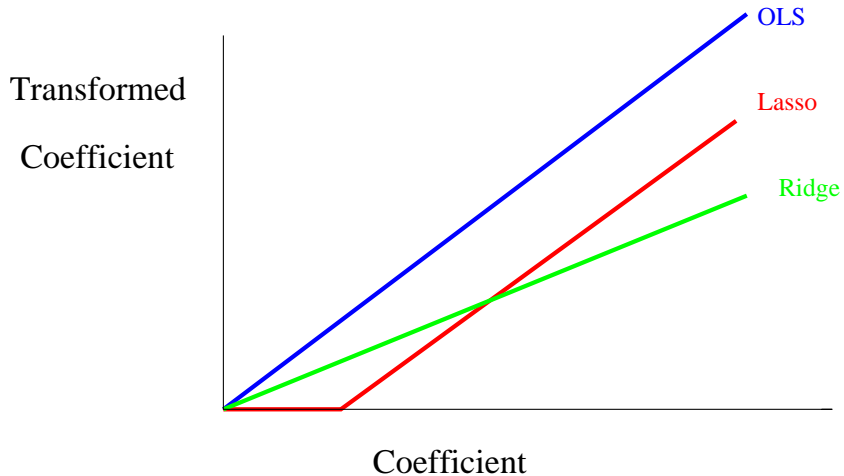
$$\hat{\beta}^{lasso} = \operatorname{argmin}_{\beta} \left\{ \|\beta - \hat{\beta}\|_2^2/2 + \lambda \|\beta\|_1 \right\}$$

are

- $\hat{\beta}_j^{ridge} = \hat{\beta}_j / (1 + \lambda)$
- $\hat{\beta}_j^{lasso} = \operatorname{sign}(\hat{\beta}_j) (|\hat{\beta}_j| - \lambda)_+$

visualization of $\hat{\beta}^{ridge}$ and $\hat{\beta}^{lasso}$ from desmos.com webpage

Ridge regression shrinks the $\hat{\beta}$ in all components/directions. Lasso translates them towards zero by a constant, truncating at zero.



Maximum number of selected covariates

The number of parameter/covariates selected by the lasso regression estimator is bounded non-trivially. The cardinality (i.e. the number of included covariates) of every lasso estimated linear regression model is smaller than or equal to $\min\{n, p\}$.

For any $\lambda > 0$, $\hat{\beta}^{lasso}(\lambda)$ has at most $\min(n, p)$ non-zeros entries. Here p is the number of parameters, and n is the number of training samples.¹

The number of all possible sub-models from variable selection (such as subset selection) is 2^p . However, the typical number of different Lasso-estimated sub-models is at the order of

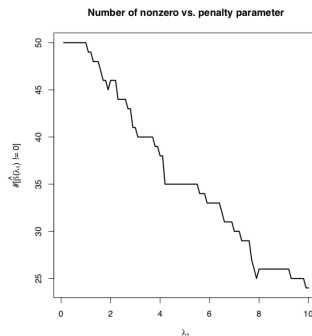
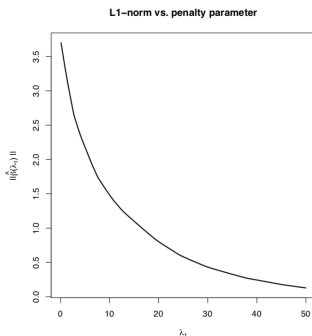
$$O(\min(n, p))$$

which is significantly smaller than 2^p if $p \gg n$.

¹Osborne, M. R., Presnell, B., and Turlach, B. A. (2000). On the lasso and its dual. *Journal of Computational and Graphical Statistics*, 9(2), 319–337.

sparsity of the lasso solution

- A large λ helps increase the sparsity ¹.
- However, $\left\| \hat{\beta}^{lasso}(\lambda) \right\|_1$ decreases monotonically as λ increases, but the number of non-zero coefficients does not.



source of figure

Solution path of LASSO

- The solution $\hat{\beta}^{lasso}(\lambda)$ is a **piecewise linear** function of λ .¹
- The path algorithm starts at $\lambda = \infty$ or $s = 0$, and traces the solution path by continuously changing λ .² Each new λ -solution is computed successively by solving the KKT conditions with a good initial guess set as the precise solution for a neighboring old λ .
- The key is to find the turning knots $\lambda_1, \dots, \lambda_T$
- An interesting reading: Efron et al. (AOS; 2003)
- Read the reference **Regularization Paths for Generalized Linear Models via Coordinate Descent**
- Path algorithms are available for many methods, such as fused lasso, trend filtering, locally adaptive regression splines, SVMs, 1-norm SVMs, relaxed maximum entropy method ...
- One can compute all possible Lasso sub-models $\cup_{\lambda>0} \{j : \hat{\beta}^{lasso}(\lambda) \neq 0\}$ with $O(np \min\{n, p\})$ operation counts.

¹Rosset, S. and Zhu, J. (2007). The Annals of Statistics, pp 1012–1030.

²homotopy method

Consistence of variable selection

- Empirical fact: The set of variables selected by Lasso contains the true (non-zero) variables in the sparse model, with high chance:

$$\left\{j : \hat{\beta}_j^{lasso} \neq 0\right\} \supseteq \left\{j : \hat{\beta}_j^{true} \neq 0\right\}$$

- The Theory(Meinshausen and Buhlmann 2006): Under some (restrictive and crucial) assumptions, if $\lambda = \lambda_n \gg \sqrt{\log(p)/n}$, then

$$\left\{j : \hat{\beta}_j^{lasso} \neq 0\right\} \rightarrow \left\{j : \hat{\beta}_j^{true} \neq 0\right\}, \quad p \gg n \rightarrow \infty$$

with prob. 1.

- For high-dimensional model selection, with strongly correlated design \mathbf{X} , the Lasso can perform very poorly for variable selection.

The motivation for the Lasso came from an interesting proposal of Breiman (1993). Breiman's non-negative Garotte minimizes

$$\min_{c_j \geq 0, \forall j} \frac{1}{2} \sum_{i=1}^n (y_i - \sum_{j=1}^p c_j \hat{\beta}_j x_{ij})^2 + \lambda \sum_{j=1}^p c_j$$

and then $\hat{\beta}_j^{ng} = \hat{c}_j \hat{\beta}_j$.

- It shrinks small $|\hat{\beta}_j|$ to zero. It is almost unbiased for large $|\hat{\beta}_j|$.
- Garotte starts with the OLS estimates and shrinks them by non-negative factors whose sum is constrained.
- In contrast, Lasso avoids the explicit use of the OLS estimates.
- Lasso is also closely related to the wavelet soft-thresholding method by Donoho and Johnstone (1994), and boosting method.

Exercise

Show when $\mathbf{X}^T \mathbf{X} = n\mathbf{I}$, the non-negative Garotte estimator is

$$\hat{\beta}_j^{ng} = \left(1 - \frac{\lambda}{2\hat{\beta}_j^2}\right)_+ \hat{\beta}_j.$$

One can just focus on the component-wise formulation

$$\min_{\beta} \frac{1}{2}(z_j - \beta_j)^2 + J(|\beta_j|)$$

- Hard thresholding: $J(|\beta_j|) = \lambda^2 - (|\beta_j| - \lambda)^2 I(|\beta_j| < \lambda)$, then¹

$$\hat{\beta}_j = z_j I(|z_j| > \lambda).$$

- Soft thresholding (Lasso): $J(|\beta_j|) = \lambda|\beta_j|$, then

$$\hat{\beta}_j = \text{sign}(z_j)(|z_j| - \lambda)_+.$$

¹the proof is exercise
Xiang Zhou

Ridge regression vs LASSO

- Both can yield a reduction in variance at the expense of a small increase in bias.
- Bayesian Interpretation for Ridge Regression and the Lasso: Figure 6.11, [ISL]. The prior distribution for β in ridge regression is Gaussian distribution; the prior distribution for β in LASSO is Laplace distribution $p(x) = \exp(-|x|/\lambda)/2\lambda$.
- Unlike ridge regression, the lasso performs variable selection when estimating the coefficients, and hence results in models easier to interpret.
- Ridge regression tends to give similar coefficient values to *correlated variables*, whereas the lasso may give quite different coefficient values to correlated variables.

Exercise

Ex. 7, Ex. 5, in Section 6.8, [ISL]

The Grouped Lasso

- In some problems, the predictors belong to pre-defined groups: genes that belong to the same biological pathway;
- Want to shrink and select the members of a group together: an entire group of predictors may drop out of the model.
- Suppose that the p predictors are divided into L groups. Group ℓ has p_ℓ number of predictors, with design matrix \mathbf{X}_ℓ and the coefficient $\beta_\ell \in \mathbb{R}^{p_\ell}$. The grouped-lasso minimizes the convex criterion:

$$\min_{\beta \in \mathbb{R}^p} \left(\left\| \mathbf{y} - \beta_0 \mathbf{1} - \sum_{\ell=1}^L \mathbf{X}_\ell \beta_\ell \right\|_2^2 + \lambda \sum_{\ell=1}^L \sqrt{p_\ell} \|\beta_\ell\|_2 \right).$$

This procedure encourages sparsity at both the group and individual levels

Zou and Hastie (2005) introduced the elastic-net penalty:

$$\lambda \sum_{j=1}^p (\alpha \beta_j^2 + (1 - \alpha) |\beta_j|)$$

a different compromise between ridge and lasso from using $\|\cdot\|_r$ norm.

The elastic-net selects variables like the lasso, and shrinks together the coefficients of correlated predictors like ridge.

One can modify the lasso penalty function so that larger coefficients are shrunk less severely (this may help reduce the unnecessary bias); the smoothly clipped absolute deviation (SCAD) penalty of Fan and Li (2005):

$$\min_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p q_{\lambda}(|\beta_j|),$$

where for some $a \geq 2$, such that the derivative $q_{\lambda}(|\beta|)$ is

$$\frac{dq_{\lambda}(\beta)}{d\beta} = \lambda \cdot \text{sign}(\beta) \left[I(|\beta| \leq \lambda) + \frac{(a\lambda - |\beta|)_+}{(a-1)\lambda} I(|\beta| > \lambda) \right]$$

- The derivatives decreases from λ to 0 as $|\beta|$ increases.
- $q_{\lambda}(\beta)$ is not convex, a drawback for computation.

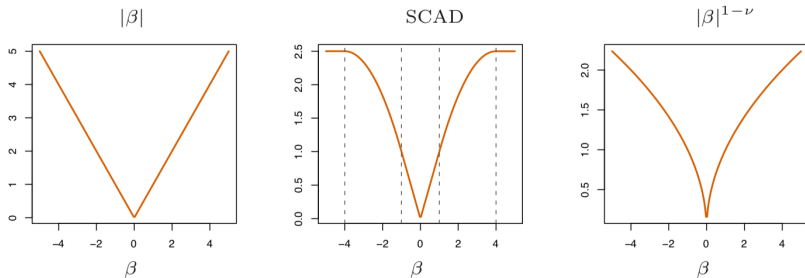
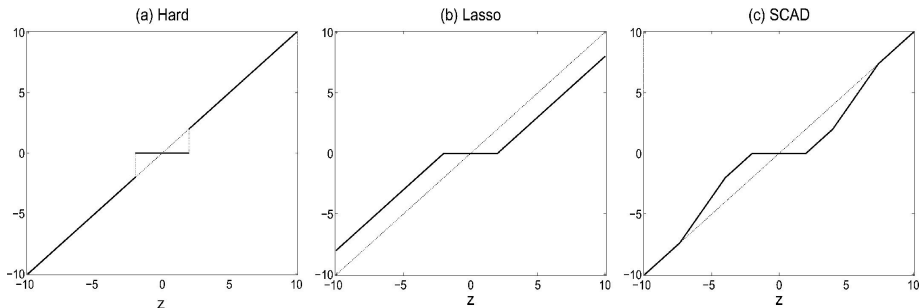


FIGURE 3.20. *The lasso and two alternative non-convex penalties designed to penalize large coefficients less. For SCAD we use $\lambda = 1$ and $a = 4$, and $\nu = \frac{1}{2}$ in the last panel.*

$$\beta \rightarrow \hat{\beta}$$



Reference: Antoniadis and Fan (2001); Fan and Li (2001)

Adaptive Lasso (Zou, 2006; JASA)

Adaptive Lasso is an important two-stage procedure to address some bias problems of the Lasso.

$$\min_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p w_j |\beta_j|,$$

where $w_j > 0$ and adjusts the penalty on each β_j .

- The resulting estimator is

$$\hat{\beta}_j = \text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda w_j)_+.$$

- A reasonable choice for w_j is obtained from the first stage:

$$w_j = |\hat{\beta}_j^{OLS}|^{-\nu}, \text{ or } |\hat{\beta}_j^{\text{lasso}}|^{-\nu}$$

with $\nu > 0$

- If $|w_j|$ is small, the adaptive Lasso employs a small regularization for the coefficient β_j which implies less bias.

What not covered here

- Optimization theory and computational method for LASSO: Least Angle Regression, Dantzig Selector
- Principal Components Regression, partial linear regression (Section 3.5 in [ESL])
- Further reading on LASSO: *Statistics for High-Dimensional Data: Methods, Theory and Applications*, by Peter Bühlmann and Sara van de Geer, Springer Series in Statistics (2011)

Revisit Bias-Variance

Revisit Bias-Variance

- In the beginning, We calculated the MSE for ordinary least square where the bias vanishes.
- We now focus on the ridge regression, $\hat{\beta}^\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$, where $\lambda > 0$ is the penalty parameter. The corresponding MSE for $\hat{\beta}^\lambda$ at the new testing point x_0 then is written as ¹

$$\begin{aligned} \text{MSE}^\lambda(x_0) &= \text{Var}_{\mathbf{y}}(x_0^\top \hat{\beta}^\lambda) + \left(\mathbb{E}_{\mathbf{y}}(x_0^\top \hat{\beta}^\lambda) - f(x_0) \right)^2 \\ &= \text{Var}_{\mathbf{y}}(x_0^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}) + \left(x_0^\top \mathbb{E}_{\mathbf{y}}(\hat{\beta}^\lambda) - f(x_0) \right)^2 \\ &= \sigma_\varepsilon^2 \|\mathbf{h}_\lambda x_0\|_2^2 + \left(x_0^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \beta - x_0^\top \beta \right)^2 \end{aligned}$$

where $\mathbf{h}_\lambda = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1}$. Exercise: with the aid of SVD, find the minimum point λ for MSE^λ .

¹again, the training data's \mathbf{X} part is fixed or conditioned.
Xiang Zhou CityU

Decomposition of average squared bias

- Note that we have assumed that the ground truth is a linear model $f(x) = x^\top \beta$.
- From now we do not assume f is linear and it could be a nonlinear function for a general consideration. This is a more realistic setting.
- The additive error model $Y = f(X) + \varepsilon$ is still assumed for data.
- The best-fitting approximation in the *linear model class*¹ is given by

$$\beta^* = \underset{\beta}{\operatorname{argmin}} \mathbb{E}_X (f(X) - X^\top \beta)^2$$

- Note that β_* satisfies the normal equation:

$$\mathbb{E}(X^\top X) \beta_* = \mathbb{E}[f(X)X]$$

¹In other words, $x \rightarrow x^\top \beta^*$ is $f_{\mathcal{H}}(x)$ defined in Topic 1 (\mathcal{H} is hypothesis space)

Decomposition of average (squared) bias

- We still consider the MSE^λ for the ridge regression. The variance part $\text{Var}(x_0^\top \hat{\beta}^\lambda)$ is unchanged. Now the squared bias becomes
$$\begin{aligned} \left(x_0^\top \mathbb{E}_{\mathbf{y}}(\hat{\beta}^\lambda) - f(x_0)\right)^2 &= \left(x_0^\top \mathbb{E}_{\mathbf{y}}(\hat{\beta}^\lambda) - x_0^\top \beta_* + x_0^\top \beta_* - f(x_0)\right)^2 = \\ &= \left(x_0^\top \mathbb{E}_{\mathbf{y}}(\hat{\beta}^\lambda) - x_0^\top \beta_*\right)^2 + \left(x_0^\top \beta_* - f(x_0)\right)^2 + \\ &\quad 2 \left(x_0^\top \mathbb{E}_{\mathbf{y}}(\hat{\beta}^\lambda) - x_0^\top \beta_*\right) \left(x_0^\top \beta_* - f(x_0)\right) \end{aligned}$$
- Taking expectation for $x_0 \sim X$ and using the normal equation, we have the average squared bias is

$$\begin{aligned} &\mathbb{E}_{x_0} \left(x_0^\top \mathbb{E}_{\mathbf{y}}(\hat{\beta}^\lambda) - f(x_0) \right)^2 \\ &= \underbrace{\mathbb{E}_{x_0} \left(x_0^\top (\mathbb{E}_{\mathbf{y}} \hat{\beta}^\lambda - \beta_*) \right)^2}_{\text{Ave (Estimation Bias}^2)} + \underbrace{\mathbb{E}_{x_0} \left(x_0^\top \beta_* - f(x_0) \right)^2}_{\text{Ave (Model Bias}^2)} \end{aligned}$$

This is equation (7.14) in [ELS].

- Model Bias, by definition, involves the ground truth, which is not accessible; Estimation Bias, by definition, involves β^* , which is not accessible either.
- However, this decomposition is conceptually inspiring to have Model Bias.
- The Model Bias corresponds to the approximation error in Topic 1 ($f_{\mathcal{H}}$ vs f^*) and the Estimation Bias is similar to the sampling error in Topic 1 ($\hat{f}_{\mathcal{D}}$ vs $f_{\mathcal{H}}$). In approximation theory viewpoint, the variance $\text{Var}_{\mathcal{D}} \hat{f}_{\mathcal{D}}$ is not considered.
- Ridge method and LASSO further restricts the linear model \mathcal{H}_{linear} to a smaller set in the form $\mathcal{H}_{\lambda} := \{f : f(x) = \beta^T x, \|\beta\| \leq \lambda\}$; they affect the estimation bias (and the variance of the prediction). The improvement of Model Bias needs go from linear to nonlinear if the ground truth is far away from being linear.