

# Regularization and Variable Selection

Wei Li

Syracuse University

Spring 2021



# Regularization and Variable Selection

Motivation

Subset Selection

Model Selection

Lasso

Ridge Regression

Consistency of LASSO

Computation of LASSO

Beyond Lasso and Remarks

# Motivation

# Predictive Accuracy

Suppose  $Y_i = \beta^T X_i + \epsilon_i$ ,  $X_i$  a  $p$ -vector, where  $E[\epsilon_i | X_i] = 0$  and  $\text{Var}(\epsilon_i | X_i) = \sigma^2$ .

Then Ordinary Least Squares has low bias. But for variance:

- ▶ When  $n \gg p$ , OLS tends to have low variance, and hence perform well on test observations.
- ▶ When  $n < p$ , then there is no longer OLS estimate, and the variance of these estimates are infinite. May be able to increase bias slightly, but decrease variance substantially via
  - ▶ regularization/shrinkage
  - ▶ features selection
  - ▶ dimension reduction
- ▶ When  $n$  is not much larger than  $p$ , LS fit can have high variance and may result in over fitting
- ▶ For a truly linear underlying model, the bias of the linear regression estimate is exactly 0, and the variance is  $p \cdot \sigma^2 / n$

The **key message**: It is possible to trade a little bias with the large reduction in variance, thus achieving higher prediction accuracy.

# Model Interpretability

- ▶ When we have a large number of variables  $X$  in the model there will generally be many that have little or no effect on  $Y$
- ▶ Leaving these variables in the model makes it harder to see the “big picture”, i.e., the effect of the “important variables”
- ▶ The model would be easier to interpret by removing (i.e. setting the coefficients to zero) the unimportant variables. This can be accomplished via variable selection, or feature selection.

# Subset Selection

# Best Subset Selection

1. For each  $k \in \{0, 1, \dots, p\}$ , find the subset of size  $k$  that gives smallest residual sum of squares
  - ▶ fit all  $\binom{p}{k}$  models that contains  $k$  predictors
  - ▶ pick the best subset (model) with smallest RSS or  $R^2$ .
2. Then among the  $p + 1$  chosen models, pick a single best model using some criterion discussed below.



# Sequential Selection

- ▶ **Forward Stepwise Selection:** Begins with the model containing no predictor, and then adds one predictor at a time that improves the model the most (say, by  $R^2$ , or RSS ) until no further improvement is possible.
- ▶ **Backward Stepwise Selection:** Begins with the model containing all predictors, and then deleting one predictor at a time that improves the model the most (say, by  $R^2$ , or RSS ) until no further improvement is possible.

Note that both procedure produce a sequence of models of different size  $k$  (i.e., different number of predictors  $k$ ).

# Model Selection

# Model Selection

Note that the larger  $k$  (model size), the smaller RSS, or higher  $R^2$ . To choose the best model size  $k^*$ , one can use a number of different criteria:

- ▶ Mallow's  $C_p$ , AIC, BIC (more in Sec 7.5)
- ▶ adjusted- $R^2$
- ▶ prediction error on the test set
- ▶ cross validation; generalized cross validation (GCV) (more in Sec 7.10)

Suppose the data are  $(\mathbf{x}_i, y_i), i = 1, \dots, n$ . A fitted linear regression model is  $\hat{f}(\mathbf{x}) = \hat{\boldsymbol{\beta}}^T \mathbf{x}$ .

- ▶ The degree of freedom (df) of  $\hat{\boldsymbol{\beta}}$  as the number of nonzero elements in  $\boldsymbol{\beta}$  (model size, say  $k$ ), including the intercept
- ▶ The residual sum of squares as  $RSS = \sum_{i=1}^n \left( y_i - \hat{\boldsymbol{\beta}}^T \mathbf{x}_i \right)^2$

$$AIC(k) = n \log(RSS/n) + 2 \cdot k$$

$$BIC(k) = n \log(RSS/n) + \log(n) \cdot k$$

- ▶ Suppose we can estimate  $\sigma^2$  using a full model, its estimate is given by  $\tilde{\sigma}^2$ ,

$$C_p(k) = \frac{1}{n} (RSS + 2k \cdot \tilde{\sigma}^2)$$

- ▶ Adjusted  $R^2 = 1 - \frac{RSS/(n-k)}{TSS/(n-1)}$

# Forward Stagewise Linear Regression

Let's assume that the features have been standardized to sample mean 0 and sample variance 1.

Forward Stagewise Linear Regression, henceforth called Stagewise, is an iterative technique that begins with  $\hat{\boldsymbol{\mu}} = 0$  and builds up the regression function in successive small steps. If  $\hat{\boldsymbol{\mu}}$  is the current Stagewise estimate, let  $\mathbf{c}(\hat{\boldsymbol{\mu}})$  be the vector of *current correlations*

$$\hat{\mathbf{c}} = \mathbf{c}(\hat{\boldsymbol{\mu}}) = \mathbf{X}^T(\mathbf{y} - \hat{\boldsymbol{\mu}})$$

so that  $\hat{c}_j$  is proportional to the correlation between covariate  $\mathbf{x}_j$  and the current residual vector. The next step of the Stagewise algorithm is taken in the direction of the greatest current correlation,

$$\hat{j} = \operatorname{argmax}_j |\hat{c}_j| \quad \text{and} \quad \hat{\boldsymbol{\mu}} \leftarrow \hat{\boldsymbol{\mu}} + \epsilon \cdot \operatorname{sign}(\hat{c}_{\hat{j}}) \cdot \mathbf{x}_{\hat{j}}$$

with  $\epsilon$  some small constant. “Small” is important here: the “big” choice  $\epsilon = |\hat{c}_{\hat{j}}|$  leads to a version of Forward Selection algorithm.

# Lasso

# Lasso

A generalized regularization framework:

$$\min_{\beta} L(\beta; \mathbf{y}, \mathbf{X}) + \lambda J(\beta), \lambda \geq 0$$

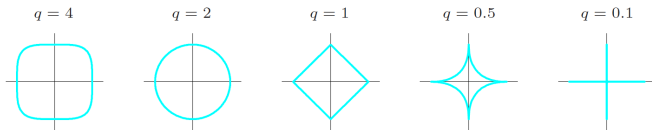
$L(\beta; \mathbf{y}, \mathbf{X})$  is the loss function

- ▶ For OLS,  $L = \sum_{i=1}^n \left[ y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right]^2$
- ▶ For MLE methods,  $L = -\log$  likelihood
- ▶ Others: hinge loss function (SVM), or exponential loss (AdaBoost)

$J(\beta)$  is the penalty function.

$$J_q(|\beta|) = \|\beta\|_q^q = \sum_{j=1}^p |\beta_j|^q, \quad q \geq 0$$

- ▶  $J_0(|\beta|) = \|\beta\|_0 = \sum_{j=1}^p 1(\beta_j \neq 0)$ . (Best subset; Donoho and Johnstone 1988.)
- ▶  $J_1(|\beta|) = \|\beta\|_1 = \sum_{j=1}^p |\beta_j|$  (LASSO; Tibshirani 1996)
- ▶  $J_2(|\beta|) = \|\beta\|_2^2 = \sum_{j=1}^p \beta_j^2$  (Ridge; Hoerl and Kennard, 1970)
- ▶  $J_\infty(|\beta|) = \|\beta\|_\infty = \max_j |\beta_j|$  (Supnorm penalty; Zhang et al. 2008)



ESL: Fig 3.12 contours of constant value of  $\sum_j |\beta_j|^q$



For the LS problem, this gives rise to

$$\begin{array}{ll} \min_{\boldsymbol{\beta} \in \mathbb{R}^p} & \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_0 \quad (\text{Best subset selection}) \\ \min_{\boldsymbol{\beta} \in \mathbb{R}^p} & \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \quad (\text{Lasso regression}) \\ \min_{\boldsymbol{\beta} \in \mathbb{R}^p} & \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \quad (\text{Ridge regression}) \end{array}$$

Another way to express the problem is

$$\min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \text{ subject to } \|\beta\|_0 \leq t \text{ (Best subset selection)}$$

$$\min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \text{ subject to } \|\beta\|_1 \leq t \text{ (Lasso regression)}$$

$$\min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \text{ subject to } \|\beta\|_2^2 \leq t \text{ (Ridge regression)}$$

$$t \geq 0.$$

# Lasso

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

In practice, we center predictors and response, and standardize predictors.

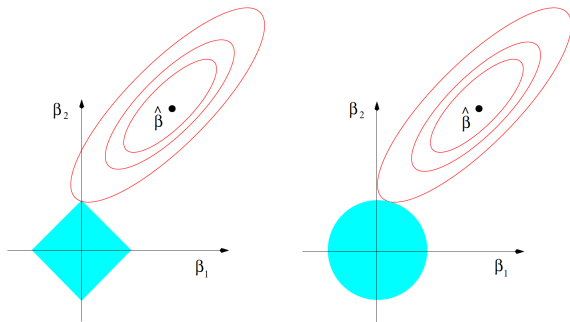
- ▶  $\lambda \geq 0$  is a tuning parameter
- ▶  $\lambda$  controls the amount of shrinkage; the larger  $\lambda$ , the greater amount of shrinkage
- ▶ What happens if  $\lambda \rightarrow 0$ ? (no penalty)
- ▶ What happens if  $\lambda \rightarrow \infty$ ?

The equivalent optimization problem is

$$\begin{aligned}\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \quad & \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \\ \text{subject to} \quad & \sum_{j=1}^p |\beta_j| \leq t\end{aligned}$$

- ▶ Explicitly apply the magnitude constraint on the parameters
- ▶ Making  $t$  sufficiently small will cause some coefficients to be exactly zero.
- ▶ The lasso is a kind of continuous subset selection

Note,  $t$  can be chosen between 0 and  $\sum_j |\hat{\beta}_j|$  where  $\hat{\beta}_j$  are OLS estimates. One can define shrinkage factor  $s = t / \sum_j |\hat{\beta}_j|$ .



ESL: Fig 3.11 Shown are contours of the error and constraint functions.

The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.

## Choose $\lambda$

The tuning parameter  $t$  or  $\lambda$  should be adaptively chosen to minimize the MSE or PE (prediction error).

- ▶ If  $t$  is chosen larger than  $t_0 = \sum_{j=1}^p |\hat{\beta}_j^{ols}|$ , then  $\hat{\beta}_j^{\text{lasso}} = \hat{\beta}_j^{ols}$ .
- ▶ If  $t = t_0/2$ , the least squares coefficients are shrunk by 50%.
- ▶ In practice, we sometimes standardize  $t$  by  $s = t/t_0$  in  $[0, 1]$ , and choose the best value  $s$  from  $[0, 1]$ .

The simplest and most popular way to estimate the prediction error.

1. Randomly split the data into  $K$  roughly equal parts
2. For each  $k = 1, \dots, K$ 
  - ▶ leave the  $k$  th portion out, and fit the model using the other  $K - 1$  parts. Denote the solution by  $\hat{\beta}^{-(k)}$
  - ▶ calculate prediction errors of  $\hat{\beta}^{-(k)}$  on the left-out  $k$  th portion
3. Average the prediction errors

Define the Index function (allocating memberships)

$$\kappa : \{1, \dots, n\} \rightarrow \{1, \dots, K\}$$

The  $K$  -fold cross-validation estimate of prediction error (PE) is

$$CV = \frac{1}{n} \sum_{i=1}^n \left( y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^{-(\kappa(i))} \right)^2$$

Typical choice:  $K = 5, 10$

For Lasso, choose a sequence of  $\lambda$  values, say  $\{\lambda_1, \dots, \lambda_m\}$ :

- ▶ For each  $\lambda$ , fit the LASSO and denote the solution by  $\hat{\beta}_\lambda^{\text{lasso}}$ .
- ▶ Compute the  $CV(\lambda)$  curve as

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left( y_i - \mathbf{x}_i^T \hat{\beta}_\lambda^{-\left(\kappa(i)\right)} \right)$$

which provides an estimate of the test error curve.

- ▶ Find the best parameter  $\lambda^*$  which minimizes  $CV(\lambda)$ .

$$\lambda^* = \underset{\lambda \in \{\lambda_1, \dots, \lambda_m\}}{\operatorname{argmin}} \quad CV(\lambda)$$

- ▶ Fit the final LASSO model with  $\lambda^*$ . The final solution is denoted as  $\hat{\beta}_{\lambda^*}^{\text{lasso}}$

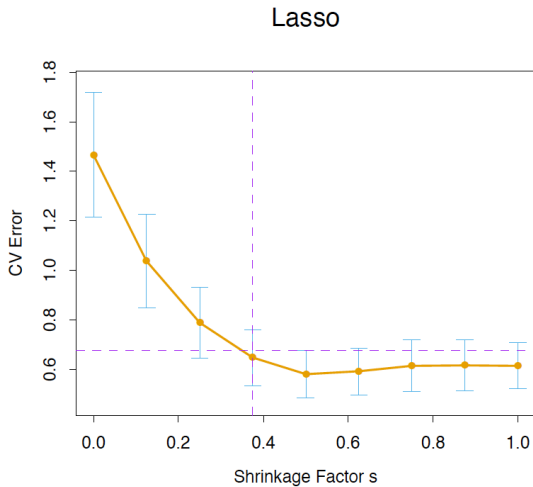


## One-standard error rule

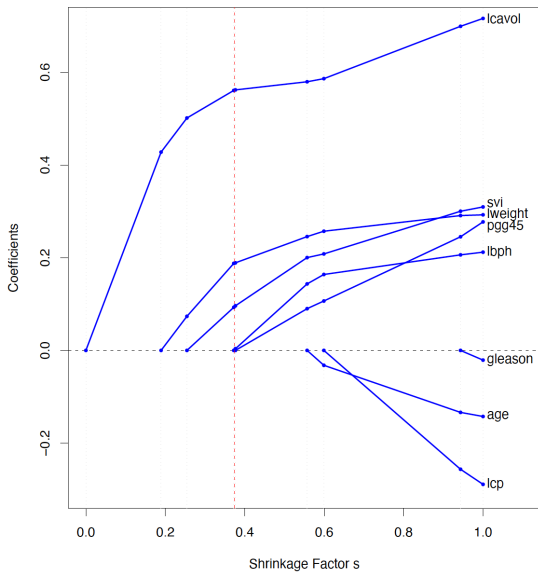
- ▶ choose the most parsimonious model with error no more than one standard error above the best error.
- ▶ used often in model selection (for a smaller model size).

$$\hat{\lambda} = \max \{ \lambda \in \{ \lambda_1, \dots, \lambda_m \} : \text{CV}(\lambda) \leq \text{CV}(\lambda^*) + \text{SE}(\text{CV}(\lambda^*)) \}$$

# Prostate Cancer Data Example



ESL Figure 3.7



ESL Figure 3.10

# Ridge Regression

# Ridge Regression

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

An equivalent way to write the ridge problem is

$$\begin{aligned} \hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \\ \text{subject to } \sum_{j=1}^p \beta_j^2 \leq t \end{aligned}$$

$$\text{RSS}(\lambda) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T\beta$$

The ridge regression solutions is

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

**Theorem:** For any design matrix  $\mathbf{X}$ , the quantity  $\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p$  is always invertible for  $\lambda > 0$ ; thus, there is always a unique solution  $\hat{\beta}^{\text{ridge}}$  (Remember,  $\mathbf{X}^T\mathbf{X}$  must be PSD, but not necessarily PD).

Also, there always exists a parameter  $\lambda^* > 0$  such that mean-squared-error of is strictly smaller than mean-squared-error of the OLS estimate.

Note the similarity to the ordinary least squares solution, but with the addition of a “ridge” down the diagonal.

- ▶ As  $\lambda \rightarrow 0$ ,  $\hat{\beta}^{\text{ridge}} \rightarrow \hat{\beta}^{\text{ols}}$
- ▶ As  $\lambda \rightarrow \infty$ ,  $\hat{\beta}^{\text{ridge}} \rightarrow \mathbf{0}$

In the special case of an orthonormal design matrix

$$\hat{\beta}_j^{\text{ridge}} = \frac{\hat{\beta}_j^{\text{ols}}}{1 + \lambda}$$

This illustrates the essential “shrinkage” feature of ridge regression

- ▶ Applying the ridge regression penalty has the effect of shrinking the estimates toward zero
- ▶ The penalty introducing bias but reducing the variance of the estimator
- ▶ Ridge estimator does not threshold, since the shrinkage is smooth (proportional to the original coefficient).

# Interpreting shrinkage

The SVD of  $\mathbf{X}$  ( $n \times p$ ) has the form

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

- ▶ Here  $\mathbf{U}$  and  $\mathbf{V}$  are  $n \times p$  and  $p \times p$  orthonormal and orthogonal matrices, with the columns of  $\mathbf{U}$  spanning the column space of  $\mathbf{X}$ , and the columns of  $\mathbf{V}$  spanning the row space.
- ▶  $\mathbf{D}$  is a  $p \times p$  diagonal matrix, with diagonal entries  $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$  called the singular values of  $\mathbf{X}$ . If one or more values  $d_j = 0$ ,  $\mathbf{X}$  is singular.

One can show that

$$\mathbf{X}\hat{\beta}^{\text{ridge}} = \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y}$$
$$\mathbf{X}\hat{\beta}^{\text{ols}} = \sum_{j=1}^p \mathbf{u}_j \mathbf{u}_j^T \mathbf{y}$$



**Principal component analysis** (PCA, Pearson 1901) is a statistical procedure that

- ▶ uses an orthogonal transformation to convert a set of observations of correlated variables into a set of linearly uncorrelated variables (called principal components)
- ▶ finds directions with maximum variability

The transformed data matrix is given as  $\mathbf{Z} = \mathbf{X}V$  where  $V$  is obtained through SVD of  $\mathbf{X}$ . Columns of  $V$  give the **principal component directions**.

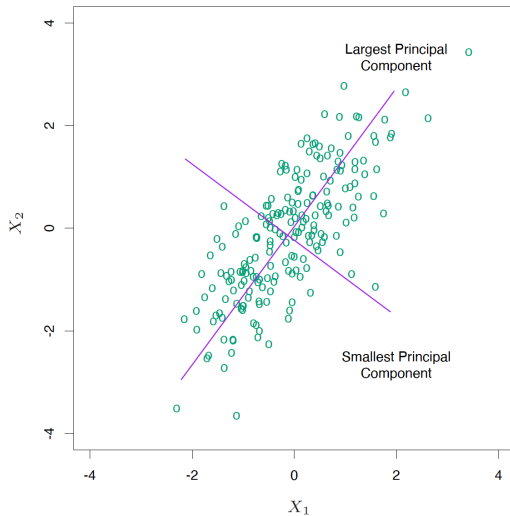
The **first principal component direction**  $v_1$  has the property that  $\mathbf{z}_1 = \mathbf{X}v_1$  has the largest sample variance amongst all normalized linear combinations of the columns of  $\mathbf{X}$ . This sample variance can be seen as

$$\text{Var}(\mathbf{z}_1) = \text{Var}(\mathbf{X}v_1) = \frac{d_1^2}{n}$$

and in fact  $\mathbf{z}_1 = \mathbf{X}v_1 = \mathbf{u}_1d_1$ . The derived variable  $\mathbf{z}_1$  is called the **first principal component** of  $\mathbf{X}$ .

Subsequent principal components  $\mathbf{z}_j$  have maximum variance  $d_j^2/n$ , subject to being orthogonal to the earlier ones. Conversely the last principal component has minimum variance.

Hence the small singular values  $d_j$  correspond to directions in the column space of  $\mathbf{X}$  having small variance, and ridge regression shrinks these directions the most.



ESL Figure 3.9

# Degrees of Freedom

For the ridge estimator, an unbiased estimate for its  $df$  (details in Sec 7.6) is

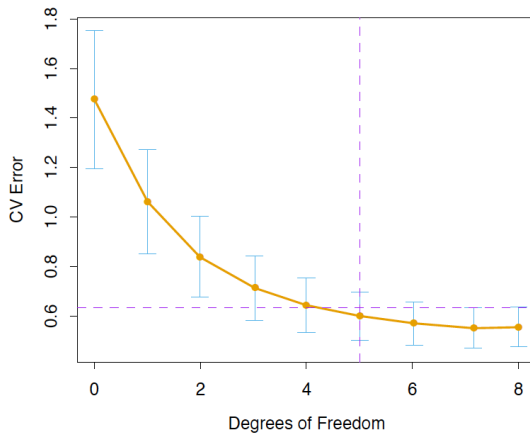
$$\begin{aligned}\widehat{df}\left(\widehat{\boldsymbol{\beta}}^{\text{ridge}}\right) &= \text{trace}\left\{\mathbf{X}\left(\mathbf{X}^T\mathbf{X}+\lambda\mathbf{I}_d\right)^{-1}\mathbf{X}^T\right\} \\ &= \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}\end{aligned}$$

Note:

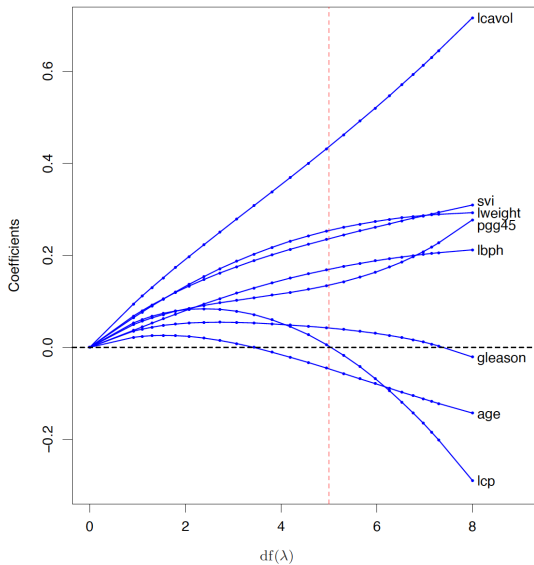
- ▶  $df(\lambda)$  is monotone decreasing functions of  $\lambda$ , called effective d.f.
- ▶  $df(\lambda) = p$  when  $\lambda = 0$ ;  $df(\lambda) \rightarrow 0$  if  $\lambda \rightarrow \infty$

# Prostate Cancer Data Example

## Ridge Regression



ESL Figure 3.7



ESL Figure 3.8

# A comparison between ridge regression and OLS

- ▶ If the design matrix  $\mathbf{X}$  is not full rank, then  $\mathbf{X}^T\mathbf{X}$  is not invertible. For any design matrix  $\mathbf{X}$ , the quantity  $\mathbf{X}^T\mathbf{X} + \lambda I_p$  is always invertible; thus, there is always a unique solution  $\hat{\beta}^{\text{ridge}}$
- ▶ There always exists a  $\lambda$  such that the MSE of  $\beta^{\text{ridge}}$  is less than the MSE of  $\hat{\beta}^{\text{ols}}$
- ▶ Recall that if the true model is approximately linear, the OLS estimates generally have low bias but can still be highly variable. The penalty term makes the ridge regression estimates biased but can also substantially reduce variance

# A comparison between ridge regression and the lasso

Ridge regression:

- ▶ is a continuous shrinkage method; achieves better prediction performance than OLS through a biase-variance trade-off.
- ▶ cannot produce a parsimonious model, for it always keeps all the predictors in the model.

The lasso regression:

- ▶ does both continuous shrinkage and automatic variable selection simultaneously. The lasso not only set coefficients to zero exactly, but it also shrinks the nonzero coefficients.



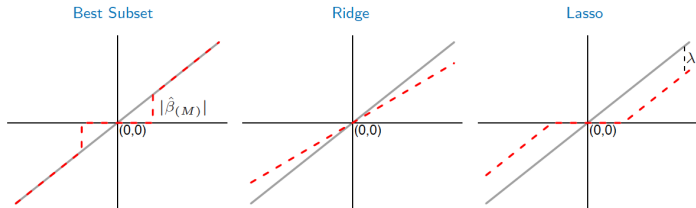
## A comparison for orthonormal design: $\mathbf{X}^T \mathbf{X} = I$

- ▶ The solution to the Lasso problem is

$$\begin{aligned}\hat{\beta}_j^{\text{lasso}} &= \text{sign} \left( \hat{\beta}_j^{\text{ols}} \right) \left( \left| \hat{\beta}_j^{\text{ols}} \right| - \frac{\lambda}{2} \right)_+ \\ &= \begin{cases} \hat{\beta}_j^{\text{ols}} - \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{ols}} > \frac{\lambda}{2} \\ 0 & \text{if } \left| \hat{\beta}_j^{\text{ols}} \right| \leq \frac{\lambda}{2} \\ \hat{\beta}_j^{\text{ols}} + \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{ols}} < -\frac{\lambda}{2} \end{cases}\end{aligned}$$

- ▶ shrinks big coefficients by a constant  $\frac{\lambda}{2}$  towards zero.
- ▶ truncates small coefficients to zero exactly.
- ▶ “soft-thresholding”
- ▶ Best subset (of size  $M$ ) :  $\hat{\beta}_j^{\text{ols}}$  if  $\text{rank} \left( \left| \hat{\beta}_j^{\text{ols}} \right| \right) \leq M$  keeps the largest coefficients; “hard-thresholding”
- ▶ Ridge regression:  $\hat{\beta}_j^{\text{ols}} / (1 + \lambda)$ : does a proportional shrinkage.

Estimator	Formula
Best subset (size $M$ )	$\hat{\beta}_j \cdot I( \hat{\beta}_j  \geq  \hat{\beta}_{(M)} )$
Ridge	$\hat{\beta}_j / (1 + \lambda)$
Lasso	$\text{sign}(\hat{\beta}_j)( \hat{\beta}_j  - \lambda)_+$



ESL Table 3.4: note that the lasso expression uses the  $RSS/2$  criterion.

# Consistency of LASSO

# Consistency of LASSO

It does not always have the oracle property, that is, it does not always perform as well in terms of variable selection as if the true underlying model has been given (Fan and Li, 2001).

Certain conditions are needed for the lasso to have the oracle property (Fan and Li, 2001; Zou, 2006).

True model:  $y_i = \beta_{00} + \sum_{j=1}^p x_{ij}\beta_{j0} + \epsilon_i$ .

important index set:  $\mathcal{A}_0 = \{j : \beta_{j0} \neq 0, j = 1, \dots, p\}$ .

unimportant index set:  $\mathcal{A}_0^c = \{j : \beta_{j0} = 0, j = 1, \dots, p\}$ .

An oracle performs as if the true model were known:

- ▶ selection consistency

$$P\left(\hat{\beta}_j \neq 0 \text{ for } j \in \mathcal{A}_0; \quad \hat{\beta}_j = 0 \text{ for } j \in \mathcal{A}_0^c\right) \rightarrow 1$$

- ▶ estimation consistency:

$$\sqrt{n}\left(\hat{\beta}_{\mathcal{A}_0} - \beta_{\mathcal{A}_0}\right) \rightarrow_d N(0, \Sigma_I)$$

$\beta_{\mathcal{A}_0} = \{\beta_j, j \in \mathcal{A}_0\}$  and  $\Sigma_I$  is the covariance matrix if knowing the true model.

Knight and Fu (2000)

- ▶ Estimation Consistency: The LASSO solution is of estimation consistency for fixed  $p$ . In other words,

$$\hat{\beta}^{\text{lasso}}(\lambda_n) \rightarrow_p \beta, \text{ as } \lambda_n = o(n)$$

It is root-  $n$  consistent and asymptotically normal.

- ▶ Model Selection Property: For  $\lambda_n \propto n^{\frac{1}{2}}$ , as  $n \rightarrow \infty$ , there is a non-vanishing positive probability for lasso to select the true model.

Zhao and Yu (2006):

Under the Irrepresentable Condition (IC), the LASSO is model selection consistent in both fixed and large  $p$  settings.

$$\max_{j \in \mathcal{S}^c} \left\| \mathbf{x}_j^T \mathbf{X}_{\mathcal{S}} (\mathbf{X}_{\mathcal{S}}^T \mathbf{X}_{\mathcal{S}})^{-1} \right\|_1 \leq (1 - \epsilon) \text{ for some } \epsilon \in (0, 1]$$

Here  $\mathcal{S}$  indexes the subset of features with non-zero coefficients in the true underlying model, and  $\mathbf{X}_{\mathcal{S}}$  are the columns of  $\mathbf{X}$  corresponding to those features. Similarly defined for  $\mathcal{S}^c$  and  $\mathbf{X}_{\mathcal{S}^c}$

# Computation of LASSO



# Computation of LASSO

- ▶ LARS solution path (Efron et al. 2004)
  - ▶ the most efficient algorithm for LASSO
  - ▶ designed for standard linear models
  - ▶ R package “lars”
- ▶ Coordinate Descent Algorithm (CDA; Friedman et al. 2010)
  - ▶ designed for GLM
  - ▶ suitable for ultra-high dimensional problems
  - ▶ R package “glmnet”

# Coordinate descent for lasso

The simple idea is that given convex, differentiable  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , if we are at a point  $x$  such that  $f(x)$  is minimized along each coordinate axis, then we found a global minimizer.

More general, for  $f(x) = g(x) + \sum_{i=1}^n h_i(x_i)$  (with  $g$  convex, differentiable and each  $h_i$  convex ) we can use coordinate descent to find a minimizer: start with some initial guess  $x^{(0)}$ , and repeat

$$\begin{aligned}x_1^{(k)} &\in \operatorname{argmin}_{x_1} f\left(x_1, x_2^{(k-1)}, x_3^{(k-1)}, \dots, x_n^{(k-1)}\right) \\x_2^{(k)} &\in \operatorname{argmin}_{x_2} f\left(x_1^{(k)}, x_2, x_3^{(k-1)}, \dots, x_n^{(k-1)}\right) \\x_3^{(k)} &\in \operatorname{argmin}_{x_3} f\left(x_1^{(k)}, x_2^{(k)}, x_3, \dots, x_n^{(k-1)}\right) \\&\dots \\x_n^{(k)} &\in \operatorname{argmin}_{x_n} f\left(x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, \dots, x_n\right)\end{aligned}$$

for  $k = 1, 2, 3, \dots$

$$\partial_{\beta_j} RSS^{\text{lasso}}(\beta) = \partial_{\beta_j} RSS^{\text{OLS}}(\beta) + \partial_{\beta_j} \lambda \|\beta\|_1$$

Using subdifferential,

$$\beta_j = \beta_j(\rho_j) = \frac{1}{z_j} \text{soft}(\rho_j; \lambda),$$

where  $\text{soft}(\rho_j; \lambda) := \text{sign}(\rho_j)(|\rho_j| - \lambda)_+$ .

**Algorithm** (Coordinate descent for lasso):

- ▶ Repeat
  - ▶ for  $j = 1, \dots, d$  do
    - ▶  $z_j = 2 \sum_{i=1}^n (x_{ij})^2$
    - ▶  $\rho_j = 2 \sum_{i=1}^n x_{ij} [y_i - \beta^T x_i + \beta_j x_{ij}]$
    - ▶  $\beta_j = \text{soft}(\frac{\rho_j}{z_j}; \frac{\lambda}{z_j})$
- ▶ until convergence

## Beyond Lasso and Remarks

## Group Lasso (Yuan and Lin 2006)

- ▶ In a regression model, a multi-level categorical predictor is usually represented by a group of dummy variables.
- ▶ In an additive model, a continuous predictor may be represented by a group of basis functions to incorporate nonlinear relationship.

Suppose the  $p$  predictors are divided into  $L$  groups, with  $p_\ell$  the number in group  $\ell$ .

$$\hat{\beta}^{\text{glasso}} = \arg \min_{\beta} \left\| \mathbf{y} - \sum_{j=1}^L \mathbf{x}_j \beta_j \right\|^2 + \lambda \sum_{j=1}^L \sqrt{p_j} \|\beta_j\|$$

where

- ▶  $\lambda \geq 0$  is a tuning parameter.
- ▶  $\|\beta_j\| = \sqrt{\beta_j^T \beta_j}$  is the  $L_2$  norm.
- ▶ The group lasso penalty encourages sparsity at the factor level.
- ▶ When  $p_1 = \cdots = p_L = 1$ , the group lasso reduces to the lasso.

## Elastic net (Zhou and Hastie 2005)

$$J_{\lambda}(\boldsymbol{\beta}) = \sum_{j=1}^p [\lambda |\beta_j| + (1 - \lambda) \beta_j^2], \lambda \in [0, 1]$$

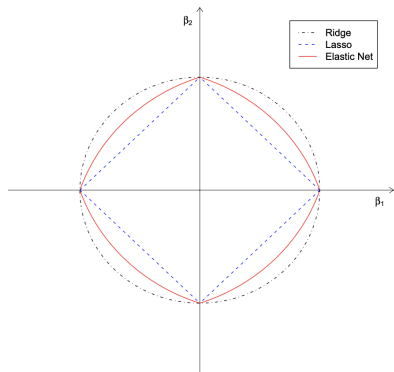


Figure: Elastic net penalty

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

$$\|\boldsymbol{\beta}\|_1 = \sum_{j=1}^p |\beta_j|, \quad \text{and} \quad \|\boldsymbol{\beta}\|_2 = \left( \sum_{j=1}^p \beta_j^2 \right)^{1/2}$$

The elastic net penalty function

- ▶ has singularity at the vertex (necessary for sparsity)
- ▶ has strict convex edges (encouraging grouping)

Then the naive Elastic Net estimator:

$$\hat{\boldsymbol{\beta}}_{\text{EN}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} L(\lambda_1, \lambda_2, \boldsymbol{\beta})$$

.

Properties: the elastic net solution

- ▶ simultaneously does automatic variable selection and continuous shrinkage
- ▶ can select groups of correlated variables, retaining ‘all the big fish’.

A correction to the double shrinkage is to use  $(1 + \lambda_2)\hat{\beta}_{\text{EN}}$  as the EN estimator.



# Bayesian interpretation of Ridge, Lasso

Model:

$$\mathbf{y} \mid \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}), \quad \beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$$

Posterior (multivariate normal regression semi-conjugate prior):

$$\beta \mid \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mu_{\beta \mid \mathbf{X}, \mathbf{y}}, \Sigma_{\beta \mid \mathbf{X}, \mathbf{y}})$$

where

$$\begin{aligned}\mu_{\beta \mid \mathbf{X}, \mathbf{y}} &= \left( \frac{\sigma^2}{\tau^2} \mathbf{I} + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \\ \Sigma_{\beta \mid \mathbf{X}, \mathbf{y}} &= (\tau^{-2} \mathbf{I} + \mathbf{X}^T \mathbf{X} / \sigma^2)^{-1}\end{aligned}$$

Obviously, MAP estimate is the same as the ridge estimate with  $\lambda = \sigma^2 / \tau^2$ .

Assume  $\beta$  has the prior distribution where the  $\beta_j$  's are independent and each having mean-zero Laplace distribution:

$$f(\beta) = \prod_{j=1}^p \frac{1}{2\tau} \exp\left(-\frac{|\beta_j|}{\tau}\right)$$

$$\begin{aligned} p(\beta \mid \mathbf{y}, \mathbf{X}) &\propto p(\beta) \cdot p(\mathbf{y} \mid \mathbf{X}, \beta) \\ &\propto \exp\left[-\frac{|\beta_1|}{\tau} - \dots - \frac{|\beta_p|}{\tau}\right] \cdot \exp\left[-(\mathbf{y} - \mathbf{X}\beta)^T \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)\right] \\ &= \exp\left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) - \frac{1}{\tau} \|\beta\|_1\right] \end{aligned}$$

The MAP is the same as

$$\underset{\beta}{\operatorname{argmin}} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \frac{2\sigma^2}{\tau} \|\beta\|_1$$

# Remark

Inference after model selection?

- ▶ the usual hypothesis test, confidence intervals become invalid if model selection is not accounted for.
- ▶ In general it is difficult and remains a very active research area, called **post-selection inference**