

# BIOS 835: Shrinkage Methods

Alexander McLain

September 7, 2023



# Introduction

- ▶ Today, we'll talk about shrinkage methods
- ▶ All shrinkage methods are biased (by design) even when the model is correct.
- ▶ They will, however, result in lower variance than the full LR model.
- ▶ Some shrinkage methods we'll discuss:
  - ▶ Ridge Regression
  - ▶ Lasso
  - ▶ Elastic Net
  - ▶ Least Angle Regression (LAR)

## Instability of LS Estimators

- ▶ Recall that  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$  can be unstable
- ▶ Further, when  $\text{rank}(\mathbf{X}) < p$   $\mathbf{X}'\mathbf{X}$  is singular
  - ▶  $\mathbf{X}$  is ill-conditioned
  - ▶ columns of  $\mathbf{X}$  are collinear
  - ▶  $p > n$
- ▶ We can measure ill-conditioning by  $\kappa = d_1/d_p$  where  $d_1$  and  $d_p$  are the largest and smallest eigenvalues, respectively, of  $\mathbf{X}$ .
- ▶ We can measure collinearity for a predictor  $X_k$  via  $VIF_k = (1 - R_k^2)^{-1}$  where  $R_k^2$  is the  $R^2$  for all other covariates regressed on  $X_k$

## Biased Regression Methods

- ▶ The instability of the LS estimates is due to  $\mathbf{X}'\mathbf{X}$  being singular or nearly singular.
- ▶ **Solution** allow  $\hat{\beta}$  to be biased.
- ▶ We'll assume that  $\mathbf{x}$  and  $\mathbf{y}$  have been centered (no  $\beta_0$ ).
- ▶ First, we discuss Principal Component Analysis (PCA) to get to PC Regression (PCR), which can be viewed as a shrinkage method.

## Introduction

- ▶ Principal components explain the variance-covariance structure of a set of variables through some linear combinations.
- ▶  $\mathbf{X}$ : variables which we want to explain the variance-covariance matrix  $\Sigma = \text{Cov}(\mathbf{X})$
- ▶  $\mathbf{Z}$ : a linear combination of the  $\mathbf{X}$ 's (this is not the response)

$$\mathbf{Z}_1 = \mathbf{X}\mathbf{v}_1 = v_{11}\mathbf{X}_1 + v_{12}\mathbf{X}_2 + \dots + v_{1p}\mathbf{X}_p$$

$$\mathbf{Z}_2 = \mathbf{X}\mathbf{v}_2 = v_{21}\mathbf{X}_1 + v_{22}\mathbf{X}_2 + \dots + v_{2p}\mathbf{X}_p$$

$$\vdots \quad \vdots \quad \vdots$$

$$\mathbf{Z}_p = \mathbf{X}\mathbf{v}_p = v_{p1}\mathbf{X}_1 + v_{p2}\mathbf{X}_2 + \dots + v_{pp}\mathbf{X}_p$$

## PC goal

- ▶ Using standard results

$$\text{Var}(\mathbf{Z}_i) = \mathbf{v}_i' \Sigma \mathbf{v}_i \quad \text{and} \quad \text{Cov}(\mathbf{Z}_i, \mathbf{Z}_k) = \mathbf{v}_i' \Sigma \mathbf{v}_k \quad \text{for } i, k = 1, 2, \dots, p$$

### Main idea of Principal Components

- ▶ Set  $\mathbf{Z}_1 = \mathbf{X} \mathbf{v}_1$  such that  $\mathbf{v}_1' \Sigma \mathbf{v}_1$  is maximized over all  $\mathbf{v}$  such that  $\mathbf{v}' \mathbf{v} = 1$ .
- ▶ Set  $\mathbf{Z}_2 = \mathbf{X} \mathbf{v}_2$  such that  $\mathbf{v}_2' \Sigma \mathbf{v}_2$  is maximized over all  $\mathbf{v}$  such that  $\mathbf{v}' \mathbf{v} = 1$  and  $\text{Cov}(\mathbf{Z}_1, \mathbf{Z}_2) = 0$ .

⋮

- ▶ Set  $\mathbf{Z}_i = \mathbf{X} \mathbf{v}_i$  such that  $\mathbf{v}_i' \Sigma \mathbf{v}_i$  is maximized over all  $\mathbf{v}$  such that  $\mathbf{v}' \mathbf{v} = 1$  and  $\text{Cov}(\mathbf{Z}_k, \mathbf{Z}_i) = 0$  for all  $k = 1, 2, \dots, i = 1$ .

⋮

## PC solution

### Theorem

Let  $\Sigma = \text{Cov}(\mathbf{X})$  where  $\mathbf{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p\}$ . Let  $\Sigma$  have eigenvalue-eigenvector pairs  $(d_1, \mathbf{v}_1), (d_2, \mathbf{v}_2), \dots, (d_p, \mathbf{v}_p)$  where  $d_1 \geq d_2 \geq \dots \geq d_p$ . Then the  $i$ th principal component is

$$\mathbf{Z}_i = \mathbf{X} \mathbf{v}_i = v_{i1} \mathbf{X}_1 + v_{i2} \mathbf{X}_2 + \dots + v_{ip} \mathbf{X}_p$$

for  $i = 1, 2, \dots, p$  where

$$\text{Var}(\mathbf{Z}_i) = \mathbf{v}_i' \Sigma \mathbf{v}_i = d_i \quad \text{and} \quad \text{Cov}(\mathbf{Z}_i, \mathbf{Z}_k) = \mathbf{v}_i' \Sigma \mathbf{v}_k = 0 \quad \text{for } i, k = 1, 2, \dots, p$$



## Properties

- ▶ Let  $\sigma_{ii}$  be the  $(i, i)$ th component of  $\Sigma$ .

$$\sum_{i=1}^p \text{Var}(\mathbf{Z}_i) = \sum_{i=1}^p \text{Var}(\mathbf{X}_i) = \sum_{i=1}^p \sigma_{ii} = \sum_{i=1}^p d_i = \text{tr}(\Sigma) = \text{tr}(\Lambda)$$

where  $\Lambda$  is a diagonal matrix of with  $i$ th element  $d_i$ , the eigenvalues.

- ▶ Recall  $d_1 \geq d_2 \geq \dots \geq d_p$  where  $\text{Var}(\mathbf{Z}_i) = d_i$ .
- ▶ The correlation between  $Z_i$  and  $X_k$  is

$$\rho_{Z_i, X_k} = \frac{v_{ik} \sqrt{d_i}}{\sqrt{\sigma_{kk}}}$$

(which can also be measured by  $v_{ik}$ ).

## Properties

- ▶ The total proportion variability that can be attributed to  $Z_k$  is

$$\frac{d_k}{\sum_{i=1}^p d_i}$$

- ▶ Suppose

$$\frac{d_1 + d_2 + d_3}{\sum_{i=1}^p d_i} = 0.97.$$

what would this imply about  $Z_4, \dots, Z_p$ ?

## How many PC's to use

- ▶ If  $\frac{\sum_{i=1}^r d_i}{\sum_{i=1}^p d_i}$  is “close to one,” then only the first  $r$  PC's are needed.
- ▶ Scree plot: a plot of  $d_i$  used for selecting  $r$
- ▶ *Kaiser's Rule* for the PC of a correlation matrix cutoff at 1.
- ▶ PC's are not invariant to scaling
- ▶ PC on standardized data = PC on a correlation matrix
- ▶ PC's are unique up to a rotation (varimax criteria).

## Connection to SVD

### Definition

The *singular value decomposition* (SVD) of an  $N \times p$  centered input matrix  $\mathbf{X}$  has the form

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}',$$

where  $\mathbf{U}$  and  $\mathbf{V}$  are  $N \times p$  and  $p \times p$  **orthogonal matrices**, and  $\mathbf{D}$  is a  $p \times p$  diagonal matrix with diagonal entries  $d_1 \geq d_2 \geq \dots \geq d_p$ . (starting to sound familiar?)

- ▶ the columns of  $\mathbf{U}$  are orthogonal unit vectors of length  $N$  called the *left singular vectors* of  $\mathbf{X}$
- ▶ the columns of  $\mathbf{V}$  are orthogonal unit vectors of length  $p$  called the *right singular vectors* of  $\mathbf{X}$
- ▶ The  $d_i$ 's are called the *singular values* of  $\mathbf{X}$

## Connection to SVD

- ▶ The sample covariance matrix of  $\mathbf{X}$  is  $\mathbf{S} = \mathbf{X}'\mathbf{X}/N$ , so using the definition of SVD we have

$$\mathbf{X}'\mathbf{X} = \mathbf{V}\mathbf{D}^2\mathbf{V}'$$

which is called the *eigen decomposition*  $\mathbf{X}'\mathbf{X}$ .

- ▶ The columns of  $\mathbf{V}$ ,  $\mathbf{v}_i$ , are equal to the eigenvectors of  $\mathbf{X}'\mathbf{X}$
- ▶ The singular values of  $\mathbf{X}$ ,  $d_i$ , are the eigenvalues of  $\mathbf{X}'\mathbf{X}$  and the diagonal values of  $\mathbf{D}^2$ .
- ▶ Note that  $\mathbf{Z} = \mathbf{X}\mathbf{V} = \mathbf{U}\mathbf{D}$  where  $\mathbf{Z}$  is an  $N \times p$  matrix of *PC scores*.
- ▶ Efficient algorithms exist to calculate the SVD of  $\mathbf{X}$ , so this is how PC's are estimated.

## Principal components regression

- ▶ Principal components regression (PCR) forms the linear combination  $\mathbf{z}_m = \mathbf{X}\mathbf{v}_m$  and then regresses  $\mathbf{y}$  on  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M$  for some  $M \leq p$ .
- ▶ Since the  $\mathbf{z}_j$ 's are orthogonal, this is just the sum of univariate regressions

$$\hat{\mathbf{y}}_{(M)}^{\text{pcr}} = \bar{y}\mathbf{1} + \sum_{m=1}^M \hat{\theta}_m \mathbf{z}_m$$

where  $\hat{\theta}_m = \langle \mathbf{z}_m, \mathbf{y} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle$ .

- ▶ The  $\mathbf{z}_m$ 's are linear combinations of the  $\mathbf{x}$ 's
- ▶ The coefficients for  $\mathbf{x}_j$  from PCR can be expressed as

$$\beta_{j(M)}^{\text{pcr}} = \sum_{m=1}^M \hat{\theta}_m v_{jm}$$

## Ridge Regression

- ▶ Potential instability in  $(\mathbf{X}'\mathbf{X})^{-1}$  could be solved by adding a constant  $\lambda$  to  $\mathbf{X}'\mathbf{X}$ .
- ▶ Adding a constant to the diagonal of  $\mathbf{X}'\mathbf{X}$  will bring it closer to a matrix with 0's on the off diagonal.
- ▶ This results in the ridge regression estimate of  $\hat{\beta}$

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

where  $\mathbf{I}$  is the  $p \times p$  identity matrix.

## Ridge Regression Properties

- ▶ Another way to view the ridge regression model is that

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \{ (\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta) + \lambda \|\beta\|^2 \},$$

where  $\|\beta\|^2 = \sum_{j=1}^p \beta_j^2$  (note: the  $\lambda$  here is the same as on the previous slide)

- ▶ (Yet) Another way to view the ridge regression model is that  $\hat{\beta}^{\text{ridge}}$  minimizes

$$ESS(\beta) = (\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta) \quad \text{such that} \quad \|\beta\|^2 \leq t,$$

which makes explicit the size constraint on the parameters.

- ▶ The relationship between  $t$  and  $\lambda$  is made through Lagrange multipliers.



## Ridge Regression and Bayes

The ridge regression estimate naturally comes up in Bayesian regression analysis.

- ▶ If we were to estimate the model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$  in a Bayesian framework with
  - ▶  $\mathbf{e} \sim MVN(\mathbf{0}, \sigma^2 \mathbf{I})$ , and
  - ▶  $\boldsymbol{\beta} \sim MVN(\mathbf{0}, \sigma_\beta^2 \mathbf{I})$  (prior distribution).
- ▶ Then the mean (or mode) of the posterior distribution of  $\boldsymbol{\beta} | \mathbf{Y}$  is  $\hat{\boldsymbol{\beta}}^{\text{ridge}}$  with  $\lambda = \sigma^2 / \sigma_\beta^2$ .

## Ridge Regression and PCR

If  $\mathbf{Z} = \mathbf{V}\mathbf{X}$  and  $\boldsymbol{\alpha} = \mathbf{V}'\boldsymbol{\beta}$  (thus  $\boldsymbol{\beta} = \mathbf{V}\boldsymbol{\alpha}$ )

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{Z}\boldsymbol{\alpha} + \boldsymbol{\epsilon} \quad (1)$$

- If we were to apply a ridge regression PCR model, the  $j$ th coefficient is

$$\hat{\alpha}_j^{\text{ridge}} = \frac{d_j^2}{d_j^2 + \lambda} \hat{\alpha}_j$$

where  $\hat{\alpha}$  is the LS estimate of (1), and  $d_j$  is the  $j$ th eigenvalue.

## Ridge Regression and PCR

If  $\mathbf{Z} = \mathbf{V}\mathbf{X}$  and  $\boldsymbol{\alpha} = \mathbf{V}'\boldsymbol{\beta}$  (thus  $\boldsymbol{\beta} = \mathbf{V}\boldsymbol{\alpha}$ )

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{Z}\boldsymbol{\alpha} + \boldsymbol{\epsilon} \quad (1)$$

- If we were to apply a ridge regression PCR model, the  $j$ th coefficient is

$$\hat{\alpha}_j^{\text{ridge}} = \frac{d_j^2}{d_j^2 + \lambda} \hat{\alpha}_j$$

where  $\hat{\boldsymbol{\alpha}}$  is the LS estimate of (1), and  $d_j$  is the  $j$ th eigenvalue.

- Then  $\hat{\boldsymbol{\beta}}^{\text{ridge}} = \mathbf{V}\hat{\boldsymbol{\alpha}}^{\text{ridge}}$ .
- Thus, the shrinkage of  $\hat{\boldsymbol{\beta}}^{\text{ridge}}$  compared to  $\hat{\boldsymbol{\beta}}$  is a function of the singular values.

## Ridge Regression Degrees of Freedom

- ▶ Degrees of freedom are a good measure of the flexibility of the model.
- ▶ The *effective degrees of freedom* for the ridge model are

$$\begin{aligned}df(\lambda) &= \text{tr}\{\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\}, \\&= \text{tr}(\mathbf{H}_\lambda) \\&= \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}\end{aligned}$$

- ▶ Note that when  $\lambda = 0$  we have  $df(\lambda) = p$ , and  $df(\lambda) \rightarrow 0$  as  $\lambda \rightarrow \infty$ .
- ▶ All ridge coefficients will be non-zero.

## Ridge Regression Properties

- The bias and variance of Ridge coefficients have nice closed forms

$$\text{Var}_\lambda(\hat{\beta}^{\text{ridge}}) = \sigma^2 \sum_{j=1}^p \frac{d_j^2}{(d_j^2 + \lambda)^2}$$

$$\text{Bias}_\lambda^2(\hat{\beta}^{\text{ridge}}) = \lambda^2 \sum_{j=1}^p \frac{(\hat{\alpha}_j^{\text{ridge}})^2}{(d_j^2 + \lambda)^2}$$

thus,

$$\text{MSE}_\lambda(\hat{\beta}^{\text{ridge}}) = \sigma^2 \sum_{j=1}^p \frac{\sigma^2 d_j^2 + \lambda (\hat{\alpha}_j^{\text{ridge}})^2}{(d_j^2 + \lambda)^2}$$

## Lasso regression model

- ▶ The Lasso regression model, results in solution  $\hat{\beta}^{\text{lasso}}$  minimizes

$$ESS(\beta) = (\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta) \quad \text{such that} \quad \sum_{i=1}^p |\beta_j| \leq t,$$

which makes explicit the size constraint on the parameters.

- ▶ The solution to this is such that

$$\hat{\beta}^{\text{lasso}} = \operatorname{argmin}_{\beta} \left\{ (\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta) + \lambda \sum_{i=1}^p |\beta_j| \right\},$$

- ▶ The relationship between these forms is made through Lagrange multipliers.

## Lasso properties

- ▶ A major difference between the Lasso and Ridge models is that making  $t$  sufficiently small (or  $\lambda$  large) will result in some  $\beta$  coefficients being exactly zero.
- ▶ Let  $t_0 = \sum_{j=1}^p |\hat{\beta}_j^{ls}|$  then if  $t = t_0/2$  the Lasso will result in coefficients that are the least squares coefficients shrunk by about 50% on average.
- ▶ There is no closed-form solution for Lasso coefficients (they are nonlinear in  $\mathbf{y}$ ).

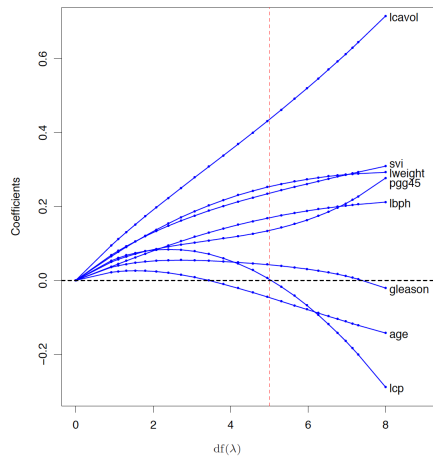
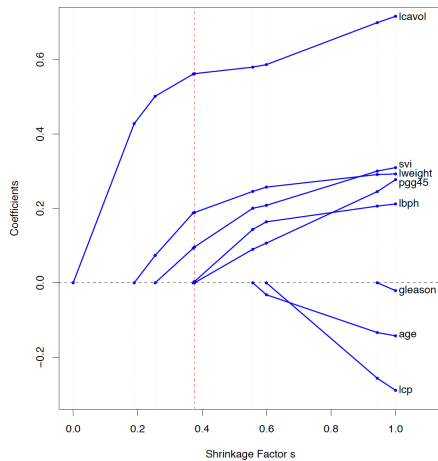


Figure: From ESL, where  $s = t / \sum_1^p |\hat{\beta}_j^{ls}|$



## Ridge vs Lasso vs Elastic Net

- The Lasso regression model,

$$\hat{\beta}^{\text{lasso}} = \operatorname{argmin}_{\beta} \left\{ (\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta) + \lambda \sum_{i=1}^p |\beta_i| \right\},$$

- The ridge regression model,

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \left\{ (\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta) + \lambda \sum_{i=1}^p \beta_i^2 \right\},$$

- The elastic net regression model (Zou and Hastie, 2005),

$$\hat{\beta}^{\text{en}} = \operatorname{argmin}_{\beta} \left[ (\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta) + \lambda \sum_{i=1}^p \{ \alpha \beta_i^2 + (1 - \alpha) |\beta_i| \} \right],$$

## Ridge vs Lasso vs Best Subset

- Suppose  $\mathbf{X}$  is a standardized orthonormal input matrix
- The estimators for Ridge, Lasso, and best subset, denoted by  $\hat{\beta}^{\text{ridge}}$ ,  $\hat{\beta}^{\text{lasso}}$  and  $\hat{\beta}^{\text{bss}}$  respectively, are

$$\hat{\beta}_j^{\text{bss}} = \hat{\beta}_j^{\text{ls}} I \left( |\hat{\beta}_j| \geq |\hat{\beta}_{(M)}| \right) \text{ of size } M$$

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

$$\hat{\beta}_j^{\text{lasso}} = \text{sign} \left( \hat{\beta}_j^{\text{ls}} \right) \left( |\hat{\beta}_j^{\text{ls}}| - \lambda \right)_+$$

where  $\hat{\beta}_{(M)}$  is the  $M$ th largest regression coefficient.

## Least Angle Regression

- ▶ Least Angle Regression (LAR) was developed in Efron *et al.* (2004).<sup>1</sup>
- ▶ LAR starts by entering the variable most closely correlated with the response. LAR moves the coefficient (continuously) towards the LS estimate until another variable has an equal correlation with the residuals of the model.
- ▶ Then, both variables are moved toward their LS estimates (while keeping their correlation tied and decreasing) until another variable has equal correlation with the residuals. Etc.
- ▶ LAR is similar to the forward stagewise regression in that only “some” of a variable’s influence on the outcome is entered into the model.

---

<sup>1</sup>Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R., 2004. Least angle regression. *The Annals of statistics*, 32(2), pp.407-499.

## LAR Algorithm

### LAR algorithm

1. Standardize predictors. Start with residual  $\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}}$  and  $\beta = 0$ .
2. Find the predictor  $\mathbf{x}_j$  most correlated with  $\mathbf{r}$ .
3. Move  $\beta_j$  from  $0 \rightarrow \langle \mathbf{x}_j, \mathbf{r} \rangle$  until another variable  $\mathbf{x}_k$  has equal absolute correlation with the residual.
4. Move  $(\beta_j, \beta_k)$  towards their joint least squares estimates of the current  $\mathbf{r}$  on  $(\mathbf{x}_j, \mathbf{x}_k)$ , until another variable  $\mathbf{x}_l$  has equal absolute correlation with the residual.
5. Continue until all  $p$  predictors have been entered.

\*LAR can be used when  $p > n$

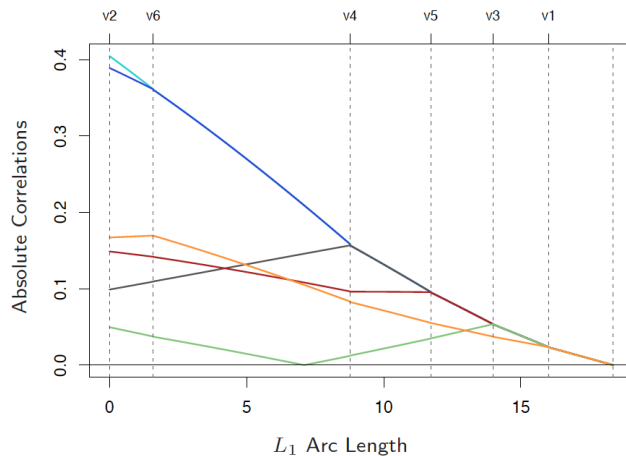


Figure: From ESL (online version)

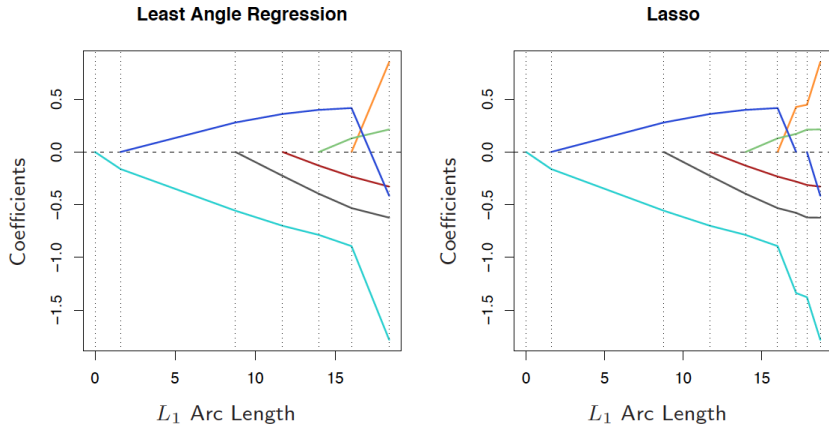


Figure: From ESL (online version), note the similarities.

## Modified LAR Algorithm

- ▶ A Modified LAR algorithm is done by adding
  - 4a. If a non-zero coefficient hits zero, drop its variable from the active set and recompute the current joint least squares direction
- ▶ The Modified LAR algorithm is an extremely efficient way of estimating the LASSO model (recall the need for quadratic programming with LASSO).
- ▶ The current joint least squares direction and the point where another variable will have equal absolute correlation can be calculated (the variables don't actually have to be moved continuously).

## Partial Least Squares

- ▶ PCR was performed by making a linear combination of the  $\mathbf{X}$ 's, then performing LR on the derived linear combinations.
- ▶ Partial Least Squares (PLS), like PCR, also constructs a linear combination of the  $\mathbf{X}$ 's and performs LR on the derived linear combinations.
- ▶ Note that when performing PCA (for PCR), the outcome ( $\mathbf{y}$ ) is never used.
- ▶ PLS does use the outcome to form the linear combinations.



# Partial Least Squares Algorithm

## PLS Algorithm

1. Standardize each  $\mathbf{x}_j$  and set  $\hat{\mathbf{y}}^{(0)} = \bar{\mathbf{y}}\mathbf{1}$ , and  $\hat{\mathbf{x}}_j^{(0)} = \mathbf{x}_j$ ,  $j = 1, 2, \dots, p$ .
2. For  $m = 1, 2, \dots, M \leq p$ 
  - 2.1  $\mathbf{z}_m = \sum_{j=1}^p \hat{\psi}_{mj} \mathbf{x}_j^{(m-1)}$ , where  $\hat{\psi}_{mj} = \langle \mathbf{x}_j^{(m-1)}, \mathbf{y} \rangle$ .
  - 2.2  $\hat{\theta}_m = \langle \mathbf{z}_m, \mathbf{y} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle$ .
  - 2.3  $\hat{\mathbf{y}}^{(m)} = \hat{\mathbf{y}}^{(m-1)} + \hat{\theta}_m \mathbf{z}_m$ .
  - 2.4 Orthogonalize each  $\mathbf{x}_j^{(m-1)}$  with respect to  $\mathbf{z}_m$ . That is, set
 
$$\mathbf{x}_j^{(m)} = \mathbf{x}_j^{(m-1)} - (\langle \mathbf{z}_m, \mathbf{x}_j^{(m-1)} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle) \mathbf{z}_m, \text{ for } j = 1, 2, \dots, p.$$
3. Output  $\{\hat{\mathbf{y}}^{(m)}\}_1^M$ . Since the  $\{\mathbf{z}_l\}_1^m$  are linear in  $\mathbf{x}_j$ , so is  $\hat{\mathbf{y}}^{(m)} = \mathbf{X} \hat{\boldsymbol{\beta}}^{\text{pls}}(m)$ .

## Partial Least Squares Algorithm (part 2)

1. Standardize each  $\mathbf{x}_j$  and set  $\hat{\mathbf{y}}^{(0)} = \mathbf{y} - \bar{\mathbf{y}}\mathbf{1}$ , and  $\hat{\mathbf{x}}_j^{(0)} = \mathbf{x}_j$ ,  $j = 1, 2, \dots, p$ .
2. For  $m = 1, 2, \dots, M \leq p$ 
  - 2.1 Compute  $\hat{\psi}_{mj} = \langle \mathbf{x}_j^{(m-1)}, \mathbf{y}^{(m-1)} \rangle$  and set  $\mathbf{z}_m = \sum_{j=1}^p \hat{\psi}_{mj} \mathbf{x}_j^{(m-1)}$ .
  - 2.2 Compute  $\hat{\theta}_m = \langle \mathbf{z}_m, \mathbf{y}^{(m-1)} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle$  and set  $\hat{\mathbf{y}}^{(m)} = \hat{\mathbf{y}}^{(m-1)} - \hat{\theta}_m \mathbf{z}_m$ .
  - 2.3 Regress  $\mathbf{x}_j^{(m-1)}$  with respect to  $\mathbf{z}_m$  for all  $j$ . That is, set

$$\alpha_{mj} = \frac{\langle \mathbf{z}_m, \mathbf{x}_j^{(m-1)} \rangle}{\langle \mathbf{z}_m, \mathbf{z}_m \rangle}$$

and then  $\mathbf{x}_j^{(m)} = \mathbf{x}_j^{(m-1)} - \alpha_{mj} \mathbf{z}_m$ , for  $j = 1, 2, \dots, p$ .

3. Then the PLSR model with  $m$  components is

$$\hat{\mathbf{y}}(k) = \bar{\mathbf{y}}\mathbf{1} + \sum_{m=1}^k \hat{\theta}_m \mathbf{z}_m.$$

## Grouped Lasso regression model

- ▶ In some problems, the predictors belong to pre-defined groups (e.g., genes belonging to some biological pathway or dummy-coded categorical variables).
- ▶ In this situation it may be desirable to shrink and select the members of a group together.
- ▶ Suppose that the  $p$  predictors are divided into  $L$  groups, with  $p_\ell$  being the number of variables in group  $\ell$
- ▶ The Grouped Lasso regression model minimizes

$$(\mathbf{y} - \mathbf{x}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{x}\boldsymbol{\beta}) + \lambda \sum_{\ell=1}^L \sqrt{p_\ell} \|\boldsymbol{\beta}_\ell\|_2,$$

where  $\boldsymbol{\beta}_\ell$  are the regression coefficients for group  $\ell$ , and  $\|\cdot\|_2$  is the Euclidean norm.

## Fused Lasso regression model

- ▶ Along with grouped methods there is the idea of fusion or merging.
- ▶ In this situation it may be desirable to shrink and select the members of a group together.
- ▶ Consider the categorical covariates; by assuming the coefficient for one group is zero, you fuse that group with the referent group.
- ▶ The idea of fusion is to fuse groups together and not have the fusion be subject to the reference group.
- ▶ The Fused Lasso regression proposed by Tibshirani et al (2005) minimizes

$$\hat{\beta} = \arg \min \{(\mathbf{y} - \mathbf{x}\beta)'(\mathbf{y} - \mathbf{x}\beta)\}$$

such that

$$\sum_{j=1}^p |\beta_j| \leq s_1 \quad \text{and} \quad \sum_{j=2}^p |\beta_j - \beta_{j-1}| \leq s_2.$$