

# Ch 3: Linear Methods for Regression

## Advanced Statistical Data Mining

Jaejik Kim

Department of Statistics  
Sungkyunkwan University

Fall 2020

# Linear Regression Models

- ▶ The linear regression model:  $f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j$ .
- ▶  $f$ : Regression function  $E[Y|X]$ .
- ▶ Sources of inputs  $X_j$ :
  - ▶ qualitative inputs,
  - ▶ transformations of quantitative inputs,
  - ▶ basis expansions leading to polynomial representation,
  - ▶ numeric or “dummy” coding of qualitative inputs,
  - ▶ interactions between variables such as  $X_3 = X_1 \cdot X_2$ .
- ▶ The model is linear in the parameters  $\beta_j$ ,  $j = 1, \dots, p$ .

# Least Squares Estimation

- ▶ The least squares method minimizes the residual sum of squares (RSS)

$$RSS(\beta) = \sum_{i=1}^N (y_i - f(x_i))^2 = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta)$$

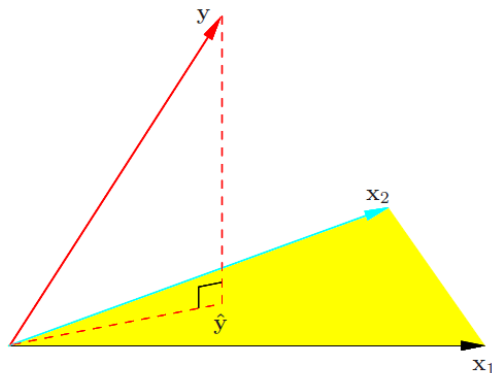
- ▶  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

- ▶ Fitted values of training inputs:  $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{H}\mathbf{y}$

- ▶  $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  (*hat matrix*).

- ▶  $\hat{\mathbf{y}}$ : The orthogonal projection of  $\mathbf{y}$  onto the column space of  $\mathbf{X}$ .

# Orthogonal Projection



**FIGURE 3.2.** The  $N$ -dimensional geometry of least squares regression with two predictors. The outcome vector  $y$  is orthogonally projected onto the hyperplane spanned by the input vectors  $x_1$  and  $x_2$ . The projection  $\hat{y}$  represents the vector of the least squares predictions

# Basic Least Squares Assumptions

- ▶ So far, no distributional assumptions have been made.
- ▶ Now we make some assumptions:
  - ▶  $x_i$ 's are fixed.
  - ▶  $\text{Cov}(Y_i, Y_j) = \begin{cases} \sigma^2 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$
- ▶  $\hat{\beta}$ : Unbiased estimator of  $\beta$ ; i.e.,  $E[\hat{\beta}] = \beta$ .
- ▶  $\text{Var}(\hat{\beta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2$
- ▶  $\hat{\sigma}^2 = \frac{RSS(\hat{\beta})}{N - p - 1}$ ;  $E[\hat{\sigma}^2] = \sigma^2$

# Multivariate Expectations and Variances

- ▶  $E \left[ \sum_{i=1}^N a_i Y_i \right] = \sum_{i=1}^N a_i E[Y_i] \iff E[\mathbf{a}^\top \mathbf{Y}] = \mathbf{a}^\top E[\mathbf{Y}]$
- ▶  $\text{Var} \left[ \sum_{i=1}^N a_i Y_i \right] = \sum_{i=1}^N \sum_{j=1}^N a_i a_j \text{Cov}[Y_i, Y_j]$   
 $\iff \text{Var}[\mathbf{a}^\top \mathbf{Y}] = \mathbf{a}^\top \text{Var}[\mathbf{Y}] \mathbf{a}$
- ▶  $\text{Cov} \left[ \sum_{i=1}^N a_i Y_i, \sum_{j=1}^N b_j Y_j \right] = \sum_{i=1}^N \sum_{j=1}^N a_i b_j \text{Cov}[Y_i, Y_j]$   
 $\iff \text{Cov}[\mathbf{a}^\top \mathbf{Y}, \mathbf{b}^\top \mathbf{Y}] = \mathbf{a}^\top \text{Var}[\mathbf{Y}] \mathbf{b}$
- ▶  $E[\mathbf{A}\mathbf{Y}] = \mathbf{A} E[\mathbf{Y}]$
- ▶  $\text{Var}[\mathbf{A}\mathbf{Y}] = \mathbf{A} \text{Var}[\mathbf{Y}] \mathbf{A}^\top$
- ▶  $\text{Cov}[\mathbf{A}\mathbf{Y}, \mathbf{B}\mathbf{Y}] = \mathbf{A} \text{Var}[\mathbf{Y}] \mathbf{B}^\top$

# Gauss-Markov Theorem

- ▶ Consider a linear combination of the form  $a^\top \hat{\beta}$  for an arbitrary  $a \in \mathbb{R}^{p+1}$ .
- ▶ LSE of  $a^\top \beta = a^\top \hat{\beta} = a^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .
- ▶  $E[a^\top \hat{\beta}] = a^\top \beta$
- ▶ **Gauss-Markov Theorem:** The LSE  $a^\top \hat{\beta}$  has variance no bigger than that of any other linear unbiased estimator  $\tilde{\theta} = \mathbf{c}^\top \mathbf{Y}$  that is unbiased for  $a^\top \beta$ .
- ▶  $\text{MSE}(\tilde{\theta}) = E(\tilde{\theta} - \theta)^2 = \text{Var}(\tilde{\theta}) + [E(\tilde{\theta}) - \theta]^2 \Rightarrow$  Gauss-Markov Theorem implies that the LSE has the smallest MSE (variance) of all unbiased linear estimators.
- ▶ Conversely, there may exist biased estimators which have smaller MSE than the LSE.

# Inferential Assumptions

- ▶ Inference  $\Rightarrow$  Stronger assumptions:
  - ▶  $Y = E(Y|X_1, \dots, X_p) + \epsilon = \beta_0 + \sum_{j=1}^p X_j \beta_j + \epsilon$
  - ▶  $\epsilon \sim \text{Normal}(0, \sigma^2)$
- ▶  $\hat{\beta} \sim \text{Normal}(\beta, (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2)$
- ▶  $\frac{(N - p - 1) \hat{\sigma}^2}{\sigma^2} \sim \chi^2_{N-p-1}$
- ▶  $\hat{\beta}$  &  $\hat{\sigma}^2$  are independent



# Hypothesis Testing

- ▶ Since the marginal distribution of a multivariate Normal is a univariate Normal,  $\hat{\beta}_j \sim \text{Normal}(\beta_j, \sigma^2 v_j)$  where  $v_j$  is the  $j$ th diagonal element of  $(\mathbf{X}^\top \mathbf{X})^{-1}$ .
- ▶  $\frac{\hat{\beta}_j - \beta_j}{\sigma \sqrt{v_j}} \sim \text{Normal}(0, 1)$
- ▶  $z_j = \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \sqrt{v_j}} \sim t_{N-p-1}$
- ▶ Reject  $H_0$  from  $H_0 : \beta_j = 0$  vs.  $H_a : \beta_j \neq 0$  if  $|z_j| > t_{N-p-1, (1-\alpha/2)}$ .

# Hypothesis Testing

- ▶ Test for a group of coefficients:

$$F = \frac{(RSS_0 - RSS_1)/(p_1 - p_0)}{RSS_1/(N - p_1 - 1)} \sim F_{p_1 - p_0, N - p_1 - 1}.$$

- ▶  $RSS_1$ : The RSS for the full model ( $p_1 + 1$  parameters):

$$Y = \beta_0 + \sum_{j=1}^{p_1} X_j \beta_j + \epsilon$$

- ▶  $RSS_0$ : RSS for the reduced model ( $p_0 + 1$  parameters,

$$p_0 < p_1): Y = \beta_0 + \sum_{j=1}^{p_0} X_j \beta_j + \epsilon$$

- ▶ Reject  $H_0$  from  $H_0 : \beta_{p_0+1} = \dots = \beta_{p_1} = 0$  vs.  $H_a$  : “ $H_0$  is not true” if  $F > F_{p_1 - p_0, N - p_1 - 1, (1-\alpha)}$ .

# Simple Regression and Orthogonality

- ▶ Univariate regression model with no intercept where  $X \in \mathbb{R}$ :

$$Y = X\beta + \epsilon.$$

- ▶ Given data  $(x_1, y_1), \dots, (x_N, y_N)$ ,

$$\hat{\beta} = \frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2} = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}.$$

- ▶  $\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{x}\hat{\beta}$ .
- ▶  $\hat{\mathbf{y}}$ : Projection of  $\mathbf{y}$  onto  $\mathbf{x}$ .
- ▶  $\mathbf{r}$ : Orthogonal complement of that projection.

- ▶ Multivariate regression: If the inputs are orthogonal

$$(\langle \mathbf{x}_j, \mathbf{x}_k \rangle = 0 \text{ if } j \neq k), \text{ then } \hat{\beta}_j = \frac{\langle \mathbf{x}_j, \mathbf{y} \rangle}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle} \text{ for all } j.$$

$\Rightarrow$  No effect on each other's parameter estimates in the model.

# Gram-Schmidt Orthogonalization

## Theorem Gram-Schmidt Orthogonalization

Let  $V$  be an inner product space and  $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$  be a linearly independent subset of  $V$ . Define  $S' = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_p\}$ , where  $\mathbf{z}_1 = \mathbf{x}_1$  and

$$\mathbf{z}_j = \mathbf{x}_j - \sum_{k=1}^{j-1} \frac{\langle \mathbf{z}_k, \mathbf{x}_j \rangle}{\langle \mathbf{z}_k, \mathbf{z}_k \rangle} \mathbf{z}_k, \quad 2 \leq j \leq p.$$

Then,  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_p$  are orthogonal (orthogonal basis of  $V$ ).

$$\text{span}(S) = \text{span}(S').$$

$\Rightarrow$  Orthogonalization does NOT change the subspace spanned by  $\mathbf{x}_j$ 's.

# Logic of Orthogonalization

- ▶ Simple linear regression:  $Y = \beta_0 + \beta_1 X + \epsilon$ .

- ▶ LSE  $\hat{\beta}_1 = \frac{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{y} \rangle}{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{x} - \bar{x}\mathbf{1} \rangle}$ .

Step 1  $\mathbf{x} = \gamma_0 \mathbf{1} \Rightarrow \hat{\gamma}_0 = \frac{\langle \mathbf{1}, \mathbf{x} \rangle}{\langle \mathbf{1}, \mathbf{1} \rangle} = \bar{x} \Rightarrow \mathbf{z} = \mathbf{x} - \bar{x}\mathbf{1}$ .

Step 2  $\mathbf{y} = \gamma_1 \mathbf{z} \Rightarrow \hat{\gamma}_1 = \frac{\langle \mathbf{z}, \mathbf{y} \rangle}{\langle \mathbf{z}, \mathbf{z} \rangle} = \frac{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{y} \rangle}{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{x} - \bar{x}\mathbf{1} \rangle} = \hat{\beta}_1$ .

- ▶ Multiple regression:  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$ .

Step 1  $\mathbf{x}_1 = \gamma_{01} \mathbf{1} \Rightarrow \hat{\gamma}_{01} = \frac{\langle \mathbf{1}, \mathbf{x}_1 \rangle}{\langle \mathbf{1}, \mathbf{1} \rangle} \Rightarrow \mathbf{z}_1 = \mathbf{x}_1 - \hat{\gamma}_{01} \mathbf{1}$ .

Step 2  $\mathbf{x}_2 = \gamma_{12} \mathbf{z}_1 + \gamma_{02} \mathbf{1} \Rightarrow \hat{\gamma}_{12} = \frac{\langle \mathbf{z}_1, \mathbf{x}_2 \rangle}{\langle \mathbf{z}_1, \mathbf{z}_1 \rangle}, \hat{\gamma}_{02} = \frac{\langle \mathbf{1}, \mathbf{x}_2 \rangle}{\langle \mathbf{1}, \mathbf{1} \rangle}$   
 $\Rightarrow \mathbf{z}_2 = \mathbf{x}_2 - \hat{\gamma}_{12} \mathbf{z}_1 - \hat{\gamma}_{02} \mathbf{1}$ .

Step 3  $\mathbf{y} = \gamma_2 \mathbf{z}_2 \Rightarrow \hat{\gamma}_2 = \frac{\langle \mathbf{z}_2, \mathbf{y} \rangle}{\langle \mathbf{z}_2, \mathbf{z}_2 \rangle} = \hat{\beta}_2$ .

# Gram-Schmidt Orthogonalization

- ▶ If the inputs are not orthogonal, we can orthogonalize the inputs and obtain the coefficient estimate for the last input.

- ▶ Gram-Schmidt Procedure:

$$0. \mathbf{z}_0 = \mathbf{x}_0 = \mathbf{1}$$

$$1. \mathbf{z}_1 = \mathbf{x}_1 - \frac{\langle \mathbf{z}_0, \mathbf{x}_1 \rangle}{\langle \mathbf{z}_0, \mathbf{z}_0 \rangle} \mathbf{z}_0$$

$$2. \mathbf{z}_2 = \mathbf{x}_2 - \frac{\langle \mathbf{z}_0, \mathbf{x}_2 \rangle}{\langle \mathbf{z}_0, \mathbf{z}_0 \rangle} \mathbf{z}_0 - \frac{\langle \mathbf{z}_1, \mathbf{x}_2 \rangle}{\langle \mathbf{z}_1, \mathbf{z}_1 \rangle} \mathbf{z}_1$$

$$\vdots$$

$$p. \mathbf{z}_p = \mathbf{x}_p - \sum_{j=0}^{p-1} \frac{\langle \mathbf{z}_j, \mathbf{x}_p \rangle}{\langle \mathbf{z}_j, \mathbf{z}_j \rangle} \mathbf{z}_j$$

- ▶ Then  $\hat{\beta}_p = \frac{\langle \mathbf{z}_p, \mathbf{y} \rangle}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle}$ .

# Meaning of G-S Orthogonalization

- ▶ Any one of  $\mathbf{x}_j$ ,  $j = 1, \dots, p$  can be in the last position.  $\Rightarrow$  All  $\beta_j$  can be estimated by the G-S orthogonalization.
- ▶  $\hat{\beta}_j$ : Simple regression coefficient of  $\mathbf{y}$  on  $\mathbf{z}^*$ , where  $\mathbf{z}^*$  is the residual after regressing

$$\mathbf{x}_j = \gamma_0 \mathbf{x}_0 + \gamma_1 \mathbf{x}_1 + \dots + \gamma_{j-1} \mathbf{x}_{j-1} + \gamma_{j+1} \mathbf{x}_{j+1} + \dots + \gamma_p \mathbf{x}_p.$$

$\Rightarrow \hat{\beta}_j$ : **Additional contribution** of  $\mathbf{x}_j$  on  $\mathbf{y}$ , after  $\mathbf{x}_j$  has been adjusted for  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{x}_{j+1}, \dots, \mathbf{x}_p$ .

# Problem of Least Squares

- ▶ When  $\mathbf{x}_p$  (the observed values for the  $p$ th variable) is highly correlated with some of the other  $\mathbf{x}_k$ 's, the residual vector  $\mathbf{z}_p$  will be close to 0.

- ▶ 
$$\text{Var}(\hat{\beta}_p) = \text{Var}\left(\frac{\mathbf{z}_p^\top \mathbf{y}}{\mathbf{z}_p^\top \mathbf{z}_p}\right) = \frac{\mathbf{z}_p^\top (\sigma^2 \mathbf{I}) \mathbf{z}_p}{(\mathbf{z}_p^\top \mathbf{z}_p)^2} = \frac{\sigma^2 \mathbf{z}_p^\top \mathbf{z}_p}{(\mathbf{z}_p^\top \mathbf{z}_p)^2} = \frac{\sigma^2}{\|\mathbf{z}_p\|^2}$$

- ▶ If  $\mathbf{x}_p$  and other inputs are highly correlated,  $\|\mathbf{z}_p\| \downarrow \Rightarrow \text{Var}(\hat{\beta}_p) \uparrow$ .



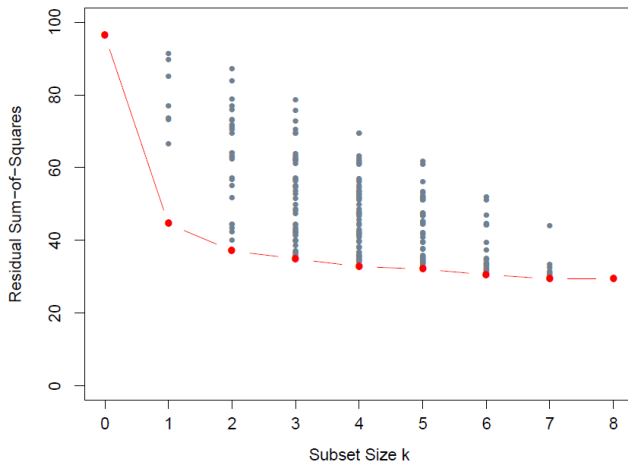
# Alternatives to Least Squares

- ▶ There are two reasons we may not be satisfied with least squares estimates:
  - ▶ *Prediction Accuracy*: LS estimates often have large variance (when predictors are correlated). A little bit of bias can be added to reduce the variance of the predicted values.
  - ▶ *Interpretation*: Often we would like to determine a smaller subset of variables that exhibit the strongest effects. To get the “big picture”, some small details can be sacrificed.
- ▶ Three methods to fix these problems:
  - ▶ Subset Selection
  - ▶ Coefficient Shrinkage
  - ▶ Methods Using Derived Input Directions

# Subset Selection

- ▶ *Best Subset Selection*: For each  $k \in \{0, 1, \dots, p\}$ , RSS values of all possible models are computed.  $\Rightarrow$  Model with the smallest RSS at each  $k$ .  $\Rightarrow$  Final model determined by AIC, BIC, or cross-validation.  $\Rightarrow$  The optimal model. If  $p$  is large, then this may be infeasible.
- ▶ *Forward Selection*: Start with the intercept, then sequentially add inputs that most improves the fit until the last input is not significant.  $\Rightarrow$  greedy algorithm  $\Rightarrow$  suboptimal (No problem for large  $p$  even when  $p \geq N$ ).
- ▶ *Backward Selection*: Start with the full model, then sequentially delete inputs that contribute least to the RSS until all inputs are significant (Only available when  $N > p$ ).
- ▶ *Stepwise Regression*: Consider forward and backward moves at each stage.  $\Rightarrow$  For large  $p$ , computationally expensive.

# Best subset



**FIGURE 3.5.** All possible subset models for the prostate cancer example. At each subset size is shown the residual sum-of-squares for each model of that size.

- ▶ Subset selection  $\Rightarrow$  Discrete process (retain or discard)  $\Rightarrow$  Sometimes high model variance.
- ▶ Shrinkage methods  $\Rightarrow$  More continuous process.
- ▶ Shrinkage methods:
  - ▶ Ridge regression.
  - ▶ Lasso, Elastic net, SCAD.
  - ▶ Least Angle Regression (LAR).

# Ridge Regression

- ▶ Another way to control the variance is to impose a size constraint on the coefficients.
- ▶ Each input variable and the output should be **standardized**; that is  $\bar{y} = 0$  and  $\sum_{i=1}^N x_{ij} = 0$  &  $\sum_{i=1}^N x_{ij}^2 / (N - 1) = 1$ .
- ▶ Ridge regression minimizes  $RSS(\beta)$  subject to the constraint  $\sum_{j=1}^p \beta_j^2 \leq s$ .
- ▶ This is equivalent to minimizing the penalized RSS:

$$PRSS(\beta; \lambda) = \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2 + \lambda \sum_{j=1}^p \beta_j^2.$$

with respect to  $\beta$ , considering  $\lambda$  to be fixed.

$$PRSS(\beta; \lambda) = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^\top \beta$$

$$\frac{\partial PRSS}{\partial \beta} = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\beta}^{\text{ridge}}) + 2\lambda \hat{\beta}^{\text{ridge}} = 0$$

$$0 = -\mathbf{X}^\top \mathbf{y} + \mathbf{X}^\top \mathbf{X} \hat{\beta}^{\text{ridge}} + \lambda \mathbf{I} \hat{\beta}^{\text{ridge}}$$

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

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

# Singular Value Decomposition

- ▶ The singular value decomposition (SVD) of  $\mathbf{X} \in \mathbb{R}^{N \times p}$  has the form

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top.$$

- ▶  $\mathbf{U}$  is a  $N \times p$  orthogonal matrix with columns which span the column space of  $\mathbf{X}$ .
- ▶  $\mathbf{V}$  is a  $p \times p$  orthogonal matrix with columns which span the row space of  $\mathbf{X}$ .
- ▶  $\mathbf{D}$  is a  $p \times p$  diagonal matrix with diagonal entries  $d_1 \geq d_2 \geq \dots \geq d_r > d_{r+1} = \dots = d_p = 0$  where  $r = \text{rank}(\mathbf{X})$ . The entries  $d_1, \dots, d_p$  are called the *singular values* of  $\mathbf{X}$ .

- ▶ Applying the singular value decomposition of  $\mathbf{X}$  to the least square regression:

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{X}\hat{\boldsymbol{\beta}}^{ls} \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{UDV}^\top [(\mathbf{UDV}^\top)^\top \mathbf{UDV}^\top]^{-1} (\mathbf{UDV}^\top)^\top \mathbf{y} \\ &= \mathbf{UU}^\top \mathbf{y}.\end{aligned}$$

- ▶  $\mathbf{U}^\top \mathbf{y}$  is the coordinates of  $\mathbf{y}$  w.r.t. the orthogonal basis  $\mathbf{U}$ .



# SVD and Ridge Regression

- ▶ Applying the singular value decomposition of  $\mathbf{X}$  to the ridge regression formulas:

$$\begin{aligned}\hat{\beta}^{\text{ridge}} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{V}(\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^\top \mathbf{y} \\ &= \sum_{j=1}^p v_j \frac{d_j}{d_j^2 + \lambda} \mathbf{u}_j^\top \mathbf{y}\end{aligned}$$

$$\begin{aligned}\hat{\mathbf{y}}^{\text{ridge}} &= \mathbf{X} \hat{\beta}^{\text{ridge}} \\ &= \mathbf{U} \mathbf{D} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^\top \mathbf{y} \\ &= \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^\top \mathbf{y}\end{aligned}$$

- ▶ Since  $\lambda \geq 0$ ,  $\frac{d_j^2}{d_j^2 + \lambda} \leq 1$ .
- ▶ Like the LSE, ridge regression computes the coordinates of  $\mathbf{y}$  w.r.t. the orthogonal basis  $\mathbf{U}$ . But, these coordinates are shrunk by  $\frac{d_j^2}{d_j^2 + \lambda}$ .

# SVD and Ridge Regression

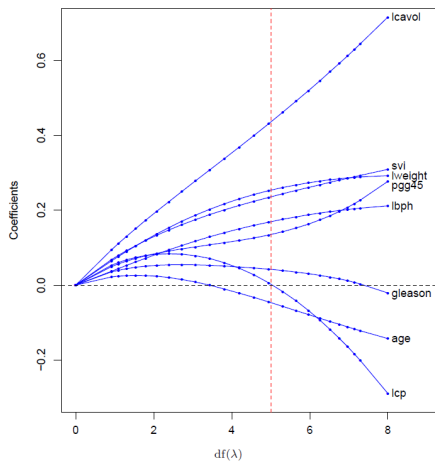
- ▶ For fixed  $\lambda$ , small  $d_j^2 \Rightarrow$  great amount of shrinkage.
- ▶ The meaning of small  $d_j^2$ :
  - ▶ Eigen decomposition of  $\mathbf{X}^\top \mathbf{X}$ :  
 $\mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top$ , where  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_p)$ , where  $\mathbf{v}_j$  is eigenvectors of  $\mathbf{X}^\top \mathbf{X}$ .
  - ▶  $\mathbf{z}_j = \mathbf{X} \mathbf{v}_j$  :  $j^{\text{th}}$  principal component.
  - ▶  $\mathbf{z}_1 = \mathbf{X} \mathbf{v}_1 = \mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{v}_1 = \mathbf{u}_1 d_1$ .
  - ▶ Sample variance of  $\mathbf{z}_1$ :  $\text{Var}(\mathbf{z}_1) = \frac{\mathbf{z}_1^\top \mathbf{z}_1}{N} = \frac{d_1^2 \mathbf{u}_1^\top \mathbf{u}_1}{N}$ .
  - ▶ Since  $d_1 \geq d_2 \geq \dots \geq d_p$ ,  
 $\text{Var}(\mathbf{z}_1) \geq \text{Var}(\mathbf{z}_2) \geq \dots \geq \text{Var}(\mathbf{z}_p)$ .
  - ▶ Small  $d_j^2$  means the directions with small variance in the column space of  $\mathbf{X} \Rightarrow$  ridge regression shrinks these directions.

# Ridge Regression: df

- ▶ The *effective degrees of freedom* of the ridge regression fit is

$$\text{df}(\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}.$$

- ▶ Effective d.f. = the # of free parameters.
- ▶ Effective d.f. of linear regression =  $\text{tr}(\mathbf{H}) = p$ , when the intercept term is removed by standardization.
- ▶ Note that  $\text{df}(0) = p$  and  $\text{df}(\infty) = 0$ .
- ▶ Note that the coefficients tend to get closer to 0 as  $\lambda$  increases.



**FIGURE 3.8.** Profiles of ridge coefficients for the prostate cancer example, as the tuning parameter  $\lambda$  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.

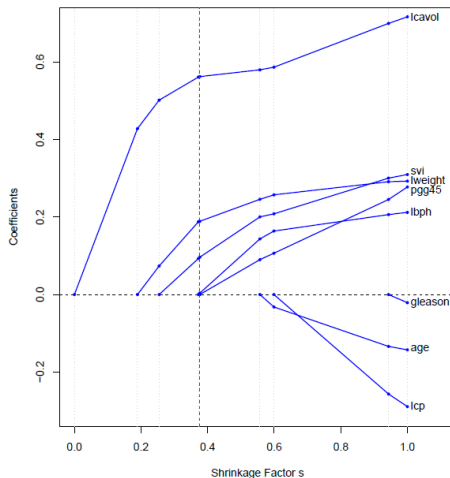
- ▶ Standardization of input variables.
- ▶ Lasso estimator:

$$\hat{\beta}^L = \arg \min_{\beta} \left[ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right].$$

$$\hat{\beta}^L = \arg \min_{\beta} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \text{ subject to } \sum_{j=1}^p |\beta_j| \leq t.$$

- ▶ Unlike ridge regression, as  $\lambda \uparrow$  (or  $t \downarrow$ ), some of  $\hat{\beta}_j^L \rightarrow 0$ .
- ▶ If  $t > \sum_{j=1}^p |\beta_j^{ls}|$ ,  $\hat{\beta}_j^L = \hat{\beta}_j^{ls}$ .

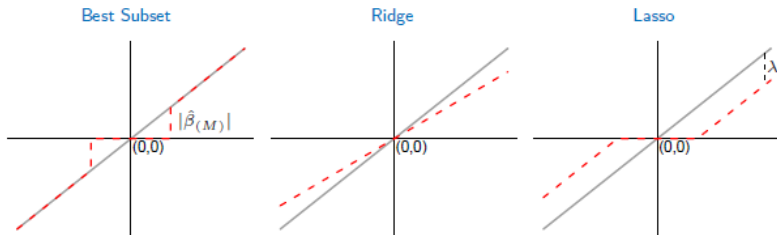
# Lasso fit vs. Restriction



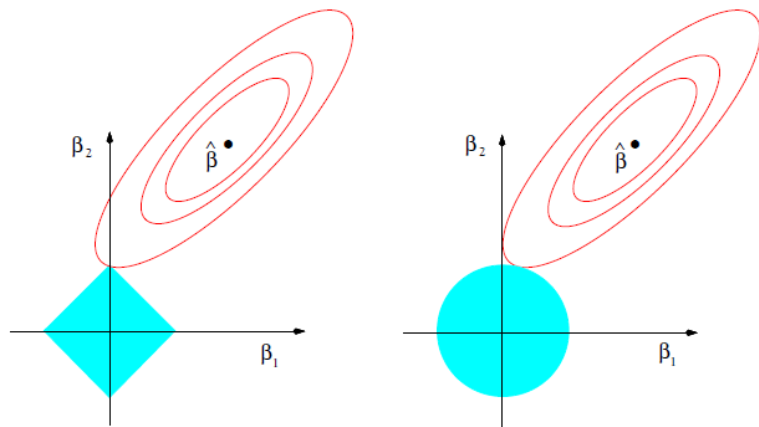
**FIGURE 3.10.** Profiles of lasso coefficients, as the tuning parameter  $t$  is varied. Coefficients are plotted versus  $s = t / \sum_1^p |\hat{\beta}_j|$ . A vertical line is drawn at  $s = 0.36$ , the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso profiles hit zero, while those for ridge do not. The profiles are piece-wise linear, and so are computed only at the points displayed; see Section 3.4.4 for details.

# Lasso vs. Ridge

When the columns of  $\mathbf{X}$  are orthonormal.



# Lasso vs. Ridge



**FIGURE 3.11.** Estimation picture for the lasso (left) and ridge regression (right). 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.



## ► Generalization

$$\hat{\beta} = \arg \min_{\beta} \left[ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right].$$

- $q = 0$ : Subset selection.
- $q = 1$ : Lasso regression.
- $q = 2$ : Ridge regression.

## ► Problems of lasso:

- No analytical solution.
- Not consistent estimator (Non-zero coefficients to be biased towards zero).

# Elastic Net & Grouped Lasso

- ▶ **Elastic net** penalty:

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

- ▶  $0 \leq \alpha \leq 1$ . Compromise between ridge & lasso.
  - ▶ It selects variables like the lasso and shrinks together the coefficients of correlated inputs.
- ▶ **Grouped lasso**: Inputs with pre-defined groups (e.g., genes with the same pathway, a set of dummy variables, etc)

$$\min_{\beta} \left[ \| \mathbf{y} - \beta_0 \mathbf{1} - \sum_{l=1}^L \mathbf{X}_l \beta_l \|_2^2 + \lambda \sum_{l=1}^L \sqrt{p_l} \| \beta_l \|_2 \right].$$

- ▶  $L$  groups of  $p$  predictors.
- ▶  $\mathbf{X}_l$ : Predictors of the  $l$ th group.
- ▶  $\sqrt{p_l}$  accounts for the group size.

- **SCAD** (Smoothly Clipped Absolute Deviation) penalty

$$\begin{cases} \lambda|\beta| & \text{if } |\beta| \leq \lambda \\ -\frac{|\beta|^2 - 2a\lambda|\beta| + \lambda^2}{2(a-1)} & \text{if } \lambda < |\beta| \leq a\lambda \\ \frac{(a+1)\lambda^2}{2} & \text{if } |\beta| > a\lambda \end{cases}$$

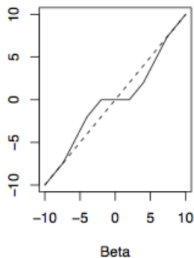
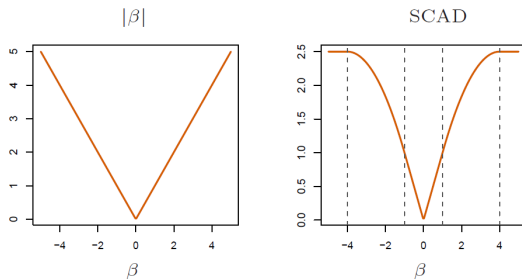
- $a > 2, \lambda > 0$ .

- Solution of the SCAD penalty:

$$\begin{cases} (|\hat{\beta}_j| - \lambda)_+ \text{sign}(\hat{\beta}_j) & \text{if } |\hat{\beta}_j| \leq 2\lambda \\ \{(a-1)\hat{\beta}_j - \text{sign}(\hat{\beta}_j)a\lambda\}/(a-2) & \text{if } 2\lambda < |\hat{\beta}_j| \leq a\lambda \\ \hat{\beta}_j & \text{if } |\hat{\beta}_j| > a\lambda \end{cases}$$

- Small coefficients being set to zero, a few other coefficients being shrunk towards zero while retaining the large coefficients. Thus, SCAD can produce sparse set of solution and approximately unbiased coefficients for large coefficients.

# Lasso vs. SCAD



- ▶ Adaptive lasso penalty:

$$\sum_{j=1}^p w_j |\beta_j|,$$

- ▶  $w_j = 1/|\hat{\beta}_j|^\nu$ , where  $\hat{\beta}_j$  is LSE and  $\nu > 0$ .
- ▶ Consistent estimator.

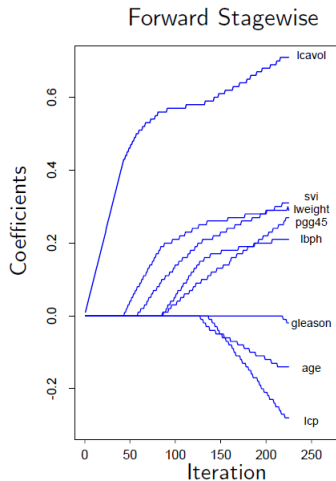
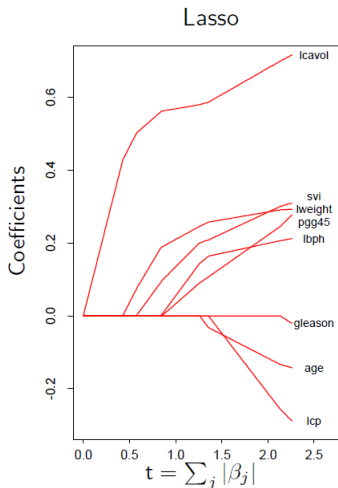
# Regularization Path Algorithm

- ▶ Lasso solution:
  - ▶ If columns of input matrix  $\mathbf{X}$  are not orthogonal, there does not exist analytical lasso solution.
  - ▶ Efficient algorithms for lasso solution.
  - ▶ Algorithms developed from boosting (ensemble learning)  
⇒ Incremental forward stagewise regression.
- ▶ Path algorithm:
  - ▶ Incremental forward stagewise regression.
  - ▶ Least angle regression (LAR).

# Incremental Forward Stagewise Regression

- ▶ Version of least squares boosting for linear regression.
- ▶ Forward selection: Discrete process.
- ▶ Incremental Forward Stagewise Regression: Almost continuous process.
- ▶ Standardized inputs are used.
- ▶ **Algorithm:**
  1. Start with the residual  $\mathbf{r} = \mathbf{y}$  and  $\beta_1, \dots, \beta_p = 0$ . All  $\mathbf{x}_j$ 's are standardized.
  2. Find the predictor  $\mathbf{x}_j$  most correlated with  $\mathbf{r}$ .
  3. Update  $\beta_j \leftarrow \beta_j + \delta_j$ , where  $\delta_j = \epsilon \cdot \text{sign}(\langle \mathbf{x}_j, \mathbf{r} \rangle)$  and  $\epsilon > 0$  is a small step size, and set  $\mathbf{r} \leftarrow \mathbf{r} - \delta_j \mathbf{x}_j$ .
  4. Repeat steps 2 and 3 many times, until the residuals are uncorrelated with all the predictors.

# Incremental Forward Stagewise Regression vs. Lasso





# Least Angle Regression (LAR)

- ▶ LAR: Democratic version of the forward selection method.
  - ▶ LAR coefficients: As much of a predictor as it deserves.
- ▶ Modified version of LAR  $\Rightarrow$  Lasso solution.
- ▶ Algorithm:
  1. Standardize inputs to have mean zero and unit norm.
  2. Start with  $\mathbf{r} = \mathbf{y} - \bar{y}\mathbf{1}$  &  $\beta_1 = \dots = \beta_p = 0$ .
  3. Find the input  $\mathbf{x}_j$  most correlated with  $\mathbf{r}$ .
  4. Move  $\beta_j$  from 0 towards its LSE  $\langle \mathbf{x}_j, \mathbf{r} \rangle$ , until some other input  $\mathbf{x}_k$  has as much correlation with the current residual as does  $\mathbf{x}_j$ .
  5. Move  $\beta_j$  &  $\beta_k$  in the direction defined by LSE of  $\mathbf{r} = \beta_j \mathbf{x}_j + \beta_k \mathbf{x}_k$ , until some other input  $\mathbf{x}_l$  has as much correlation with the current residual.
    - 5-1. If a non-zero coefficient hits zero, drop its variable from the active set of variables & recompute the current joint least square direction. (optional for lasso solution)
  6. Continue in this way until all  $p$  inputs have been entered.

- ▶ Derived input direction methods:
  - ▶ Principal component regression (PCR): Directions with high variance in  $\mathbf{X}$ .
  - ▶ Partial Least Squares (PLS): Directions with high variance in  $\mathbf{X}$  and high correlation with  $\mathbf{y}$ .

# Principal Components

- ▶ Applying the SVD of  $\mathbf{X}$  to obtain an expression for  $\mathbf{X}^\top \mathbf{X}$ , we obtain

$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{U} \mathbf{D} \mathbf{V}^\top = \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top.$$

- ▶ The columns of  $\mathbf{V}$  are the eigenvectors of  $\mathbf{X}^\top \mathbf{X}$ .
- ▶ The  $j$ th principal component of  $\mathbf{X}$  is defined as

$$\mathbf{z}_j = \mathbf{X} \mathbf{v}_j = \mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{v}_j = d_j \mathbf{u}_j.$$

- ▶ The first principal component has the largest sample variance among all normalized linear combinations of the columns of  $\mathbf{X}$ . The last principal component has the minimum variance.

# Principal Components Regression

- ▶ Regress the output  $\mathbf{y}$  on the first  $M(\leq p)$  principal components.
- ▶  $\hat{\mathbf{y}}^{\text{pcr}} = \bar{y}\mathbf{1} + \sum_{m=1}^M \hat{\theta}_m \mathbf{z}_m$ ,  $\hat{\theta}_m = \frac{\langle \mathbf{z}_m, \mathbf{y} \rangle}{\langle \mathbf{z}_m, \mathbf{z}_m \rangle}$
- ▶ Since each  $\mathbf{z}_m$  is a linear combination of the original  $\mathbf{x}$ 's, we can give an estimate of the original slope coefficients:

$$\hat{\beta}^{\text{pcr}}(M) = \sum_{m=1}^M \hat{\theta}_m \mathbf{v}_m.$$

- ▶ If  $M = p$ , then  $\hat{\beta}^{\text{pcr}}(p) = \hat{\beta}$ .
- ▶ The proportion variation in the inputs explained by first  $M$  principal components is  $\frac{\sum_{j=1}^M d_j^2}{\sum_{j=1}^p d_j^2}$ .

# Partial Least Squares (PLS)

- ▶ Unlike PCR, PLS uses  $\mathbf{y}$  to derive input directions.
- ▶ Algorithm:
  1. Standardize inputs to have mean zero and unit variance.
  2. Set  $\hat{\mathbf{y}}^{(0)} = \bar{y}\mathbf{1}$  &  $\mathbf{x}_j^{(0)} = \mathbf{x}_j$ ,  $j = 1, \dots, p$ .
  3. For  $m = 1, \dots, p$ ,
    - 3-1.  $\mathbf{z}_m = \sum_{j=1}^p \hat{\phi}_{mj} \mathbf{x}_j^{(m-1)}$ , where  $\hat{\phi}_{mj} = \langle \mathbf{x}_j^{(m-1)}, \mathbf{y} \rangle$ .
    - 3-2.  $\hat{\theta}_m = \frac{\langle \mathbf{z}_m, \mathbf{y} \rangle}{\langle \mathbf{z}_m, \mathbf{z}_m \rangle}$ .
    - 3-3.  $\hat{\mathbf{y}}^{(m)} = \hat{\mathbf{y}}^{(m-1)} + \hat{\theta}_m \mathbf{z}_m$ .
    - 3-4. Orthogonalize each  $\mathbf{x}_j^{(m-1)}$  w.r.t.  $\mathbf{z}_m$  (i.e.,
$$\mathbf{x}_j^{(m)} = \mathbf{x}_j^{(m-1)} - \frac{\langle \mathbf{z}_m, \mathbf{x}_j^{(m-1)} \rangle}{\langle \mathbf{z}_m, \mathbf{z}_m \rangle} \mathbf{z}_m, \quad j = 1, \dots, p).$$
  4. Output: The fitted vector  $\hat{\mathbf{y}}^{(m)}$ .