

Linear Methods for Regression

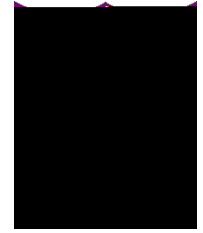
Lijun Zhang

zlj@nju.edu.cn

<http://cs.nju.edu.cn/zlj>



Outline



Introduction

Linear Regression Models and Least Squares

Subset Selection

Shrinkage Methods

Methods Using Derived Input Directions

Discussions

Summary

Introduction

Let $X = [X_1, \dots, X_n]$ be a data point, a linear regression model assumes

$$E(Y | X)$$

is a linear function of X_1, \dots, X_n

Advantages

- They are simple and often provide an adequate and interpretable description

- They can sometimes outperform nonlinear models

 - Small numbers of training cases, low signal-to-noise ratio or sparse data

- Linear methods can be applied to transformations of the inputs

Outline

Introduction

**Linear Regression Models and
Least Squares**

Subset Selection

Shrinkage Methods

Methods Using Derived Input
Directions

Discussions

Summary

Linear Regression Models



The Linear Regression Model

()

β 's are unknown coefficients

The variable X could be

Quantitative inputs

Transformations of quantitative inputs

Log, square-root or square

Basis expansions ($X = X$, $X = X$)

Numeric coding of qualitative inputs

Least Squares

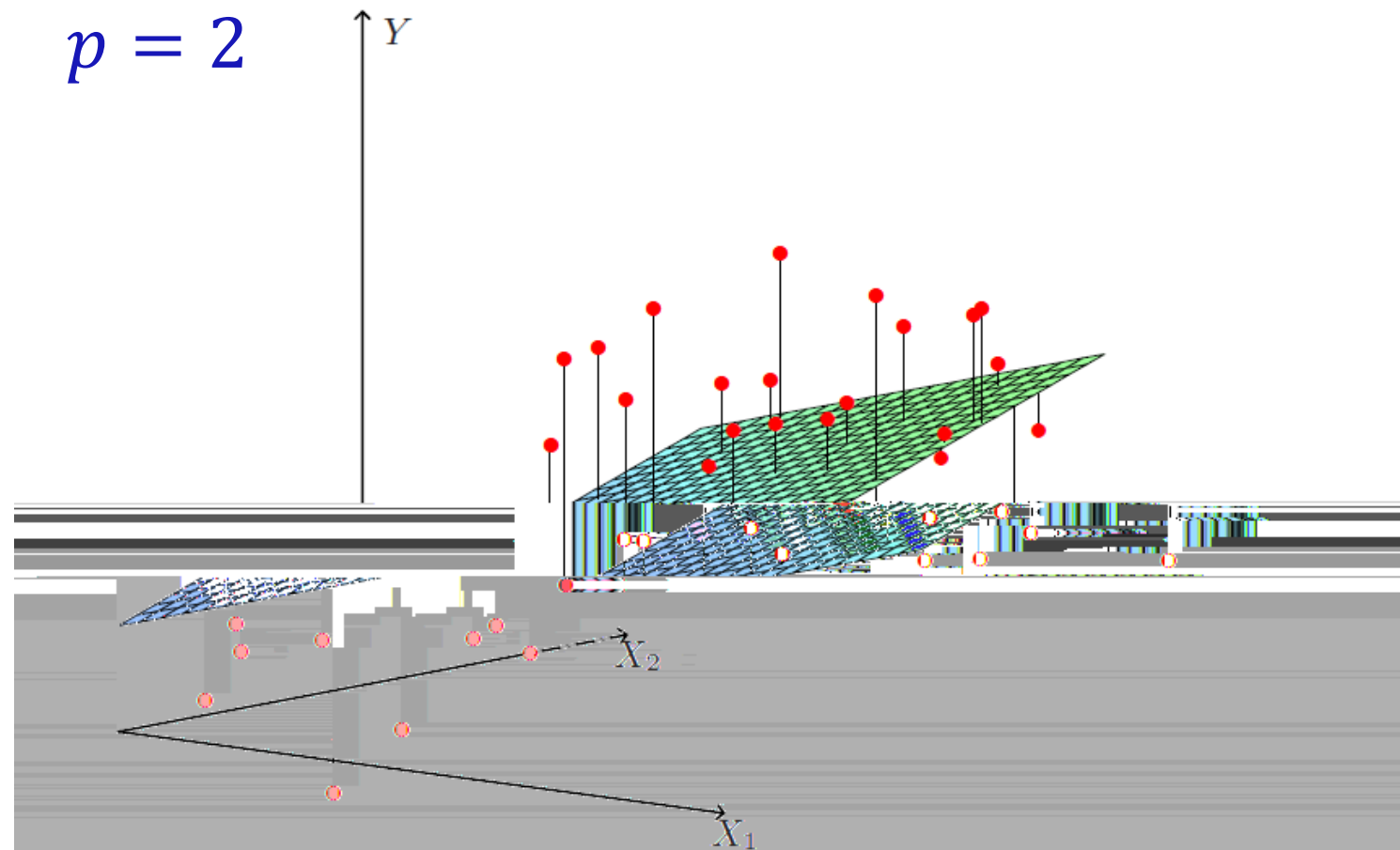
Given a set of training data $(x_1, y_1) \dots (x_N, y_N)$ where $x_i = [x_{i1}, x_{i2}, \dots, x_{ip}]$

Minimize the Residual Sum of Squares

$$RSS(\beta) = \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2$$

Valid if the y 's are conditionally independent given the inputs x

A Geometric Interpretation



3.1. Linear least squares fitting with $X \in \mathbb{R}^2$. We seek the linear function of X that minimizes the sum of squared residuals from Y .

FIGURE
function of

Optimization (2)

Differentiate with respect to β

$$\frac{\partial \text{RSS}}{\partial \beta} = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta)$$

Set the derivative to zero

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) = 0$$

Assume \mathbf{X} is invertible

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

Predictions

The Prediction of x

$$\hat{f}(x_0) = (1 : x_0)^T \hat{\beta}$$

The Predictions of Training Data

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

Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$

$$\hat{\beta} = \operatorname{argmin} \|\mathbf{y} - \mathbf{X}\beta\|$$

$\hat{\mathbf{y}}$ is the orthogonal projection of \mathbf{y} onto the subspace spanned by $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$

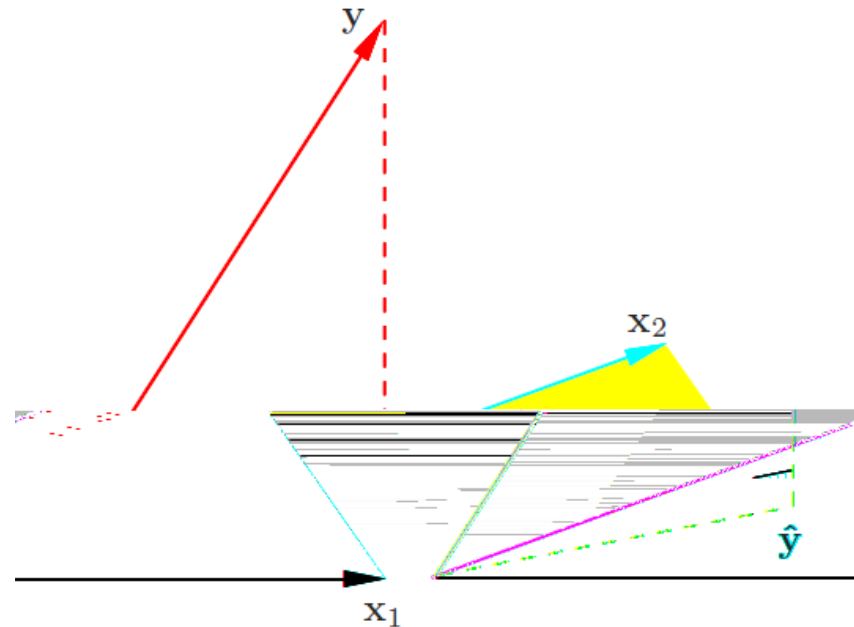
Predictions

The Prediction of x

$$\hat{f}(x_0) = (1 : x_0)^T \hat{\beta}$$

The Predictions of Training Data

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



Understanding (1)



Assume the linear model is right, but the observation contains noise

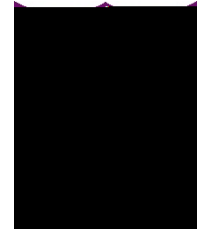
$$\begin{aligned} Y &= E(Y|X_1, \dots, X_p) + \varepsilon \\ &= \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon, \end{aligned}$$

Where $\varepsilon \sim N(0, \sigma^2)$

Then

$$\begin{aligned} & \left(\begin{array}{c} Y \\ \vdots \\ Y \end{array} \right) = \left(\begin{array}{c} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{array} \right) \left[\begin{array}{ccc} 1 & X_{11} & \dots \\ \vdots & \vdots & \vdots \\ 1 & X_{1n} & \dots \end{array} \right] \\ & \quad + \left(\begin{array}{c} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{array} \right) \end{aligned}$$

Understanding (2)



Since $\epsilon = [\epsilon_1, \dots, \epsilon_n]$ is a Gaussian random vector, thus

$$+ \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

is also a Gaussian random vector

$$\mathbb{E} \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_n \end{pmatrix} = \mathbb{E} \left(\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \right) = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} + \mathbb{E} \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

$$\text{Cov} \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_n \end{pmatrix} = \text{Cov} \left(\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \right) = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} \text{Cov} \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \text{Cov} \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Thus $\hat{\beta} \sim N(\beta, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2)$

Expected Prediction Error (EPE)

Given a test point x , assume

$$\tilde{\beta} \sim (0, \Sigma)$$

The EPE of $\tilde{f}(x_0) = x_0^T \tilde{\beta}$ is

$$\frac{\mathbb{E}(\|y - \tilde{f}(x_0)\|^2 | x_0) - \mathbb{E}(\|y - f(x_0)\|^2 | x_0)}{\mathbb{E}(\|y - \tilde{f}(x_0)\|^2 | x_0)} = \sigma^2 + \text{MSE}$$

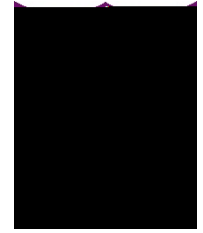
The Mean Squared Error (MSE)

$$\text{MSE}(\tilde{\beta}) = \mathbb{E}(\|y - \tilde{f}(x_0)\|^2 | x_0) - \mathbb{E}(\|y - f(x_0)\|^2 | x_0)$$

$$= \mathbb{E}(\|y - \tilde{f}(x_0)\|^2 | x_0) - \mathbb{E}(\|y - \mathbb{E}(\tilde{f}(x_0))\|^2 | x_0) + \mathbb{E}(\|y - \mathbb{E}(\tilde{f}(x_0))\|^2 | x_0) - \mathbb{E}(\|y - f(x_0)\|^2 | x_0)$$

$$= \text{Variance}(\tilde{f}(x_0)) + \text{Bias}(\tilde{f}(x_0))$$

EPE of Least Squares



Under the assumption that

$$\begin{pmatrix} y \\ \mathbf{X} \end{pmatrix} \sim (0, \Sigma)$$

The EPE of $\hat{f}(x) = x \hat{\beta}$ is

$$\begin{aligned} E\left(\begin{pmatrix} y \\ \mathbf{X} \end{pmatrix} \right) &= E\left(\begin{pmatrix} y \\ \mathbf{X} \end{pmatrix} \right) \\ &= \text{MSE}\left(\begin{pmatrix} y \\ \mathbf{X} \end{pmatrix} \right) \end{aligned}$$

The Mean Squared Error (MSE)

$$\begin{aligned} \text{MSE}\left(\begin{pmatrix} y \\ \mathbf{X} \end{pmatrix} \right) &= E\left(\begin{pmatrix} y \\ \mathbf{X} \end{pmatrix} \right) \\ &= E\left(E\left(\begin{pmatrix} y \\ \mathbf{X} \end{pmatrix} \right) \right) \\ &= \text{Var} \end{aligned}$$

The Gauss–Markov Theorem



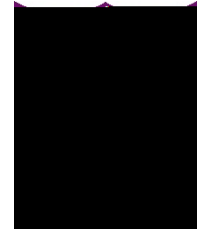
$\hat{\beta}$ has the **smallest variance** among all linear **unbiased** estimates.

We aim to estimate $f(x) = x\beta$, the estimation of $\hat{f}(x) = x\hat{\beta}$ is

From previous discussions, we have

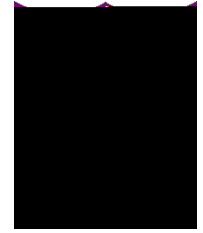
and for all $c \in \mathbb{R}^n$

Multiple Outputs (1)



Suppose we aim to predict K outputs

Multiple Outputs (2)



The Residual Sum of Squares

$$\begin{aligned}\text{RSS}(\mathbf{B}) &= \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - f_k(x_i))^2 \\ &= \text{tr}[(\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB})]\end{aligned}$$

The Solution

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

It is equivalent to performing K independent least squares

Large-scale Setting



The Problem

$$DCC(\beta) = \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2$$

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

Sampling

Faster least squares approximation

Outline

Introduction

Linear Regression Models and Least Squares

Subset Selection

Shrinkage Methods

Methods Using Derived Input Directions

Discussions

Summary

Subset Selection



Limitations of Least Squares

Prediction Accuracy: the least squares estimates often have **low bias** but large variance

Interpretation: We often would like to determine a smaller subset that exhibit the strongest effects

Shrink or **Set** Some Coefficients to Zero

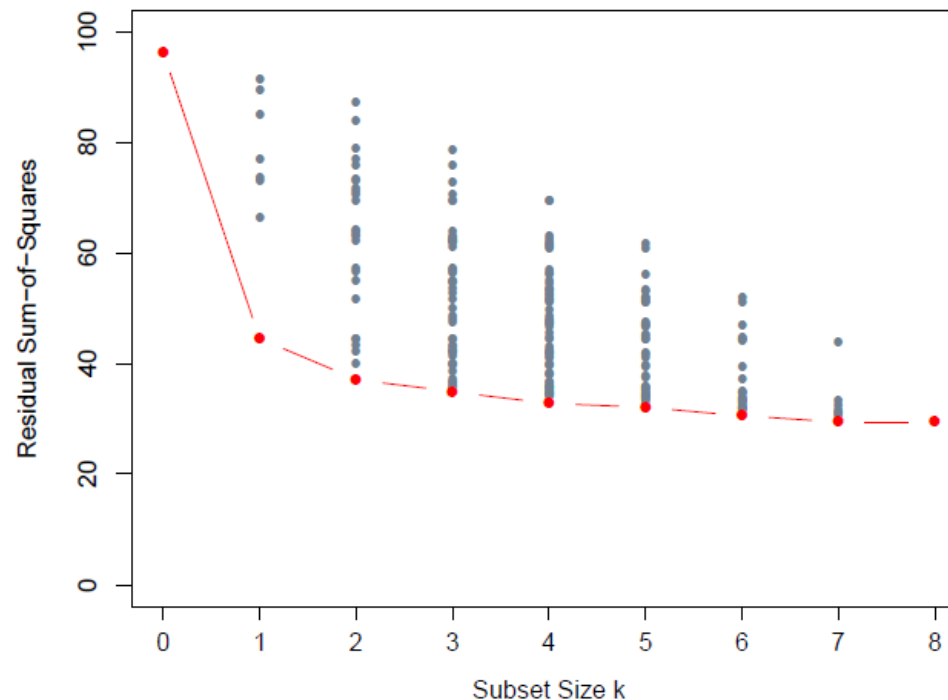
We sacrifice a little bit of bias to reduce the variance of the predicted values

Best-Subset Selection

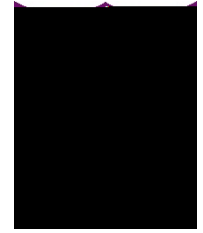
Select the subset of variables (features) such that the RSS is minimized

$$\text{RSS}(\beta) = \sum_{i=1}^N (y_i - f(x_i))^2$$

∞



Forward- and Backward- Stepwise Selection



Forward-stepwise Selection

1. Start with the intercept
2. **Sequentially add** into the model the predictor that most improves the fit

Backward-stepwise Selection

1. Start with the full model
2. **Sequentially delete** the predictor that has the least impact on the fit

Both are **greedy** algorithms

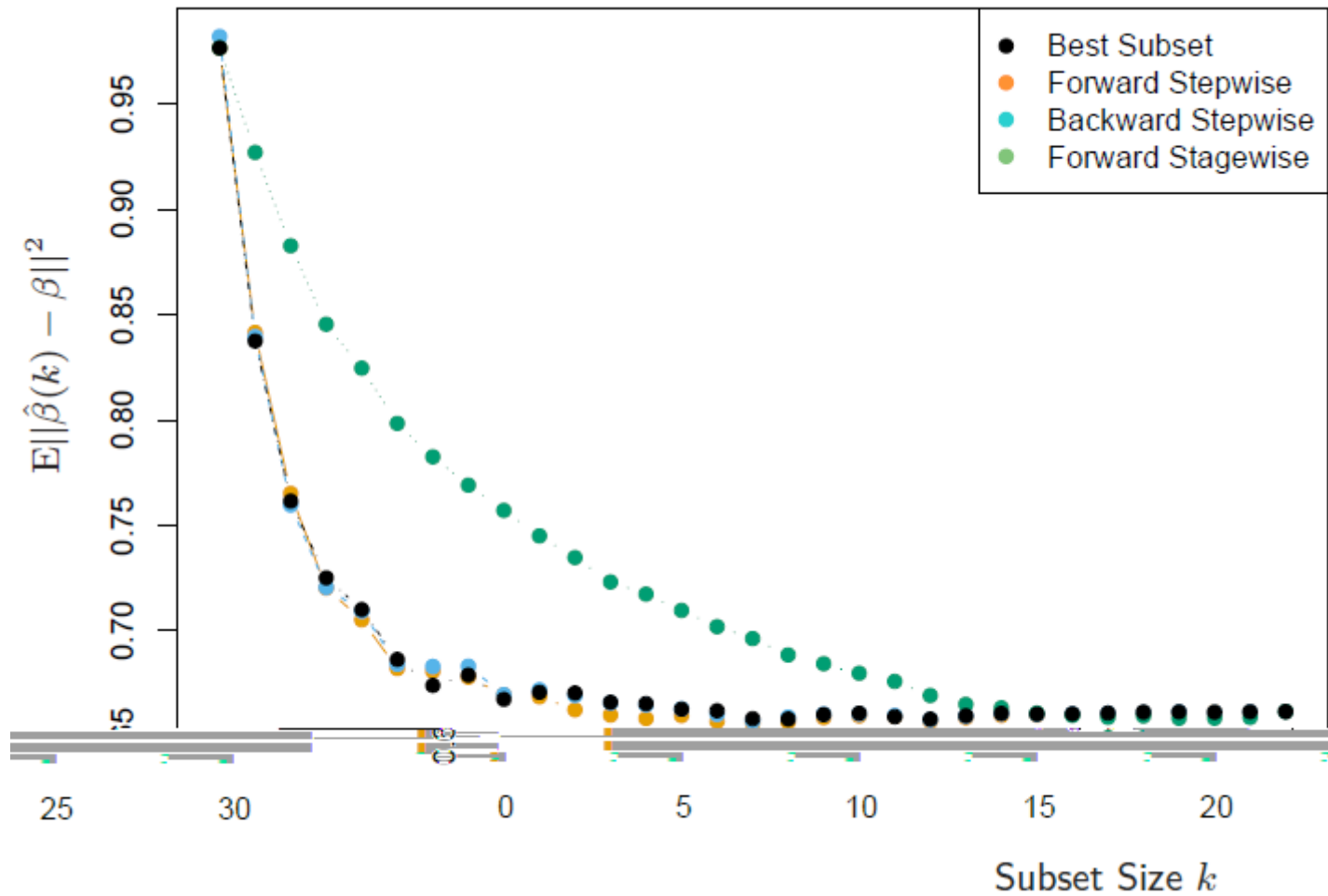
Both can be solved quite **efficiently**

Forward-Stagewise Regression

1. Start with an intercept equal to \bar{y} and centered predictors with coefficients initially all 0
2. Identify the variable **most correlated** with the current residual
3. Compute the **simple** linear regression coefficient of the residual on this chosen variable

None of the other variables are **adjusted** when a term is added to the model

Comparisons



Outline

Introduction

Linear Regression Models and Least Squares

Subset Selection

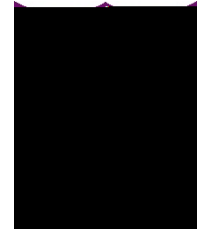
Shrinkage Methods

Methods Using Derived Input Directions

Discussions

Summary

Shrinkage Methods



Limitation of Subset Selection

A **discrete** process—variables are either retained or discarded

It often exhibits **high variance**, and so doesn't reduce the prediction error

Shrinkage Methods

More continuous, low variance

Ridge Regression

The Lasso

Least Angle Regression

Ridge Regression

Shrink the regression coefficients

By imposing a penalty on their size

The Objective

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

$\lambda \geq 0$ is a complexity parameter

An Equivalent Form

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

Coefficients cannot be too large even when variables are correlated

Optimization (1)

Let \mathbf{X} be a matrix with each row an input vector

$$\begin{bmatrix} \\ \\ \end{bmatrix}$$

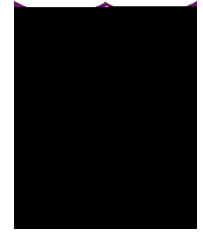
$$\boldsymbol{\beta} = [\beta, \dots, \beta] \quad \text{and} \quad \mathbf{y} = [y, \dots, y]$$

The Objective Becomes

$$\min_{\boldsymbol{\beta}} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|$$

Where $\mathbf{1} = [1, \dots, 1] \in \mathbb{R}$

Optimization (2)



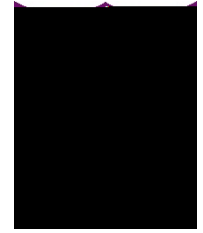
Differentiate with respect to β and set it to zero

$$2 \begin{pmatrix} \beta \\ 1 \end{pmatrix} - \begin{pmatrix} \beta \\ 1 \end{pmatrix} = 0$$

Differentiate with respect to β and set it to zero

$$2 \begin{pmatrix} \beta \\ 1 \end{pmatrix} - 2 \begin{pmatrix} \beta \\ 1 \end{pmatrix} = 0$$
$$\begin{pmatrix} \beta \\ 1 \end{pmatrix} = \begin{pmatrix} \beta \\ 1 \end{pmatrix}$$

Optimization (3)



The Final Solution

Let $H = \mathbf{I} - \dots$

Understanding (1)

Assume \mathbf{X} is centered, then

Let the SVD of \mathbf{X} be

$\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}$, contains the left singular
vectors

$\mathbf{\Sigma}$ is a diagonal matrix with diagonal entries
0

Then, we examine the prediction of
training data

Understanding (2)

Least Squares

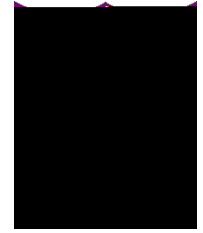
$$\begin{aligned}\mathbf{X}\hat{\beta}^{\text{ls}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \\ &= \mathbf{U}\mathbf{U}^T\mathbf{y},\end{aligned}$$

Ridge Regression

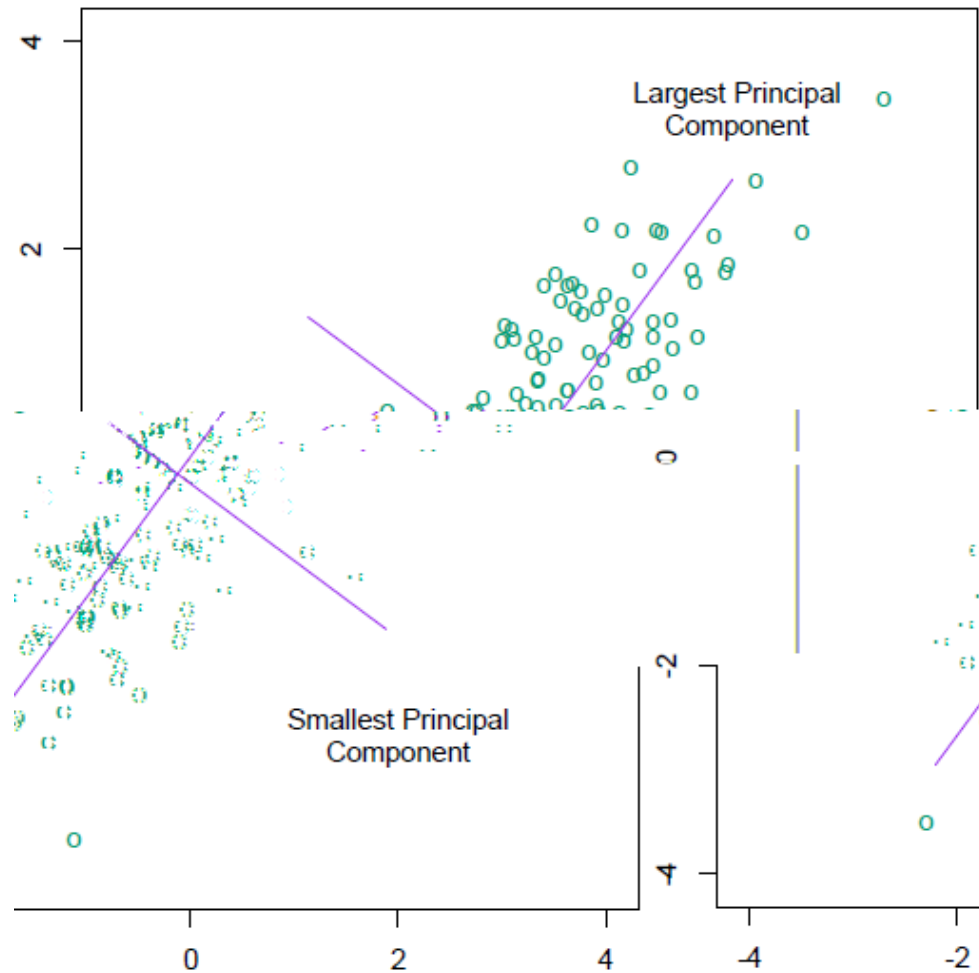
$$\begin{aligned}\mathbf{X}\hat{\beta}^{\text{ridge}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y} \\ &= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y} \\ &= \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T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

Shrink the coordinates by $\frac{d_j^2}{d_j^2 + \lambda} \leq 1$

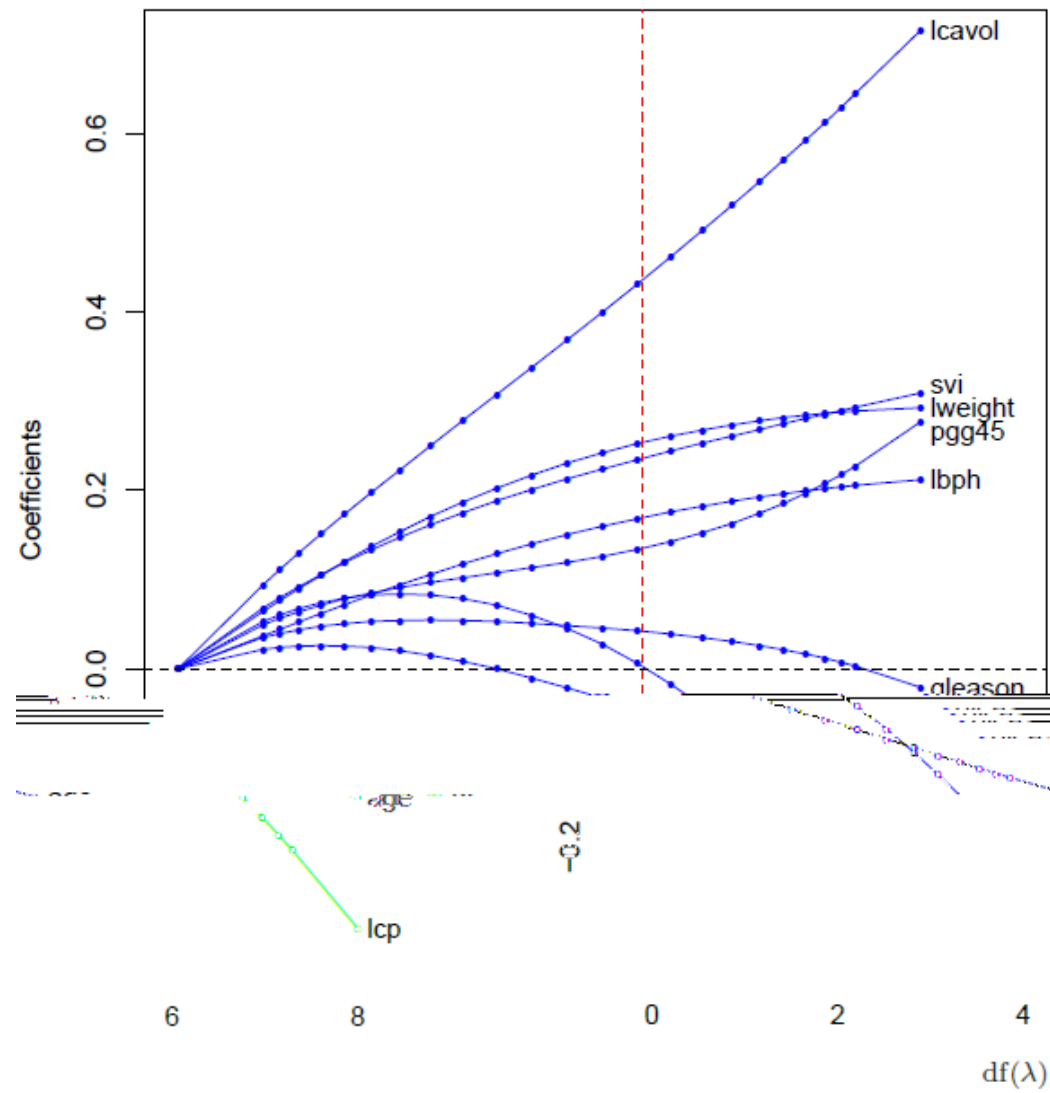
Understanding (3)



Connection with PCA



An Example



Optimization



The First Formulation

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

subject to $\sum_{j=1}^p |\beta_j| \leq t.$

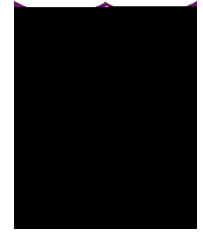
Gradient descent followed by Projection [1]

The Second Formulation

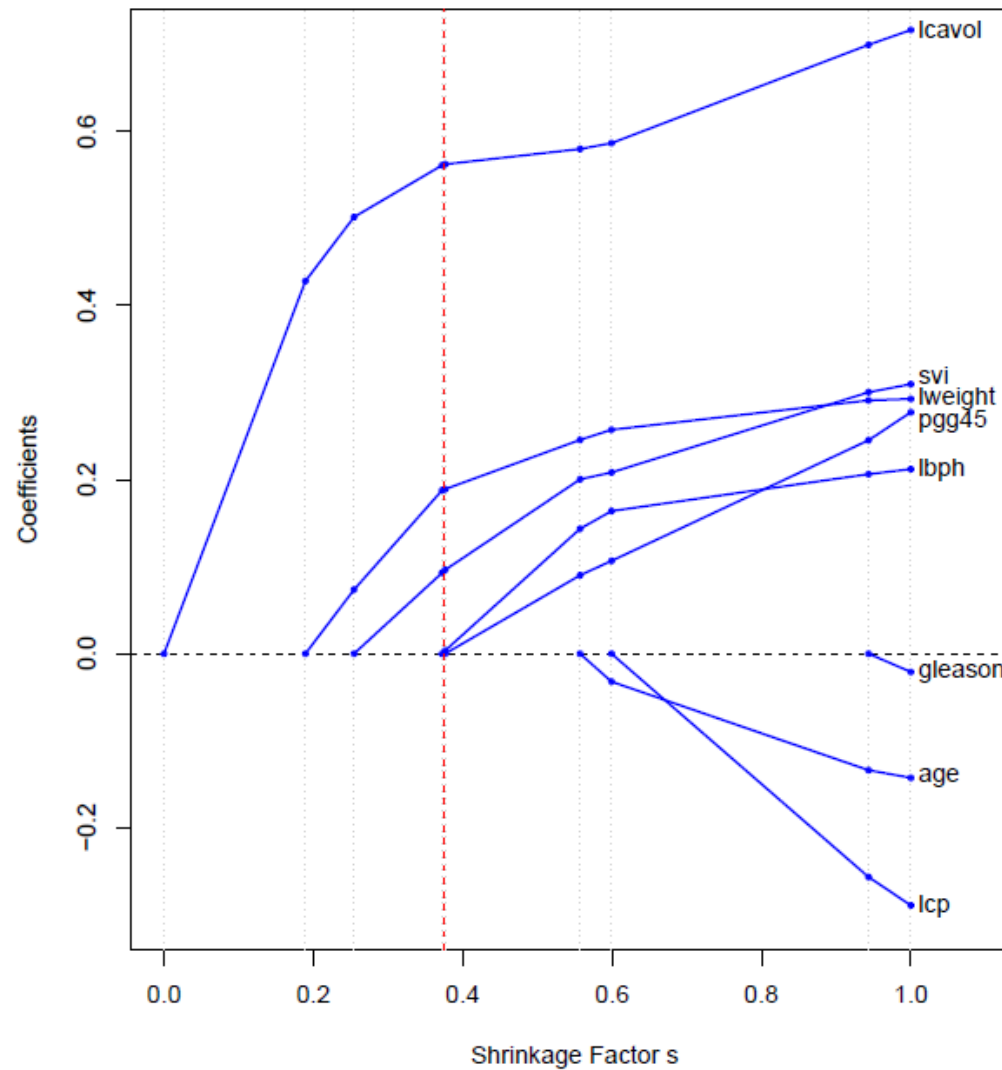
$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

Convex Composite Optimization [2]

An Example



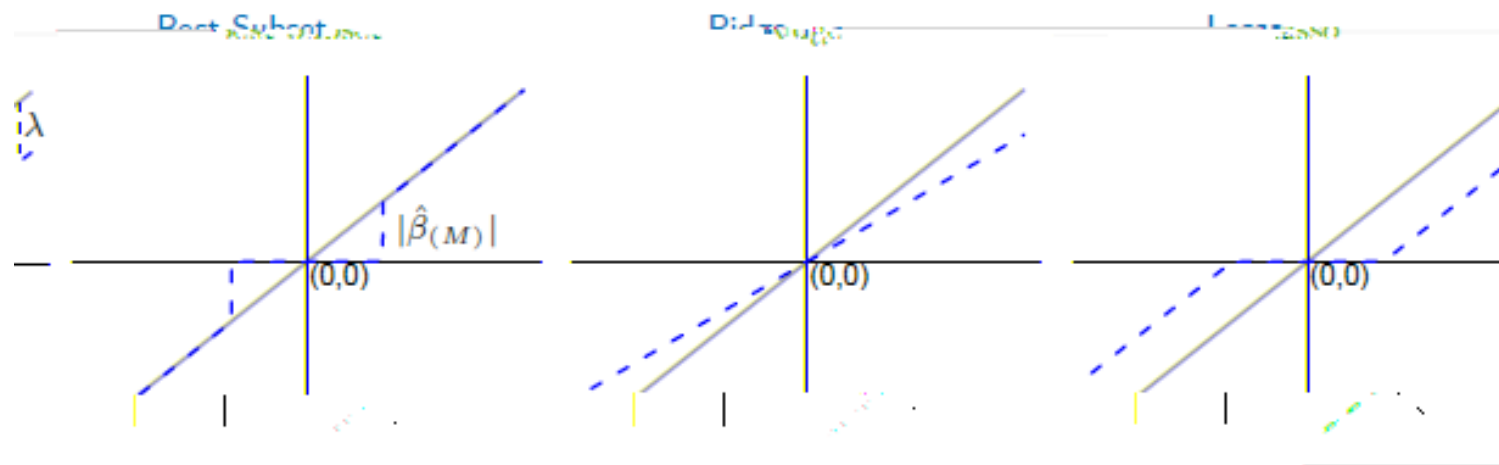
Hit 0
Piece-wise linear



Subset Selection, Ridge, Lasso

Columns of \mathbf{X} are orthonormal

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



Hard-thresholding

Scaling

Soft-thresholding

Ridge v.s. Lasso (1)

Ridge Regression

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

subject to $\sum_{j=1}^p \beta_j^2 \leq t$

ℓ_2 -norm appears in the constraint

Lasso

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

subject to $\sum_{j=1}^p |\beta_j| \leq t$.

ℓ_1 -norm appears in the constraint

Ridge v.s. Lasso (2)

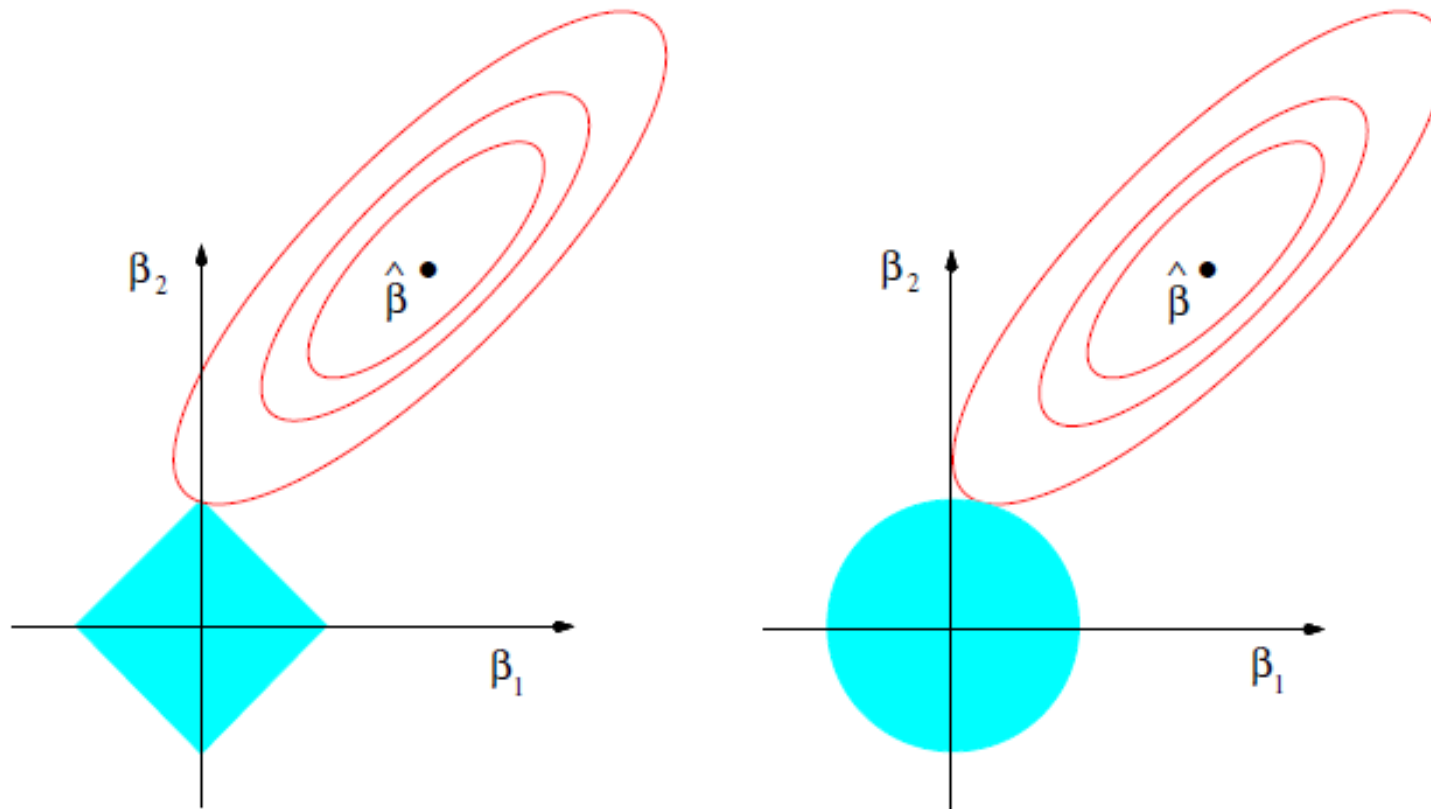


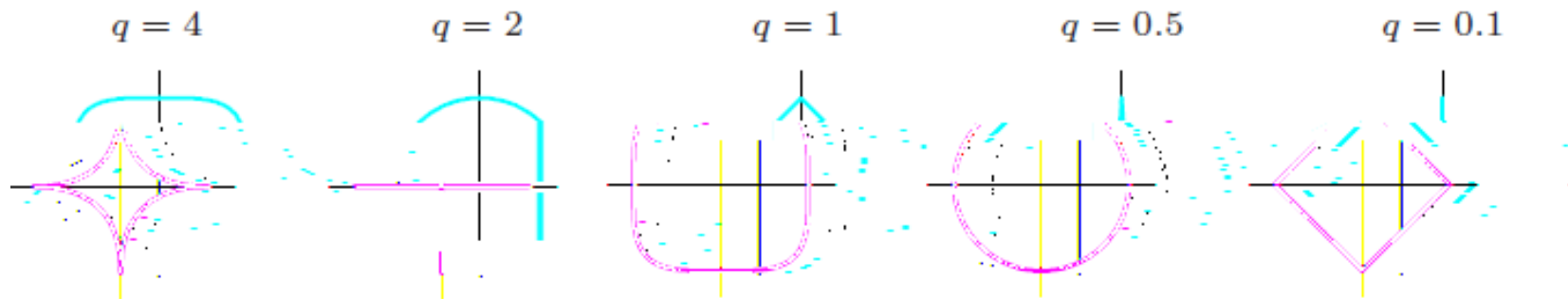
FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \leq t$ and $\beta_1^2 + \beta_2^2 \leq t^2$, respectively, while the red ellipses are the contours of the least squares error function.

Generalization (1)

A General Formulation

$$\tilde{\beta} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right\}$$

Contours of Constant Value of $\sum |\beta|$



Generalization (2)

A Mixed Formulation

The *elastic-net* penalty

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

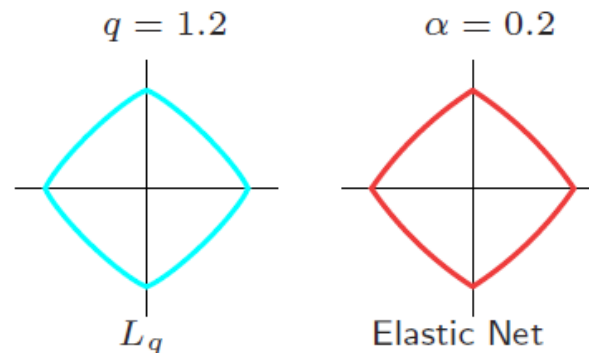


FIGURE 3.13. Contours of constant value of $\sum_j |\beta_j|^q$ for $q = 1.2$ (left plot), and the elastic-net penalty $\sum_j (\alpha \beta_j^2 + (1 - \alpha) |\beta_j|)$ for $\alpha = 0.2$ (right plot). Although very similar, the elastic-net has sharp (non-differentiable) corners, while the $q = 1.2$ penalty does not.

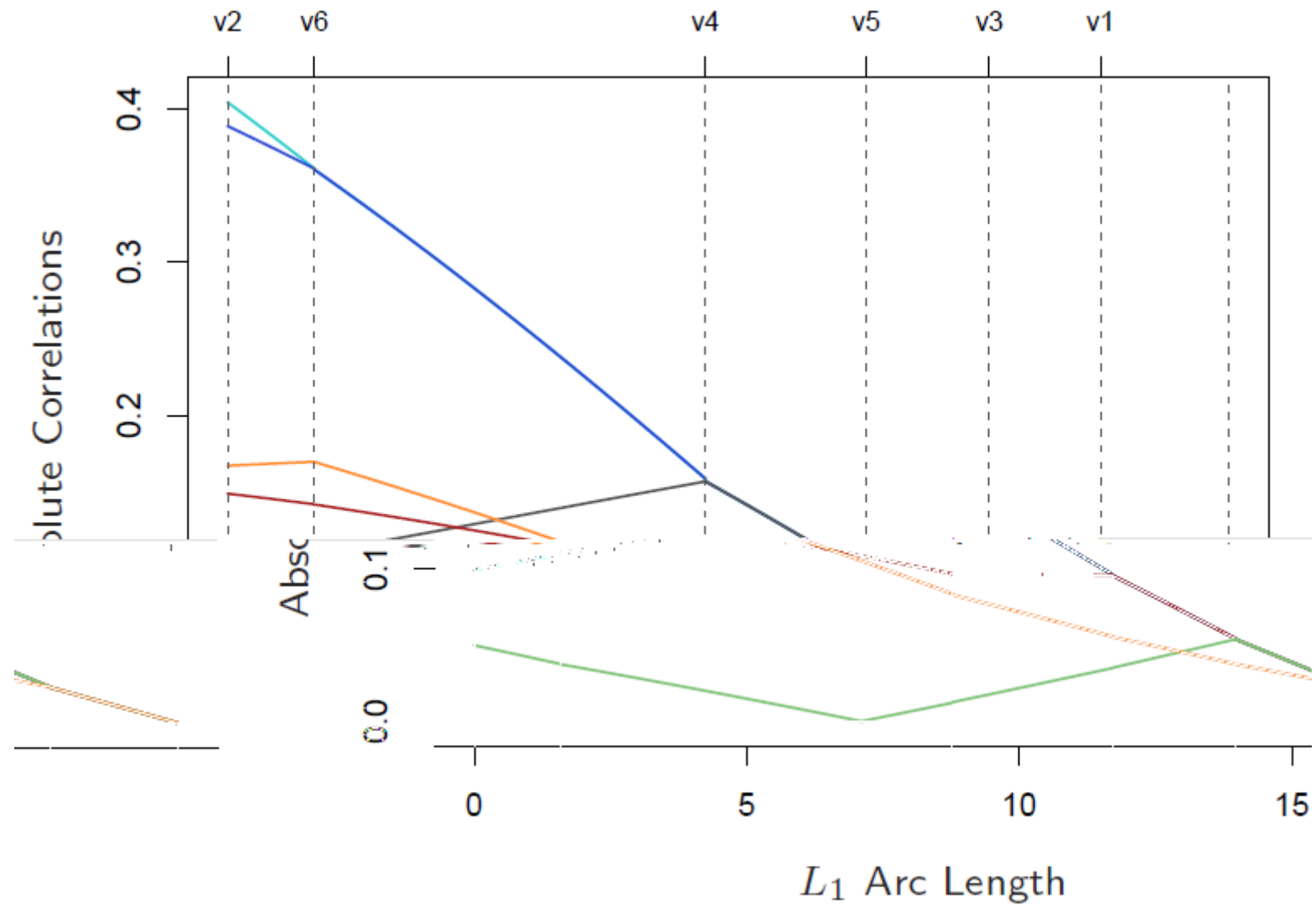
Least Angle Regression (LAR)



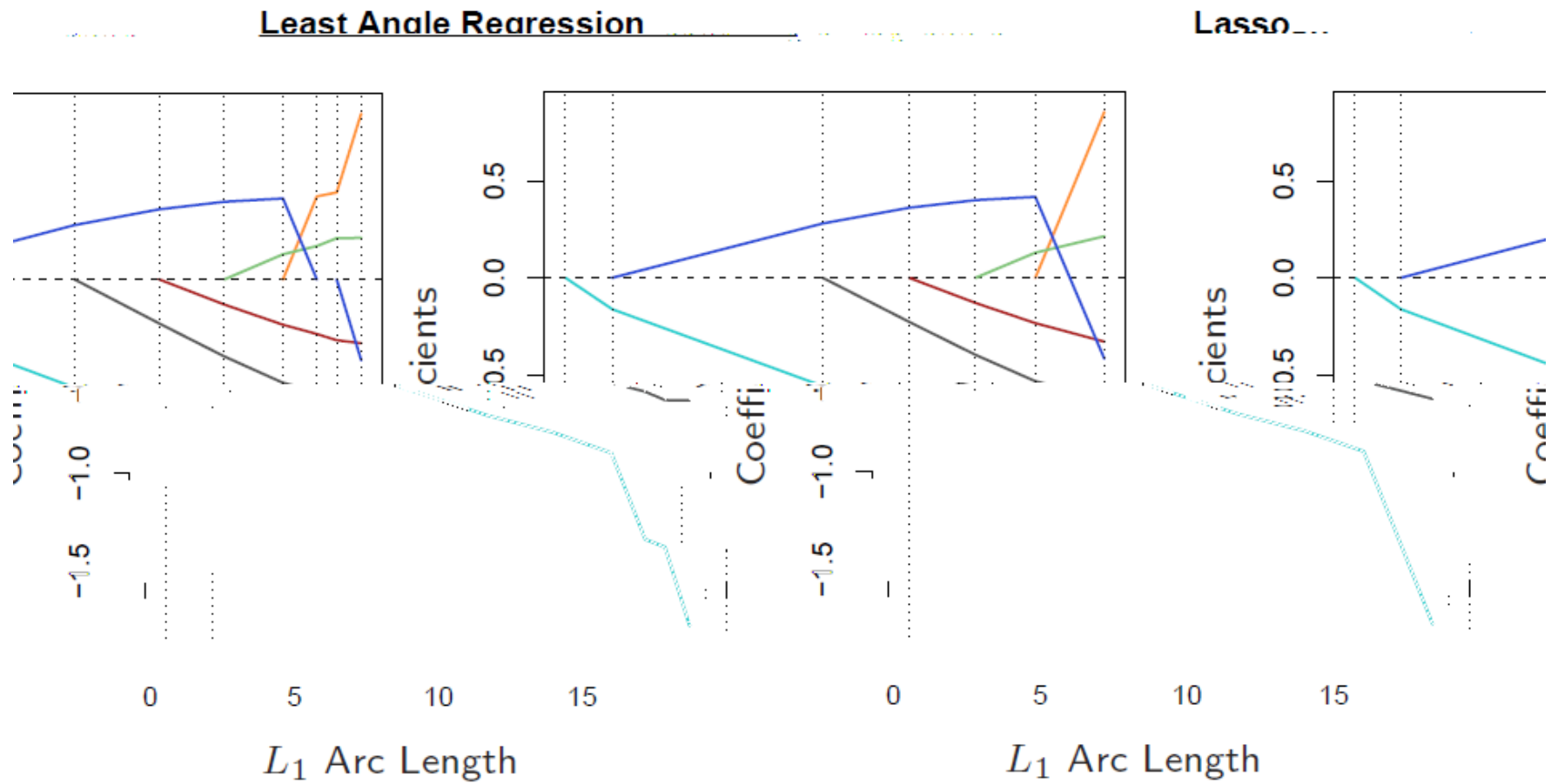
The Procedure

1. Identify the variable most **correlated** with the response
2. Move the coefficient of this variable **continuously** toward its least squares value
3. As soon as another variable "**catches up**" in terms of correlation with the residual, the process is paused
4. The second variable then joins the active set, and **their coefficients are moved together** in a way that keeps their correlations tied and decreasing

An Example



LAS v.s. Lasso





Outline



Introduction

Linear Regression Models and Least Squares

Subset Selection

Shrinkage Methods

Methods Using Derived Input Directions

Discussions

Summary

Methods Using Derived Input Directions

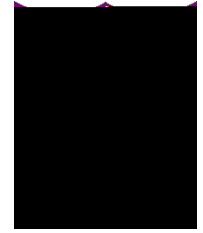


We have a large number of inputs

Often very correlated

1. Generate a small number of linear combinations
 z_1, \dots, z_k
of the original inputs X
 2. Use Z in place of X as inputs in the regression
- Linear Dimensionality Reduction + Regression

Principal Components Regression (PCR)



The linear combinations Z are generated by PCA

\mathbf{X} is centered, and v is the m -th right singular vector

Since \mathbf{z} 's are orthogonal

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

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

PCR v.s. Ridge

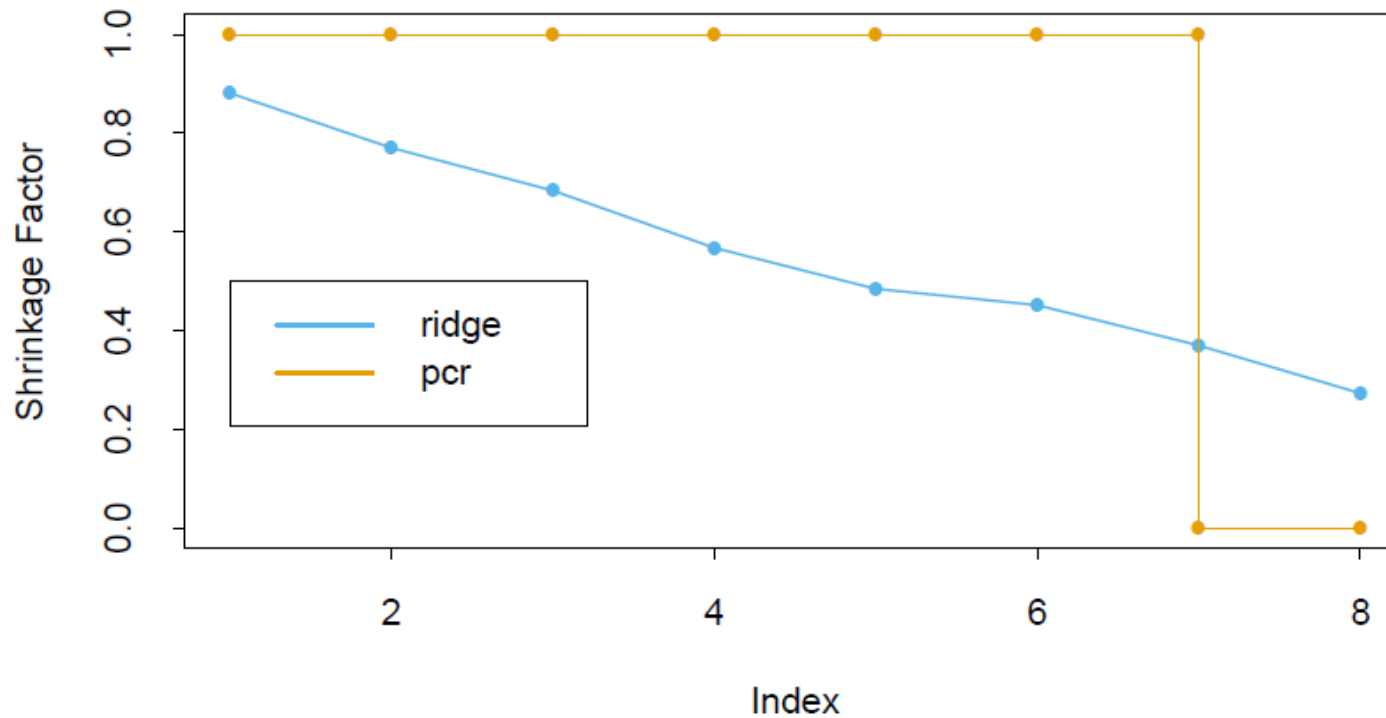


FIGURE 3.17 Ridge regression shrinks the regression coefficients of the principal components, using shrinkage factors $d_j^2 / (d_j^2 + \lambda)$ as in (3.17). Principal component regression truncates them. Shown are the shrinkage and truncation patterns corresponding to Figure 3.7, as a function of the principal component index.

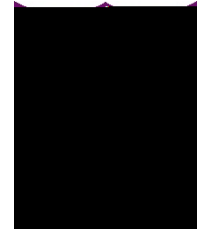
Partial Least Squares (PLS)



The Procedure

1. Compute $\hat{\phi} = \langle \mathbf{x}, \mathbf{y} \rangle$ for each feature \mathbf{x}
2. Construct the 1st derived input $\mathbf{z} = \sum \hat{\phi} \mathbf{x}$
3. \mathbf{y} is regressed on \mathbf{z} giving coefficient $\hat{\theta}$
4. Orthogonalize $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to \mathbf{z}
5. Repeat the above process

Outline



Introduction

Linear Regression Models and Least Squares

Subset Selection

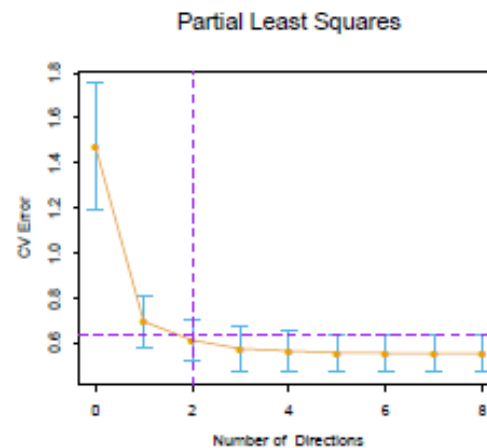
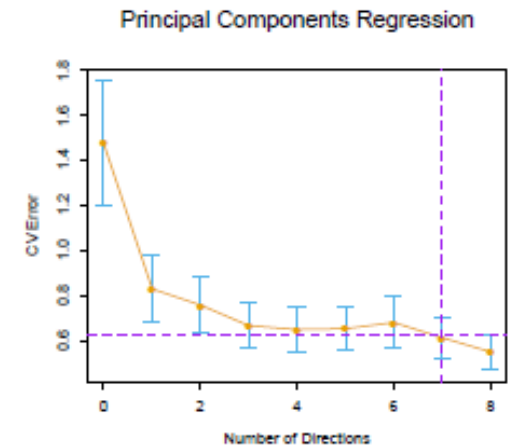
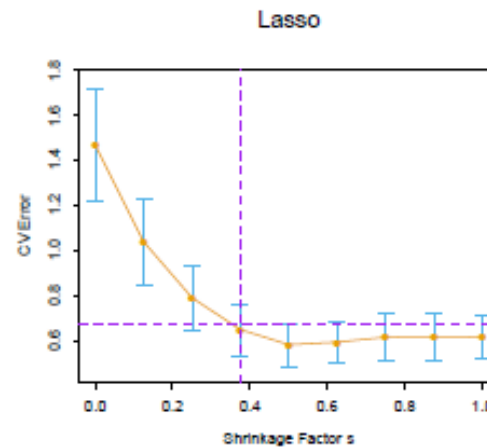
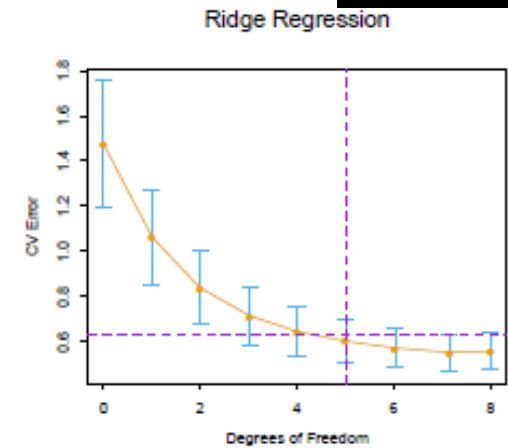
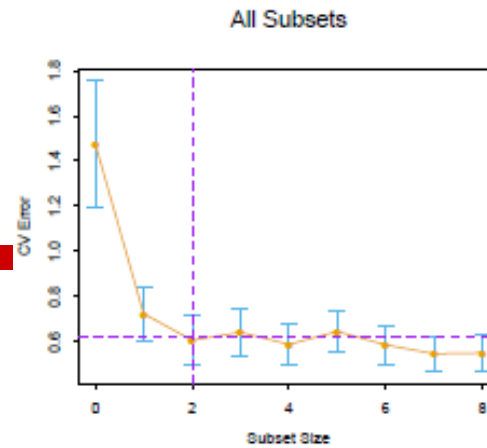
† Square7

†

†

Discussions (1)

Model complexity increases as we move from left to right.

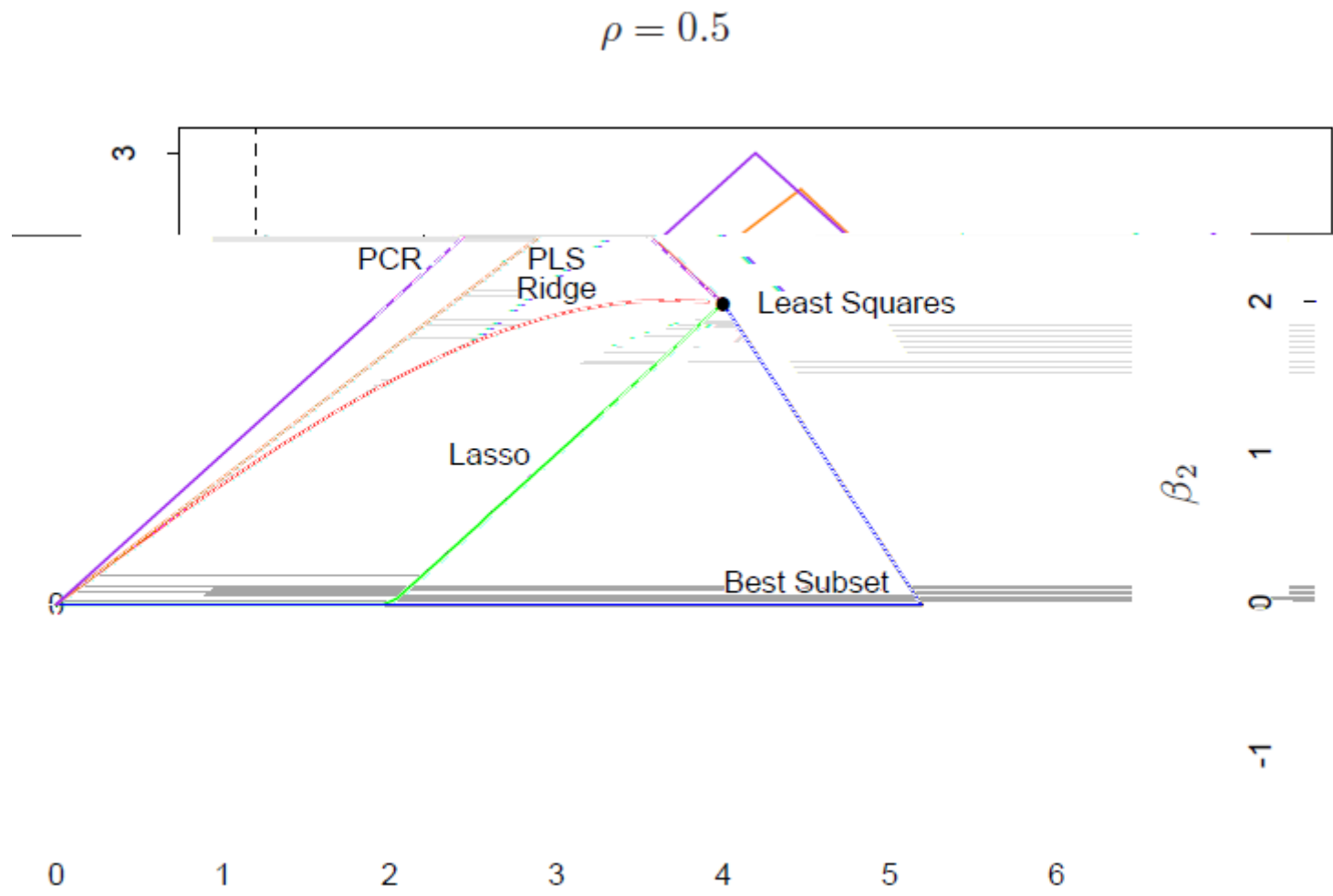


Discussions (2)

TABLE 2.2: Estimated coefficients and test error test statistics, diff. coefficient test, and shrinkage methods applied to the prostate data. The blank entries correspond to variables omitted.

| Term | LS | Best Subset | Ridge | Lasso | PCR | PLS |
|------------|--------|-------------|--------|--------|--------|--------|
| Intercept | 2.465 | 2.477 | 2.452 | 2.468 | 2.497 | 2.452 |
| lcavol | 0.680 | 0.740 | 0.420 | 0.533 | 0.543 | 0.409 |
| lweight | 0.263 | -0.316 | 0.238 | 0.169 | 0.289 | -0.344 |
| age | -0.141 | | -0.046 | | 0.152 | -0.026 |
| lbp1h | 0.210 | | 0.162 | 0.002 | 0.214 | 0.220 |
| svi | 0.305 | | 0.227 | -0.094 | 0.315 | 0.243 |
| lcp | -0.288 | | 0.000 | | -0.051 | 0.079 |
| gleason | -0.021 | | 0.040 | | 0.232 | 0.011 |
| pgg45 | 0.267 | | 0.133 | | -0.056 | 0.068 |
| Test Error | 0.521 | 0.492 | 0.492 | 0.479 | 0.449 | 0.521 |
| Std. Error | 0.152 | 0.179 | 0.143 | 0.165 | 0.164 | 0.152 |

Discussions (3)



Outline



Introduction

Linear Regression Models and Least Squares

Subset Selection

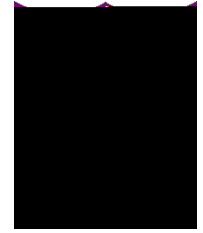
Shrinkage Methods

Methods Using Derived Input Directions

Discussions

Summary

Summary



Linear Regression Models

Least Squares

Shrinkage Methods

- Ridge Regression

- Lasso

- Least Angle Regression (LAR)

Methods Using Derived Input
Directions

- Principal Components Regression (PCR)

- Partial Least Squares (PLS)

Reference

[1] Duchi et al. Efficient projections onto the ℓ -ball for learning in high dimensions. In Proceedings of the 25th international conference on Machine learning, pp. 272-279, 2008.

[2] Nesterov. Gradient methods for minimizing composite functions. Mathematical Programming, 140(1): 125-161, 2013.