

Then select a random permutation of k folds.

HD46

ex) $n=26, k=5$ folds:

4 2 3 5 3 3 (1) 5 2 2 5 (1) 2 (1) 3 4 2 (1) 5 5 (1) 4 (1) 4 4 3
1 2 3 4 5 6 (7) 8 9 10 11 (12) 13 (14) 15 16 17 (18) 19 20 (21) 22 (23) 24 25 26

you add the case numbers

Cases 7, 12, 14, 18, 21 and 23 are in fold 1.

The other folds have 5 cases.

12) variable selection: put standardized predictors into \underline{W} and center the response $\underline{z} = \underline{y} - \bar{y}$.

Do RR on $\underline{z} = \underline{W}\underline{\hat{\beta}} + \underline{e}$

where there is no intercept ($\underline{\hat{\beta}}$ and $(\beta_2, \dots, \beta_p)^T$ are closely related).

Keep the k_i variables with the largest $|\hat{\beta}_i|$, $i=1, \dots, J$, use CV to select k_i^* .

PKD topic

13) Another ridge regression estimator $\tilde{\beta}_{RR}$ minimizes the

criterion $Q_{RR}(\beta) = \text{RSS}(\beta) + \lambda \sum_{j=2}^p \beta_j^2$.

λ is not 1

46.5

$\beta \rightarrow 0$
then
 $Q_{RR} \rightarrow \bar{y}$

§ 3.6 and 3.7 lasso

14) The lasso MLR estimator $\tilde{\beta}_L$

minimizes $Q_L(\beta) = \frac{\text{RSS}(\beta)}{a} + \frac{\lambda n}{a} \sum_{j=2}^p |\beta_j|$.

There is not a simple formula for $\tilde{\beta}_L$, but the criterion is convex, so fast algorithms exist.

15) Convex criterion often have an equivalent dual problem.

Suppose $\underline{z} = \underline{W} \underline{\beta} + \underline{e}$ where $\underline{z} = \underline{y} - \bar{y}$ and the model has no constant.

Then $Q_L(\beta) = \text{RSS}(\beta) + \lambda n \sum_{j=1}^p |\beta_j|$,

and the dual problem is

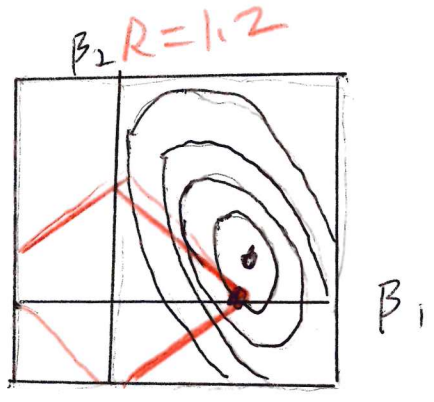
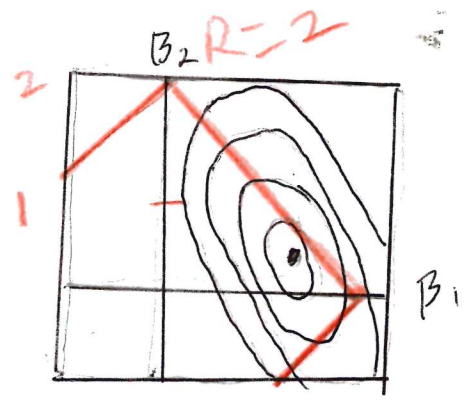
$\hat{\beta}_L$ minimizes

$Q_{LR}(\beta) = \|z - W\beta\|^2$ subject to

$\left\{ \beta \in \mathbb{R}^p \text{ with } \sum_{j=1}^p |\beta_j| \leq R \right\}$ where

each value of $R > 0$ corresponds to a $\lambda = \lambda_m > 0$.

Hyperellipsoids correspond to $RSS(\beta) = \text{constant}$. The \bullet to the OLS estimator.



Girard P 93

$\hat{\beta}_L = \beta_{OLS}$ minimizes

$Q_{LR}(\beta)$
 $L, R = 2$

$\hat{\beta}_L = \begin{pmatrix} 1.2 \\ 0 \end{pmatrix}$ minimizes

$Q_{LR}(\beta)$
 $L, R = 1.2$

16) The lasso estimator does variable selection: as $\lambda \uparrow$ from 0, some $\hat{\beta}_{iL}$ values = 0.

} see

17] LassoCLT: Assume p is fixed

47.5

and γ holds, $\underline{y} = \underline{X}\underline{\beta} + \underline{e}$.

a) If $\frac{\hat{\lambda}_{\min}}{\sqrt{n}} \xrightarrow{P} 0$, then

$$\sqrt{n}(\hat{\underline{\beta}}_L - \underline{\beta}) \xrightarrow{D} N_p(\underline{0}, \sigma^2 \underline{V}).$$

b) If $\frac{\hat{\lambda}_{\min}}{\sqrt{n}} \xrightarrow{P} \tau \geq 0$ and $S_n \xrightarrow{P} \underline{s}$

← see 18]

$$\text{then } \sqrt{n}(\hat{\underline{\beta}}_L - \underline{\beta}) \xrightarrow{D} N_p(-\frac{\tau}{2} \underline{V} \underline{s}, \sigma^2 \underline{V}).$$

c) If $\frac{\hat{\lambda}_{\min}}{n} \xrightarrow{P} 0$, $\hat{\underline{\beta}}_L \xrightarrow{P} \underline{\beta} = \underline{\beta}_{OLS} = \underline{\beta}_F$.

18] Here $S_{in} \in [-1, 1]$, $S_{in} = \text{sign}(\hat{\beta}_{iL})$ if $\hat{\beta}_{iL} \neq 0$.

$$\text{Sign}(\beta_i) = 1 \text{ if } \beta_i > 0 \quad = \beta_{iL} \\ \quad \quad \quad -1 \text{ if } \beta_i < 0.$$

19] Since $\hat{\underline{\beta}}_L \xrightarrow{P} \underline{\beta} = \underline{\beta}_{OLS, F}$, $P(S \subseteq I_{\min}) \rightarrow 1$

corresponds to

where I_{\min} is $\hat{\underline{\beta}}_L$ chosen by k -fold CV
on grid $0 < \lambda_1 < \lambda_2 < \dots < \lambda_J$.

Thus the lasso variable selection estimator is \sqrt{n} consistent.

20) Lasso variable selection is poor if $\frac{\lambda}{\sqrt{n}} \rightarrow 0$. Typically (glmnet) $\lambda_1 \propto n^{3/4}$. Hence $\hat{\beta}_L$ from K-fold CV is at best $n^{1/4}$ consistent; far worse than \sqrt{n} consistent.

21) If $p > n$, "lasso variable selection" uses p predictors for λ very close to 0, but at most $n+1$ predictors, including a constant, once some $\beta_i = 0$. The value λ_J is the smallest value of λ such that $\hat{\beta}_2 = \dots = \hat{\beta}_p = 0$.

(If $y = \alpha + x^T \beta + e$, MLR2, the smallest value of $\lambda \ni \beta = 0$. Hence $\hat{\beta}_{\lambda_i} \neq 0$ for $i=1, 2, \dots, J-1$.)

22) In low dimensions, $\lambda_1 \ll \dots \ll \lambda_J$, behavior near λ_1 is important, and $J=100$ is fine. In high dimensions, using larger J gives more models, making it more likely to get a model good for prediction (multitude of models).

23) simulations! To compare variable selection estimators and MLR estimators, we will count the number of times $S \subseteq I_{\min}$ and make prediction intervals for Y_t .

24) § 2.3, 3.12 know for E2. Let $\underbrace{z_1, \dots, z_n}_{\text{training } z_i}, z_t$ be iid. \uparrow test

The $\text{Shorth}(c)$ interval is the shortest closed interval containing at least c of the z_i . For small ordered data sets, compute the lengths of the intervals containing c cases!

$$[z_{(1)}, z_{(c)}], [z_{(2)}, z_{(c+1)}], \dots, [z_{(n-c+1)}, z_{(n)}].$$

$$z_{(c)} - z_{(1)} \quad z_{(c+1)} - z_{(2)} \quad z_{(n)} - z_{(n-c+1)}.$$

Then $\text{Shorth}(c) = [z_{(s)}, z_{(s+c-1)}]$ is the interval with shortest length.

ex] Find $\text{Shorth}(4)$.

$$0, 1, 3, 6, 9, 10, 11$$

$$6 - 0 = 6$$

$$9 - 1 = 8$$

$$10 - 3 = 7$$

$$11 - 6 = 5$$

$$\boxed{\text{Shorth}(4) = [6, 11]}$$

see HW6

25) Let $\lceil x \rceil =$ smallest integer $\geq x = \text{ceiling}(x)$.

So $\lceil 7.7 \rceil = \lceil 8 \rceil = 8$. If $\frac{c_n}{n} \rightarrow 1-\delta$,

the $\text{shorth}(c_n)$ interval is a large sample

100(1- δ)% PI for Z_F . $c_n = k_n = \lceil n(1-\delta) \rceil$

contains \approx 100(1- δ)% of the training data.

The maximum undercoverage $\lambda \approx 1.12 \sqrt{\frac{\delta}{n}}$.

at uniform (a,b) dist's

so use $c_n = \min\left(n, \left\lceil n \left[1 - \delta + 1.12 \sqrt{\frac{\delta}{n}} \right] \right\rceil \right)$ 149.5

26) This large sample $100(1-\delta)\%$ Shorth(c_n) PI is for iid RVs.

27) For MLR with $\hat{\beta}_E$ (E an estimator like OLS, OPLS, MMLE, RR, lasso; forward selection, lassos vs), and $\underline{y} = \underline{X} \underline{\beta} + \underline{e}$,

find the residuals $\underline{r} = \underline{y} - \hat{\underline{y}} = \underline{y} - \underline{X} \hat{\beta}_E$.

Roughly, find Shorth(c_n) of the r_i ,

say $[r_{(s)}, r_{(n-c+1)}]$, which is a

PI for e_f . Then $[\hat{y}_f + r_{(s)}, \hat{y}_f + r_{(n-c+1)}]$

is a large sample $100(1-\delta)\%$ PI for y_f

where $y_f = \underline{x}_f^T \underline{\beta}_E$ if $\hat{\beta}_E$ is a consistent estimator of $\underline{\beta}_E$ and

$\underline{y} = \underline{x}_E^T \underline{\beta}_E + \underline{e}$ where \underline{e} depends on E .

$$28) \quad Y = X\beta + \varepsilon, \quad \underline{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}.$$

The elastic net estimator $\hat{\underline{\beta}}_{EN}$

minimizes the criterion

$$Q_{EN}(\underline{\beta}) = \frac{1}{2} \text{RSS}(\underline{\beta}) + \lambda_{in} \left[\frac{1}{2} (1-\alpha) \|\underline{\beta}_2\|_2^2 + \alpha \|\underline{\beta}_2\|_1 \right]$$

$$\text{or } Q_2(\underline{\beta}) = \text{RSS}(\underline{\beta}) + \lambda_1 \|\underline{\beta}_2\|_2^2 + \lambda_2 \|\underline{\beta}_2\|_1$$

where $0 \leq \alpha \leq 1$, $\lambda_1 = (1-\alpha)\lambda_{in}$, $\lambda_2 = 2\alpha\lambda_{in}$.

Then $\alpha=1$, $\lambda_1=0$ "corresponds to lasso" (using $a=0.5$ or $a=1$)

$\alpha=0$, $\lambda_2=0$ corresponds to $\tilde{\underline{\beta}}_{RR}$

$\lambda_{in}=0$, $\lambda_1=\lambda_2=0$ corresponds to $\underline{\beta}_{OLS}$.

29] EN does variable selection, can be used in HD, and has large sample theory similar to that of RR and lasso.

§ 3.11 K-component estimators

30] Let $y = \alpha + \underline{x}^T \underline{\beta} + e$.

The k -component MLR estimators use p linear combinations $\underline{n}_1^T \underline{x}, \underline{n}_2^T \underline{x}, \dots, \underline{n}_p^T \underline{x}$.

Then there are p conditional distributions

$$\begin{aligned}
 & y | \underline{n}_1^T \underline{x} \\
 & y | (\underline{n}_1^T \underline{x}, \underline{n}_2^T \underline{x}) \\
 & \vdots \\
 & y | (\underline{n}_1^T \underline{x}, \dots, \underline{n}_p^T \underline{x}).
 \end{aligned}$$

↑ want k small for dimension reduction

Estimating the \underline{n}_i and performing the OLS reg of y on $(\hat{\underline{n}}_1^T \underline{x}, \dots, \hat{\underline{n}}_k^T \underline{x})$ gives the k -component estimator $\hat{\underline{\beta}}_{kE}$, eg the k -component PLS estimator $\hat{\underline{\beta}}_{kPLS}$ or the k -component PCR estimator $\hat{\underline{\beta}}_{kPCR}$.

31] Let $y = \alpha + \underline{x}^T \underline{\beta} + e$ and $\underline{X} = (\underline{x}_1)$.

Let $\underline{v}_i = \hat{A}_{kn} \underline{x}_i = \begin{pmatrix} \hat{n}_1^T \underline{x}_i \\ \vdots \\ \hat{n}_k^T \underline{x}_i \end{pmatrix}$, $\hat{A}_{kn} = \begin{pmatrix} \hat{n}_1^T \\ \vdots \\ \hat{n}_k^T \end{pmatrix}$.

$k \times 1$ $k \times p$ $p \times 1$

not $p \times 1$

Let $\underline{c}_i = \sum_1 \hat{\underline{n}}_i = \begin{pmatrix} \underline{x}_1^T \hat{\underline{n}}_i \\ \vdots \\ \underline{x}_n^T \hat{\underline{n}}_i \end{pmatrix}$ be the HD 51

i -th component vector for $i=1, \dots, p$. Let

$$\underline{V}_k = (\underline{c}_1, \dots, \underline{c}_k) = \begin{pmatrix} \underline{v}_1^T \\ \vdots \\ \underline{v}_n^T \end{pmatrix} = \sum_1 \hat{\underline{A}}_{kn}^T$$

for $k=1, \dots, p$.

Let the working OLS model for regressing y on $\underline{v} = \hat{\underline{A}}_{kn} \underline{x}$ be

$$\underline{y} = \alpha_k \underline{1} + \underline{V}_k \underline{\gamma}_k + \underline{\varepsilon} \quad \text{where}$$

$\underline{\varepsilon}$ depends on the model.

The OLS regression of y on $\underline{v} = \hat{\underline{A}}_{kn} \underline{x}$

$$\text{gives } \underline{\gamma}_k = \underline{\hat{\gamma}}_{\underline{v}} = \underline{\hat{\gamma}}_{\underline{v}y} = (\hat{\underline{A}}_{kn} \hat{\underline{\Sigma}}_x \hat{\underline{A}}_{kn}^T)^{-1} \hat{\underline{A}}_{kn} \hat{\underline{\Sigma}}_{xy}$$

$$\therefore \underline{\hat{\beta}}_{KE} = \hat{\underline{A}}_{kn}^T \hat{\underline{\gamma}}_k =$$

$$\hat{\underline{A}}_{kn}^T (\hat{\underline{A}}_{kn} \hat{\underline{\Sigma}}_x \hat{\underline{A}}_{kn}^T)^{-1} \hat{\underline{A}}_{kn} \hat{\underline{\Sigma}}_{xy} = \hat{\underline{\Lambda}}_k \hat{\underline{\Sigma}}_{xy}$$

$$= \hat{\underline{A}}_{kn}^T (\hat{\underline{A}}_{kn} \hat{\underline{\Sigma}}_x \hat{\underline{A}}_{kn}^T)^{-1} \hat{\underline{A}}_{kn} \hat{\underline{\Sigma}}_x \underline{\hat{\beta}}_{OLS} = \hat{\underline{\Lambda}}_k \hat{\underline{\Sigma}}_x \underline{\hat{\beta}}_{OLS}$$

if $\hat{\underline{A}}_{kn}^T \hat{\underline{\Sigma}}_x \hat{\underline{A}}_{kn} = \underline{A}_k$

$$\underline{y}_i = \alpha_k + \underline{x}_i^T \hat{\underline{A}}_{kn} \underline{\gamma}_k + \varepsilon_i \quad \underline{y}_i = \alpha_k + \underline{x}_i^T \underline{\hat{\beta}}_{KE} + \varepsilon_i$$

If $\hat{m}_i \xrightarrow{P} m_i$ and $\hat{A}_k \xrightarrow{P} A_k = \begin{pmatrix} m_1^T \\ \vdots \\ m_k^T \end{pmatrix}$ 51.5
 then $\hat{\beta}_{KFE} \xrightarrow{P} \beta_{KFE} =$

$$A_k^T (A_k \Sigma_x A_k^T)^{-1} A_k \Sigma_x \beta_{OLS}(x, y) = A_k \Sigma_x \beta_{OLS}(x, y).$$

32) $\hat{\beta}_{KFE} = \hat{\beta}_{OLS}(x, y)$ if the inverse matrices exist (also if $p=1$). see HW7

33) If $\beta_{OLS} = \sum_{j=1}^m c_{ij} m_{ij}$ where $1 \leq m < p$,

and $c_{ij} \neq 0$, then $\beta_{KFE} = \beta_{OLS}$ if k is large enough to contain the m vectors m_{ij} . This condition can hold to a good approx in low dimensions.

34) In high dimensions, want $k < n-1$ and likely want $n \geq Jk$ with $J \geq 5, 10$ etc.

35) Use model selection, such as 10-fold CV, to select k^* . Call the estimator

$$\hat{\beta}_{MS, E}.$$