

39) Ridge regression, lasso, and elastic net can be used if $n \geq p$ and if $p > n$. Lasso and elastic net do variable selection:

$\hat{\beta}_L$ and $\hat{\beta}_E$ have 0's. Fit OLS to the predictors, including a constant, with $\hat{\beta}_E \neq 0$.

ex) $\hat{\beta}_L = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T \rightarrow$ fit OLS regression of y on $(x_1=1, x_3)$.

40) For $\underline{y} = X\underline{\beta} + \underline{e}$ and $\underline{z} = W\underline{\eta} + \underline{e}$

the OLS estimator $\hat{\underline{\eta}} = (W^T W)^{-1} W^T \underline{z}$.

41) The ridge regression estimator $\hat{\underline{\eta}}_R$

minimizes $Q_R(\underline{\eta}) = \frac{1}{a} (\underline{z} - W\underline{\eta})^T (\underline{z} - W\underline{\eta}) + \frac{\lambda_{in}}{a} \sum_{i=1}^{p-1} \eta_i^2$
 $= \frac{1}{a} RSS_{W(\underline{\eta})} + \frac{\lambda_{in}}{a} \|\underline{\eta}\|_2^2$ with $a = 1, 2, n, 2n$ common.

If $\lambda_{in} = 0$, $\hat{\underline{\eta}}_R = \hat{\underline{\eta}}_{OLS}$. As $\lambda_{in} \rightarrow \infty$,

$\hat{\underline{\eta}}_R \rightarrow 0$ and $\hat{\underline{y}} \rightarrow \bar{y}$. So RR is a shrinkage method.

$$\hat{\underline{\eta}}_R = \underbrace{(W^T W + \lambda_{in} I_{p-1})^{-1}}_{\text{inverse exists if } \lambda_{in} > 0} W^T \underline{z}.$$

Warning: The literature typically uses $\lambda = \frac{\lambda_n}{n}$.

42) Usually a grid of m λ_n values is used
 $0 \leq \lambda_1 < \lambda_2 < \dots < \lambda_m$ where $\lambda_i = \lambda_{ni}$.

10 fold CV is used to select $\lambda_g = \hat{\lambda}_m$.

43) Suppose $n > p$ and $(W^T W)^+$ exists.

$$\text{Then } \hat{\underline{m}}_R = \overbrace{\left(W^T W + \lambda_n I_{p-1} \right)^{-1}}^{A_n} \underbrace{W^T W (W^T W)^+}_{I_{p-1}} W^T \underline{z} \overbrace{\hat{\underline{m}}_{OLS}}^{n}$$

So $\hat{\underline{m}}_R = A_n \hat{\underline{m}}_{OLS}$. HW 10 shows

$$A_n = B_n = I_{p-1} - \lambda_n (W^T W + \lambda_n I_{p-1})^{-1}$$

44) Suppose $\frac{W^T W}{n} \xrightarrow{P} V^{-1}$ and the OLS CLT holds:

$$\sqrt{n} (\hat{\underline{m}}_{OLS} - \underline{\eta}) \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V).$$

If $\frac{\lambda_n}{n} \rightarrow 0$, then $\frac{W^T W + \lambda_n I_{p-1}}{n} \xrightarrow{P} V^{-1}$ and

$n (W^T W + \lambda_n I_{p-1})^{-1} \xrightarrow{P} V$. Also $A_n = \frac{n}{n} A_n =$

$$\left(\frac{W^T W + \lambda_n I_{p-1}}{n} \right)^{-1} \frac{W^T W}{n} \xrightarrow{P} V V^{-1} = I_{p-1}.$$

45} RR CLT Assume p is fixed and

the OLS CLT holds for $\underline{z} = W\underline{\eta} + \underline{\epsilon}$.

a) If $\frac{\hat{\lambda}_{in}}{\sqrt{n}} \xrightarrow{p} 0$, $\sqrt{n}(\hat{\underline{\eta}}_R - \underline{\eta}) \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V)$.

So OLS full model and RR are asymptotically equivalent.

b) If $\frac{\hat{\lambda}_{in}}{\sqrt{n}} \xrightarrow{p} \gamma > 0$, then

$$\sqrt{n}(\hat{\underline{\eta}}_R - \underline{\eta}) \xrightarrow{D} N_{p-1}(\underbrace{-\gamma V \underline{\eta}}_{\text{interior to full model OLS unless this term} = 0}, \sigma^2 V)$$

interior to full model OLS unless this term = 0

c) $\hat{\underline{\eta}}_R$ is a consistent estimator of $\underline{\eta}$ if $\frac{\hat{\lambda}_{in}}{n} \xrightarrow{p} 0$ as $n \rightarrow \infty$.

§7.9

46} $\underline{y} = X\underline{\beta} + \underline{\epsilon}$, $\underline{z} = W\underline{\eta} + \underline{\epsilon}$. The lasso estimator

$\hat{\underline{\eta}}_L$ minimizes the criterion

$$Q_L(\underline{\eta}) = \frac{1}{a} (\underline{z} - W\underline{\eta})^T (\underline{z} - W\underline{\eta}) + \frac{\lambda_{in}}{a} \sum_{i=1}^p |\eta_i| =$$

$$\frac{1}{a} \text{RSS}_W(\underline{\eta}) + \frac{\lambda_{in}}{a} \|\underline{\eta}\|_1, \quad \leftarrow L_1 \text{ norm.}$$

Lasso is a shrinkage method and a variable selection method; often some $\hat{\eta}_i = 0$.

47) is used. λ_m is the smallest value of λ such that $\hat{\eta}_{\lambda} = \underline{0}$. Hence $\hat{\eta}_{\lambda_i} \neq \underline{0}$ for $i < m$.

48) By the Karush-Kuhn-Tucker (KKT) conditions for convex optimality (see Math 471)

$$-W^T (\underline{z} - W \hat{\underline{\eta}}_L) + \frac{\lambda_n}{2} \underline{s}_n = \underline{0} \text{ where } s_n \in [-b, b].$$

$$\text{Thus } \hat{\underline{\eta}}_L = \underbrace{(W^T W)^{-1} W^T \underline{z}}_{\hat{\underline{\eta}}_{OLS}} - n (W^T W)^{-1} \frac{\lambda_n}{n} \underline{s}_n$$

49) Lasso CLT: Assume p is fixed and the OLS CLT holds for $\underline{z} = W \underline{\eta} + \underline{e}$,

a) If $\frac{\hat{\lambda}_n}{\sqrt{n}} \xrightarrow{P} 0$, then $\sqrt{n} (\hat{\underline{\eta}}_L - \underline{\eta}) \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V)$,

b) If $\frac{\hat{\lambda}_n}{\sqrt{n}} \xrightarrow{P} \tau > 0$, then $\sqrt{n} (\hat{\underline{\eta}}_L - \underline{\eta}) \xrightarrow{D} N_{p-1}(-\frac{\tau}{2} V \underline{s}_w, \sigma^2 V)$,
 and $\underline{s}_w \in B_{\tau} \underline{s}$

c) If $\frac{\hat{\lambda}_n}{n} \xrightarrow{P} 0$, then $\hat{\underline{\eta}}_L$ is a consistent estimator of $\underline{\eta}$,

§ 7.11 90) $\underline{y} = X \underline{\beta} + \underline{e}$, $\underline{z} = W \underline{\eta} + \underline{e}$

The elastic net estimator

$\hat{\beta}_{EN}$ minimizes

the criterion $Q_{EN}(\beta) = RSS_{w(\beta)} + \lambda_1 \|\beta\|_2^2 + \lambda_2 \|\beta\|_1$,

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

so $\alpha = 1$ corresponds to lasso using $2\lambda_{in}$ and

$\alpha = 0$ corresponds to RR. The EN estimator has a CLT similar to that of lasso, and

EN is a shrinkage estimator that does variable selection.

51) The VS CLT showed

$\hat{\beta}_{VS}$ is \sqrt{n} consistent when lasso or EN are consistent where I_{min} corresponds to the lasso or EN $\beta_i \neq 0$.

52) know Consider choosing $\hat{\beta}$ to minimize

$Q(\beta) = \frac{1}{n} (z - w\beta)^T (z - w\beta) + \frac{\lambda_{in}}{n} \sum_{i=1}^{p-1} |\beta_i|^j$ where

$\lambda_{in} \geq 0$, $n > 0$ and $j > 0$ are known constants.

Then $j=2$ corresponds to RR $\hat{\beta}_R$, $j=1$ to lasso $\hat{\beta}_L$,

and $\lambda_{in} = 0$ to OLS $\hat{\beta}_{OLS}$. Do a similar problem with EN criterion from 50).

53} Get outlier resistant estimators by replacing $\hat{\beta}_{OLS}$ by the alternative estimator $\hat{\beta}_A$.

i) RMVN, RFCH, COUMB2 set D applied to nontrivial predictors, then fit $\hat{\beta}_A$ on cases corresponding to D.

ii) tuning using $\hat{\beta}_A$ instead of $\hat{\beta}_{OLS}$

iii) RMVN set U from (Y_i, X_i)

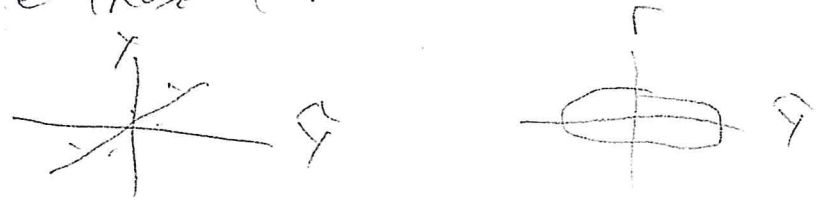
Fit $\hat{\beta}_A$ to cases in U.

Ch 8 §8.1

1} The additive error regression (AER) model

is $Y = m(X) + e$ where $m(X)$ is a real valued function, $E(e) = 0$, $V(e) = \sigma^2$, $\hat{Y} = \hat{m}(X)$, $r = Y - \hat{Y}$.
 The AER model is a 1D reg model with $SP = m(X) = \hat{m}(X)$, $ESP = \hat{m}(X)$

2} If the error distribution is not highly skewed, the response and residual plots look exactly like those for MLR

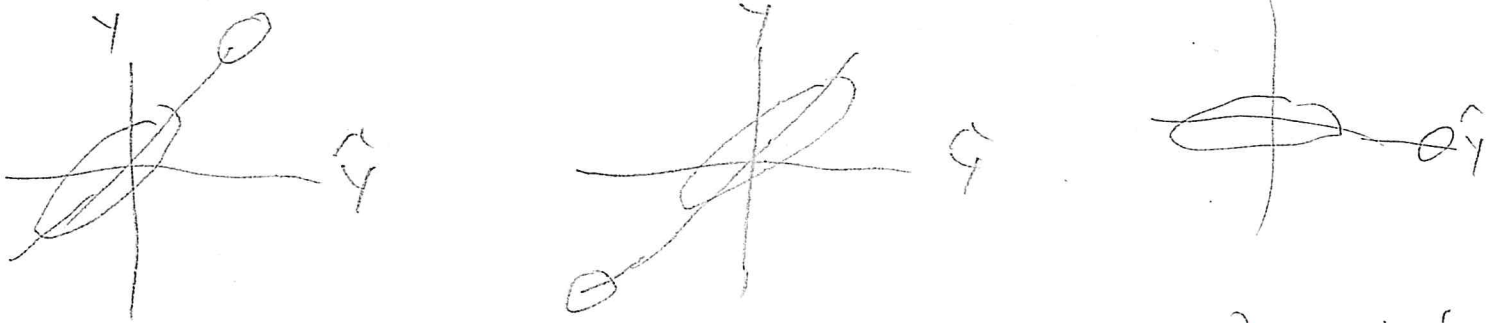


3) MLR has $\max | = \underline{x}^T \underline{\beta}$. Non-linear regression, Nonparametric regression and the AE Gam

58

$Y = \alpha + \sum_{i=1}^p s_i(x_i) + e = AP + \epsilon$, are also special cases of AER.

4) Compared to the inflexible hyperplane $\underline{x}^T \underline{\beta}$, several methods for finding $\hat{m}(x)$ are flexible. Hence $\hat{y}_i \approx y_i$ often occurs even if y_i is an outlier. Thus outlier detection is more difficult.



Look for gaps in the response and residual plots.

5) PIs are similar. Let $[r_L, r_U]$ be a PI for \mathbb{R} .

Then a PI for \mathbb{Y} is $[\hat{Y} + r_L, \hat{Y} + r_U]$.

ch 93 13 Good references

983 59

Cook and Weisberg 1999a

and Olive 2010a,
2017a

2) Know $p > 2$ Regression is the study of the conditional distribution $Y | \underline{x}$ of the response Y given the $(p-1) \times 1$ vector of non-trivial predictors \underline{x} (now \underline{x} does not contain a constant 1).

want to know how Y changes as the predictors \underline{x} change. ($Y = Y(\underline{x})$)

suppressed

NO 59.5

3) Know: In a LD regression model ... 60

Y is conditionally independent of X given $\alpha + c\beta^T X$ often $h(x) = \alpha + c\beta^T x$ a single linear combination of the predictors, written $Y \perp\!\!\!\perp X \mid \beta^T X$.

4) Know p 450: If $Y \perp\!\!\!\perp X \mid \beta^T X$, then

$Y \perp\!\!\!\perp X \mid (\alpha + c\beta^T X)$ for any constants

α and $c \neq 0$. The quantity $\alpha + c\beta^T X$ is called a sufficient predictor (SP).

A sufficient summary plot (SSP) is a plot

of SP vs Y . An estimated sufficient predictor (ESP) is $\tilde{\alpha} + \tilde{\beta}^T X$ where

$\tilde{\beta}$ is an estimator of $c\beta$ for some $c \neq 0$.

An estimated sufficient summary plot (ESSP)

or response plot is a plot of ESP vs Y .

5) * p 453 The most used statistical regression models are LD models.

ex] $y = g(\alpha + \beta^T X, e)$

$$y = m(\alpha + \beta^T x) + e$$

(often the function m is unknown)

ex) $y = \alpha + \beta^T x + e$ (MLR)

ex) $y = A^{-1}(\alpha + \beta^T x + e)$ (transformation model)

ex) logistic regression in ch 10

ex) poisson regression in ch 10

6) Any model with a single (non-trivial) predictor x ($p-1=1$) is a 1D model. In this case x is both a SP and ESP and a plot of x vs y is both a SSP and an ESSP.

7) If $p-1 > 1$ the SP is unknown and an SSP can not be made (except for simulated data). For the MLR model $y = \alpha + \beta^T x + e$ the ESP = $\hat{y} = \hat{\alpha} + \hat{\beta}^T x$, and the ESSP = forward response plot: \hat{y} vs x .