$\widehat{Y}_f$ depends on $Y_1, ..., Y_n$ so $\widehat{Y}_f \perp\!\!\!\perp Y_f$

$$= V(\underline{x}_f' \hat{\underline{\beta}}) + \sigma^2 = V(\underline{x}_f' (\underline{X}'\underline{X})^{-1} \underline{X}' \underline{Y}) + \sigma^2$$

$$= \left[ \underline{x}_f' (\underline{X}'\underline{X})^{-1} \underline{X}' \sigma^2 I \underline{X} (\underline{X}'\underline{X})^{-1} \underline{x}_f \right] + \sigma^2$$

$$= \sigma^2 \left( \underline{x}_f' (\underline{X}'\underline{X})^{-1} \underline{x}_f + 1 \right) = \sigma^2 (1 + h_f)$$

$$\underbrace{\qquad\qquad\qquad}_{h_f}$$

16} $^{p132}$ If the $\varepsilon_i$ are iid $N(0, \sigma^2)$, then

a $100(1-\delta)\%$ prediction interval (PI)

for the random variable $Y_f$ is

$$\widehat{Y}_f \pm t_{n-p, 1-\frac{\delta}{2}} \sqrt{MSE} \sqrt{1 + h_f} \qquad \text{(closed interval)}$$

want $h_f \equiv \max h_i$, $h_i = \underline{x}_i' (\underline{X}'\underline{X})^{-1} \underline{x}_i$

$= i$th diagonal **entry** of $H = P = \underline{X}(\underline{X}'\underline{X})^{-1}\underline{X}'$

17} As $n \to \infty$ the PI in 16} estimates

$$\left[ E(Y_f | \underline{x}_f) - z_{1-\frac{\delta}{2}} \sigma \text{ or } E(Y_f | \underline{x}_f) + z_{1-\frac{\delta}{2}} \sigma \right], \text{ the}$$

highest $1-\delta$ density region if
$$Y_f \mid \underline{x_f} \sim N\left(E(Y_f \mid \underline{x_f}), \sigma^2\right).$$

18) ~~know~~ A *large sample* $100(1-\delta)\%$ PI

$[L_n, U_n]$ satisfies $P\{Y_f \in [L_n, U_n]\} \to 1-\delta$

as $n \to \infty$. $\qquad$ actually $Y_f \mid \underline{x_f}$ but suppress $\underline{x_f}$

19) Suppose $Y_f \mid \underline{x_f}$ (... has pdf $f_f(y)$ and cdf $F_f(y)$

Want $[L_n, U_n] \xrightarrow{P} [L, U]$ where $F_f(U) - F_f(L) = 1-\delta$

and $U-L$ is short.

20) The highest density region is found
by moving a horizontal line down from
the top of the pdf. The line will intersect
the pdf or boundaries of the support
of the pdf at $[a_1, b_1], ..., [a_M, b_M]$ for
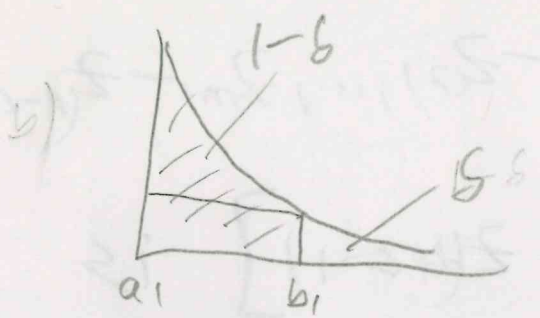some $M \geq 1$. Stop moving the line when the
areas under the pdf corresponding to
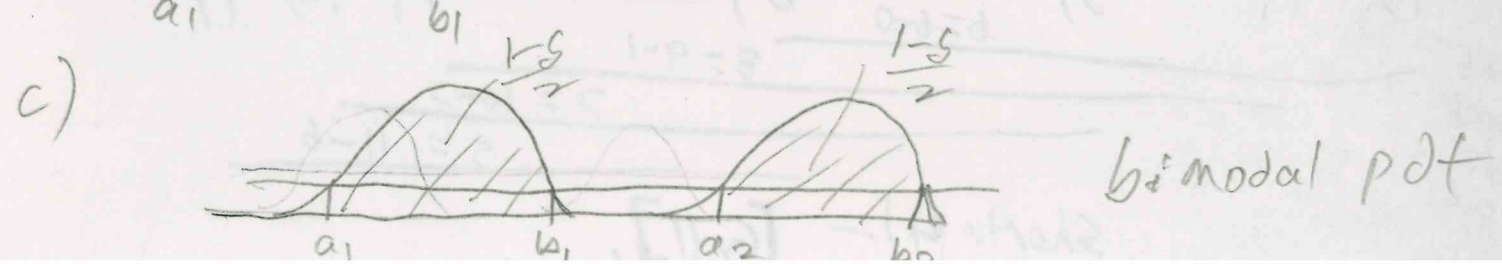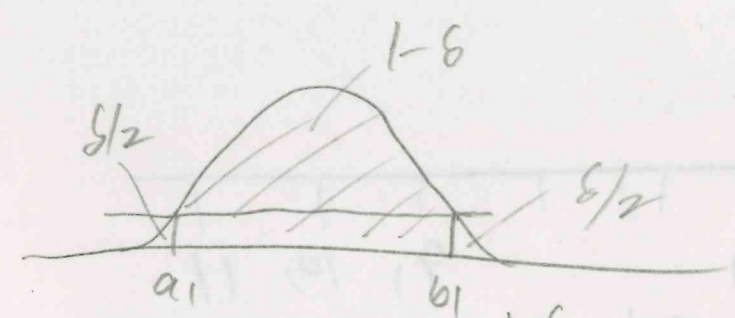the intervals is $(1-\delta)$. Often the pdf

is unimodal and decreases rapidly as $y$ moves away from the mode. Then $H=1$ and the highest density region is an interval.

ex) a) If $Y_f$ has an exponential distribution, then the highest density region is

$$[0, \xi_{1-\delta}] \quad \text{where} \quad P(Y_f \in \xi_\delta) = \alpha$$

$$\uparrow$$
$$\xi_i \ (z_i)$$



b) For a symmetric unimodal distribution, the highest density region is $[\xi_{\delta/2}, \xi_{1-\delta/2}]$.



c)



bimodal pdf

2) Suppose you have data $z_1, \ldots, z_n$. The order statistics $z_{(1)} \leq z_{(2)} \leq \ldots \leq z_{(n)}$
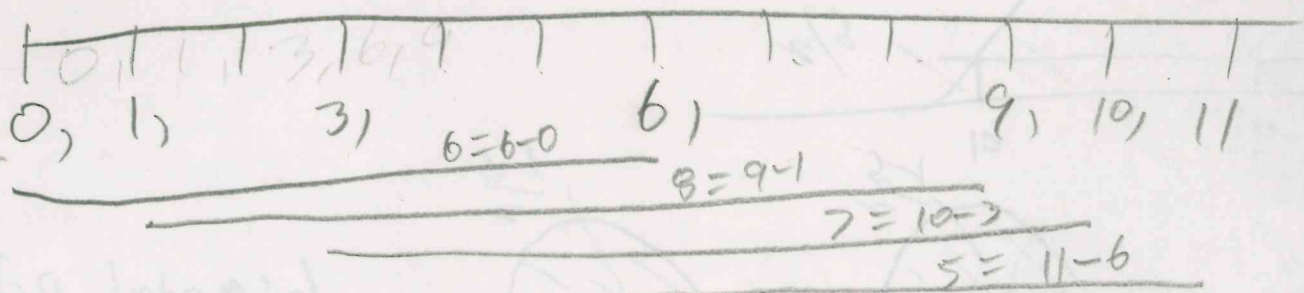
$$z_{(1)} = \min(z_1, \ldots, z_n)$$
$$z_{(n)} = \max(z_1, \ldots, z_n),$$

Consider intervals that contain $c$ cases $[z_{(1)}, z_{(c)}], [z_{(2)}, z_{(c+1)}], \ldots, [z_{(n-c+1)}, z_{(n)}]$.

Compute $z_{(c)} - z_{(1)}, z_{(c+1)} - z_{(2)}, \ldots, z_{(n)} - z_{(n-c+1)}$.

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

ex) know for quiz Let $C=4$. Data below has $n=7$.



0, 1, 3, 6=6-0 6, 9, 10, 11

8 = 9-1
7 = 10-3
5 = 11-6

Shorth(4) = [6,11].

22) If $Y_1, ..., Y_n$ are iid

$$\underline{Y} = \underline{1} \, \beta_0 + \underline{e} \quad \text{and} \quad \frac{c}{n} \to 1-\delta$$

eg $c = k_n = \lceil n(1-\delta) \rceil$, then the

shorth $(c)$ interval estimates the

highest density $100(1-\delta)\%$ region if that

region is an interval. Then the "shorth$(c)$

interval can be used as as a large sample

$100(1-\delta)\%$ PI

for $Y_f$. If $c = k_n$ then for large $n \delta$

and iid data) the shorth PI has max.

undercoverage $\approx 1.12 \sqrt{\delta/n}$. So using

$$c = \lceil n [1-\delta + 1.12 \sqrt{\tfrac{\delta}{n}}] \rceil \quad \text{works better}$$

than $c = k_n$. (Frey 2013)

23) Let $a_n = \left(1 + \frac{15}{n}\right) \sqrt{\frac{n}{n-p}} \sqrt{1 + h_f}$

Let $c = k_n$ and find the shorth estimator

applied to the residuals ie $e_1, ..., e_n$.

So short$(C) = \left[ e_{(d)}, e_{(d+1)} \right] = \left[ \tilde{\xi}_{\delta_1}, \tilde{\xi}_{1+\delta_2} \right]$.

Let $\underline{Y} = \underline{X}\underline{\beta} + \underline{\xi}$, $E\underline{\xi} = 0$, $cov(\underline{\xi}) = \sigma^2 I$.

Then a large sample $100(1+\delta)\%$ PI for

$Y_f$ is $\left[ \hat{q}_f + a_n \tilde{\xi}_{\delta_1}, \hat{q}_f + a_n \tilde{\xi}_{1+\delta_2} \right]$.

This PI is asymptotically optimal (short)

if the $x_i$ are bounded in probability and the

iid $\xi_i$ come from a large class of

zero mean unimodal distributions.

---

Skip § 5.3.2, 5.4

---

ch 9  1) $\underline{Y} = \underline{X}\underline{\beta} + \underline{\xi}$        $\underline{X}$ full rank

$E(\underline{\xi}) = \underline{0}$        $cov(\underline{\xi}) = \sigma^2 I$.

We have i) assumed the model is correct,

at least approximately. In practice this

assumption is often violated.

2} i) Could have $E Y \neq X \beta$ if

$X$ is missing important predictor variables.

ii) Could have $cov(\underset{\sim}{\varepsilon}) \neq \sigma^2 I$

iii) The $\varepsilon_i$ could be correlated instead of iid (uncorrelated).

§ 9.2  3} Suppose we fit model

$$Y = X \underset{\sim}{\beta} + \underset{\sim}{\varepsilon},$$ but the true model

is $$\underline{Y} = X \underline{\beta} + W \underset{\sim}{\gamma} + \underline{\varepsilon}$$
$\uparrow$

columns of $W$ should be in the model

Where the columns of the $n \times q$ full rank matrix $W$ are linearly independent of the columns of the full rank $n \times p$ matrix $X$.

Then $E \hat{\underset{\sim}{\beta}} = E \left[ (X'X)^{-1} X' \underline{Y} \right] =$

$(X'X)^{-1} X' \left[ \underline{X \beta} + W \underset{\sim}{\gamma} \right] = \underset{\sim}{\beta} + (X'X)^{-1} X' W \underset{\sim}{\gamma}$

$= \underset{\sim}{\beta} + L \underset{\sim}{\gamma}.$  So $\hat{\underset{\sim}{\beta}}$ is

a biased estimator of $\beta$ with bias $L\underset{\sim}{\gamma}$. This bias could be quite large, but $L\underset{\sim}{\gamma}=\underset{\sim}{0}$ if $\underline{X}'\underline{W}=0$ ie if the columns of $\underline{W}$ are orthogonal to the columns of $\underline{X}$.

4) Leaving out important predictors can destroy the linearity of the model

ex) $\overset{Fit}{Y_i}=\beta_0+\beta_1 x_{1i}+\xi_i$ when the true model is $Y_i=\beta_0+\beta_1\underbrace{x_{1i}}+\beta_2\underbrace{x_{1i}^2}+\xi_i$

$$x_{i1}=x_i \quad x_{i2}=x_{i1}^2$$

cook and weisberg P264-5

ex) Suppose $\underline{x}=\begin{pmatrix}x_1\\x_{p-1}\end{pmatrix}$, $\underline{\beta}=\begin{pmatrix}\beta_1\\\beta_{p-1}\end{pmatrix}$, and

$E(Y|\underline{x})=\underline{\beta_1}^T\underline{x_1}+\beta_{p-1}x_{p-1}\underline{x}=\underline{\beta}^T\underline{x}$.

consider regressing $Y$ on $\underline{x_1}$, so without $x_{p-1}$.

Then $E(Y|\underline{x_1}) = E\left[\overline{E(Y|\underline{x_1}, x_{p-1})}|\underline{x_1}\right]$

nontrivial fact: $EW = E(E[W|x_{p-1}])$

but $E(W|x_{p-1}) = E(Y|\underline{x_1}, x_{p-1})$    take $W = Y|\underline{x_1}$

$= E\left(\underline{\beta_1}^T \underline{x_1} + \beta_{p-1} x_{p-1} \mid \underline{x_1}\right)$

$= \underline{\beta_1}^T \underline{x_1} + \beta_{p-1} \underline{E(x_{p-1}|\underline{x_1})}$

$x_{p-1}$ gets replaced by $E(x_{p-1}|\underline{x_1})$

when $x_{p-1}$ is omitted from the LS model.

ex} $E(Y|\underline{x}) = 1 + 2x_1 + 3x_2,$    $V(Y|\underline{x}) = \sigma^2$.

Then $E(Y|x_1) = 1 + 2x_1 + 3 E(x_2|x_1)$

a) If $x_1 \perp\!\!\!\perp x_2$ then $E(x_2|x_1) = E(x_2)$

and $E(Y|x_1) = (1 + 3 E(x_2)) + 2x_1$

is linear.    The coefficient for $x_1$

does not change but the intercept does.

b) Suppose $E(x_2|x_1) = \alpha_0 + \alpha_1 x_1$

Then $E(Y|x_1) = 1 + 2x_1 + 3\alpha_0 + 3\alpha_1 x_1$

$$= (1 + 3\alpha_0) + (2 + 3\alpha_1)x_1 \quad \text{which is}$$

again linear, but both the intercept and slope have changed.

c) If $E(x_2 | x_1) = \alpha_0 + \alpha_1 \exp(\alpha_2 x_1)$, then

$$E(y | x_1) = 1 + 2x_1 + 3\alpha_0 + 3\alpha_1 \exp(\alpha_2 x_1)$$

$$= (1 + 3\alpha_0) + 2x_1 + 3\alpha_1 \exp(\alpha_2 x_1)$$

which is a nonlinear mean function.

6) Under the conditions of 5),

$$V(Y | \underline{x_1}) = E\left[\overbrace{V(Y | \underline{x_1}, x_{p-1})}^{\text{correct linear model}} | \underline{x_1}\right] + V\left[\overbrace{E(Y | \underline{x_1}, x_{p-1})}^{} | \underline{x_1}\right]$$

$$= E(\sigma^2 | \underline{x_1}) + V\left[\underbrace{(\beta_1^T \underline{x_1}}_{\text{constant given } \underline{x_1}} + \beta_{p_1} x_{p-1}) | \underline{x_1}\right]$$

$$= \sigma^2 + \beta_{p_1}^2 V(x_{p-1} | x_1).$$

Hence deleting a term from the model may result in a nonconstant variance function.

For a linear model when $x_{p-1}$ is omitted,

7} want $E(x_{p-1} \mid \underline{x_1}) = \underline{\gamma}^T \underline{x_1}$

and $V(x_{p-1} \mid \underline{x_1}) = \tau^2$

so $E(Y \mid \underline{x_1}) = \underline{\beta_1}^T \underline{x_1} + \beta_{p-1} \underline{\gamma}^T \underline{x_1} = \underbrace{\underline{\eta}^T \underline{x_1}}_{\eta}$

$$\underline{\eta} = \underline{\beta_1} + \beta_{p-1} \underline{\gamma}$$

and $V(Y \mid \underline{x_1}) = \sigma^2 + \beta_{p-1}^2 V(x_{p-1} \mid \underline{x_1}) = 0$

$$\sigma^2 + \beta_{p-1}^2 \tau^2 = \theta^2 \qquad \text{say.}$$

8} If $\underline{x} = (1, \underbrace{x_2, \ldots, x_{p-1}}_{\underline{w'}})' = (1 \quad \underline{w'})'$

and $\begin{pmatrix} Y \\ x_2 \\ \vdots \\ x_{p-1} \end{pmatrix} \sim N_p \left( \underline{\mu}, \begin{pmatrix} \sigma_Y^2 & \Sigma_{Yw} \\ \Sigma_{wY} & \Sigma_w \end{pmatrix} \right)$

then $Y \mid x_{i1}, \ldots, x_{ik}$ follows a linear

model with constant variance

$$Y_i = \beta_{0k} + \beta_{1k} x_{i1} + \cdots + \beta_{kk} x_{ik} + \varepsilon_{ik}$$

$V(\varepsilon_{ik}) = \sigma_k^2$.   Models with lower $\sigma_k^2$ are better,

9) Having too many predictors (if $n \ge 10$) is much less serious than having too few, since the $\hat{\beta_i}$ for unneeded $x_i$ tend to have $\hat{\beta_i} \xrightarrow{P} 0$.

10) P230 Overfitting:
Suppose $E(\underline{Y}) = \underset{n \times k}{X_1 \beta_1}$ where

$$X = (X_1 \; X_2) \quad \text{and} \quad \underline{B} = \begin{pmatrix} \beta_1 \\ \underline{0} \end{pmatrix} \text{ since } \underline{\beta_2} = \underline{0}.$$

Then $X_1 \beta_1 = XB = (X_1 \; X_2)\begin{pmatrix} \beta_1 \\ \underline{0} \end{pmatrix}$,

So $E(\hat{\underline{B}}) = (X'X)^{-1}X' X_1 \beta_1 =$

$$(X'X)^{-1}X' X \begin{pmatrix} \beta_1 \\ \underline{0} \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \underline{0} \end{pmatrix} = \underline{B}.$$

Also $E\hat{\underline{Y}} = E(X\hat{\underline{B}}) = X\begin{pmatrix} \beta_1 \\ \underline{0} \end{pmatrix} = X\underline{B} = X_1 \underline{\beta_1}$

and $E(MSE) = \sigma^2$.

However $R^2$ is too high and the last $k$ diagonal elements of $(X'X)^{-1}$

are larger than the diagonal elements of $(X_1'X_1)^{-1}$. So CIs for $\beta_i$ using $X$ are longer than the CIs for $\beta_i$ using $X_1$, for $i = 1, ..., k$.

11) Basically overfitting is a correct linear model with one or more $\beta_i = 0$. So large sample inference is correct but not as precise as using the model that omits the predictors with $\beta_i = 0$.
want $n \geq 10k$ if $Y = \underset{n \times k}{X_1 \beta} + \varepsilon$, $n \geq 10P$ if $Y = \underset{n \times p}{X \beta} + \varepsilon$

9.3 12) Suppose $\underline{Y} = X\underline{\beta} + \underline{\varepsilon}$

$E(\underline{\varepsilon}) = \underline{0}$ but $Cov(\underline{\varepsilon}) = \sigma^2 V$ instead of $\sigma^2 I_n$. Then $\hat{\underline{\beta}} \overset{P}{\to} \underline{\beta}$ but

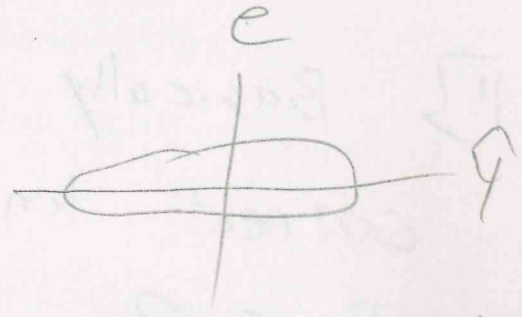$Cov(\hat{\underline{\beta}}) = \sigma^2 (X'X)^{-1} X' V X (X'X)^{-1} \neq \sigma^2 (X'X)^{-1}$.

Typically $E(MSE) \neq \sigma^2$.

Remedy: use GLS if $V$ is known, sandwich estimator

§9.4 13, Find outliers (cases far away from the bulk of the data) with response plots and residual plots.
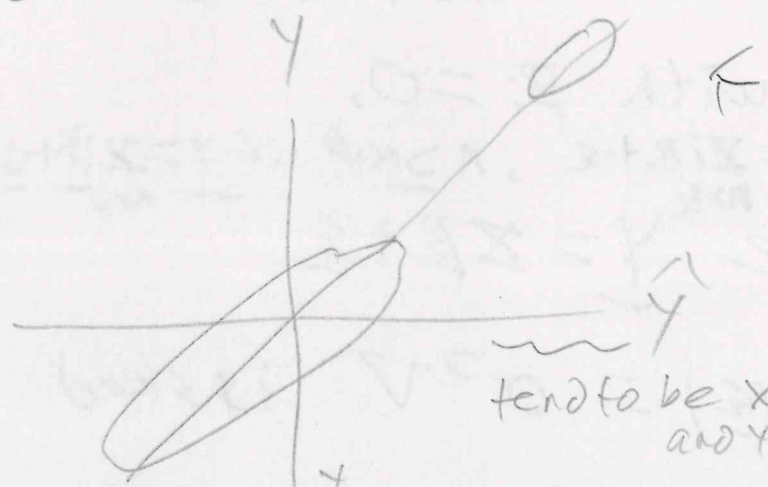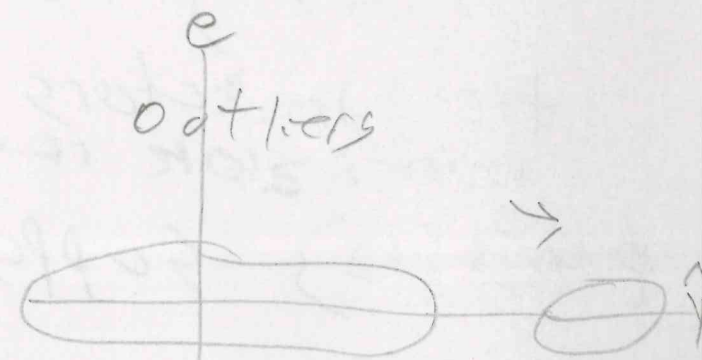
good response plot

good residual plot
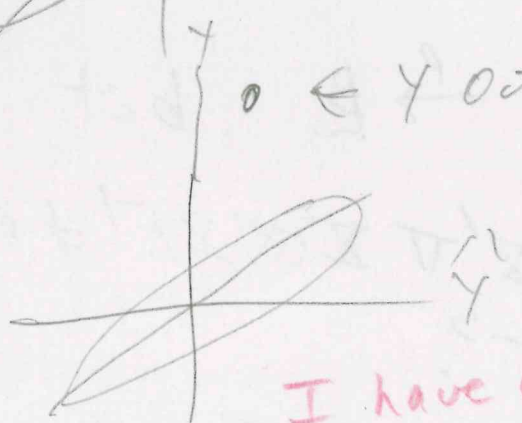
Outliers will often have gaps

outliers

tend to be X outliers and Y outliers

could be good leverage

← Y outlier far from bulk of Y's points

I have an R impact function rmreg2

rmreg2

14] Robust estimators can also be used,