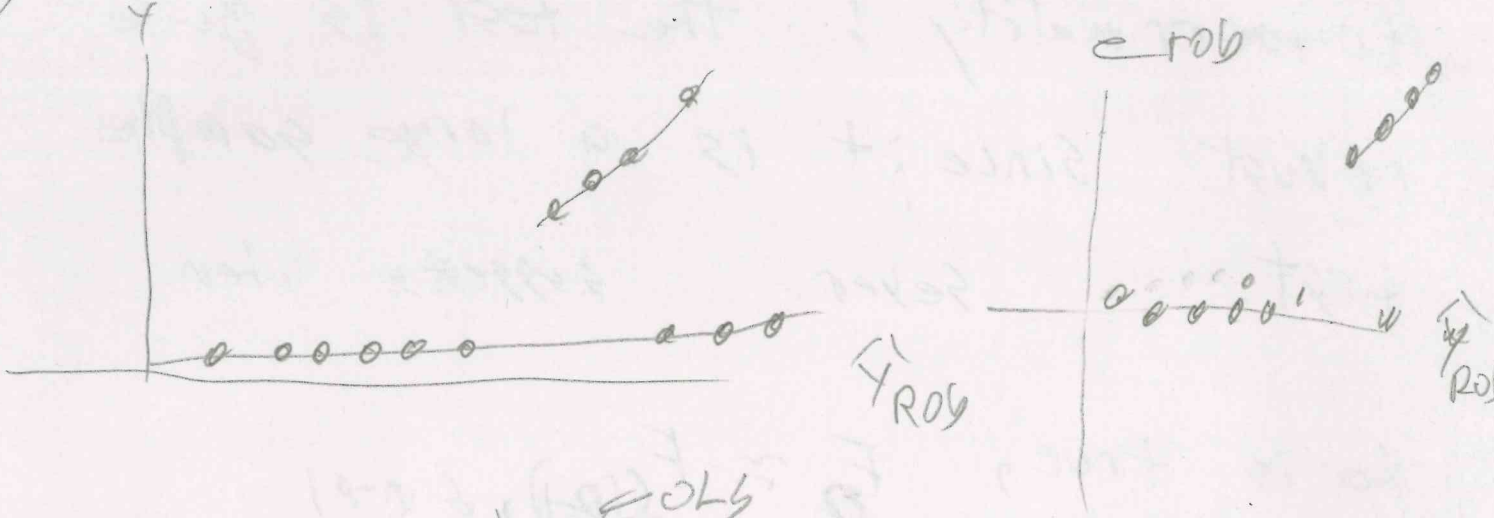
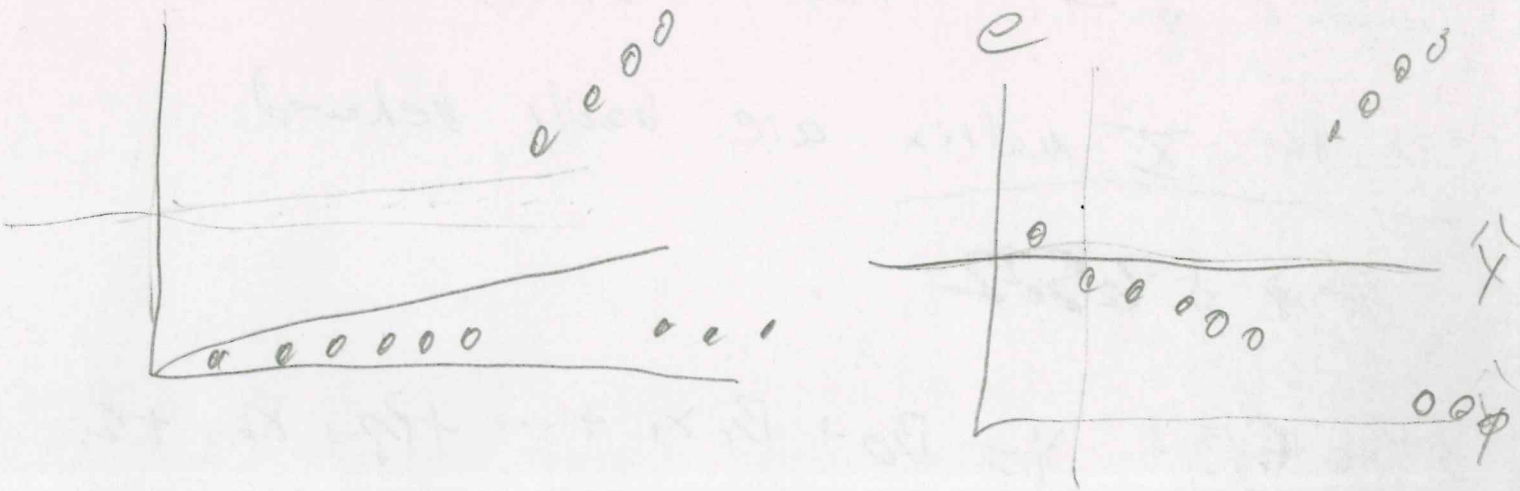


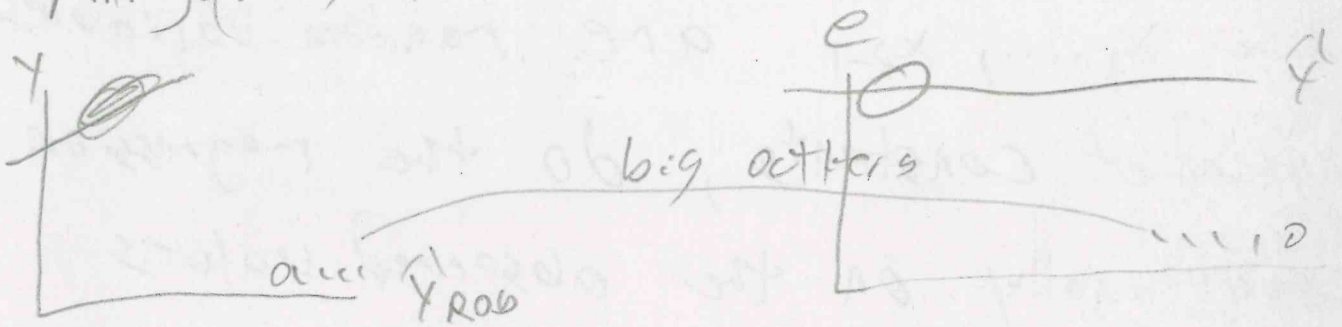
ex) $\text{lmreg2}(bex, bely)$



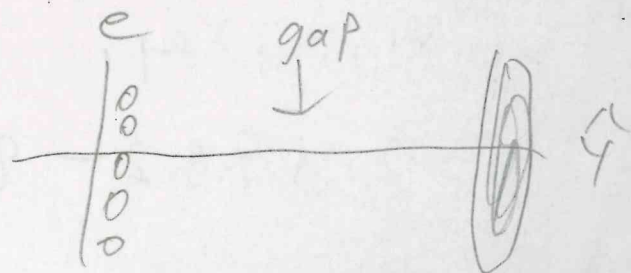
$\text{mlrplot4}(bex, bely)$ e_{OLS}



ex) $\text{lmreg2}(boxx, boxy)$



$\text{mlrplot4}(boxx, boxy)$
 $\nearrow \text{gap}$ \nwarrow \hat{e}_{it} goes through outliers



§ 9.5 15¹²³⁵ Robustness of ANOVA F test

to nonnormality; the test is quite robust since it is a large sample test. Seber suggests when

$$H_0 \text{ is true, } F_0 \approx F_{\delta(p-1), \delta(n-p)}$$

where $\delta \rightarrow \infty$ fast, unless the variables in the X matrix are badly behaved

Step § 9.5.2

§ 9.6 14 It $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1} + \epsilon_i$

where x_1, \dots, x_{p-1} are random variables instead of constants, do the regression conditionally on the observed values of x_1, \dots, x_{p-1} .

Step § 9.6.2 - 9.6.6

§9.7 17) p249 Collinearity occurs when

at least one column of $X = [v_0 \dots v_{p-1}]$ is highly correlated with a linear combination

of the other columns: $\text{corr}(v_j, \sum_{\substack{i=0 \\ i \neq j}}^{p-1} \delta_i v_i) = R_j$

for at least one j .

is high. Collinearity makes $X'X$ nearly singular.

18) Regress v_j on $(v_0, v_1, \dots, v_{j-1}, v_{j+1}, \dots, v_{p-1})$

and let R_j^2 be the coefficient of determination from the regression. R_j^2 is a measure of the collinearity of v_j with the other predictors.

19) $1 - R_j^2 \approx 0$ means $v_j \approx$ a linear

combination of the other columns of X

so $X'X$ is nearly singular.

20) p254-5 The j th variance inflation factor

$$VIF_j = \frac{1}{1 - R_j^2}$$
 Since $\text{Cov}(\beta) = \sigma^2 (X'X)^{-1}$
↑
p.249

$$SE(\hat{\beta}_j) = \sqrt{MSE} d_{jj} = \sqrt{MSE} (X'X)^{-1}_{j+j+1}$$

see p119 p.42

$$= \frac{\sqrt{MSE}}{SD(\tilde{v}_j) \sqrt{n-1}} \sqrt{\frac{1}{1-R_j^2}} = \frac{\sqrt{MSE} \sqrt{VIF_j}}{SD(\tilde{v}_j) \sqrt{n-1}}$$

$$SD(\tilde{v}_j) = \sqrt{\frac{\sum_{i=1}^n (v_{ij} - \bar{v}_j)^2}{n-1}} = \sqrt{\hat{V}(\tilde{v}_j)}$$

↑
sample variance

is large if $R_j^2 \approx 1$ or if VIF_j is large.

Then the CI for β_j will be long, and

it is harder to reject $H_0: \beta_j = 0$ so

the test has low power. Note that $VIF_j =$

$$\frac{\sum_{i=1}^n (v_{ij} - \bar{v}_j)^2}{\sum_{i=1}^n (v_{ij} - \bar{v}_j)^2} (X_j' X_j)^{-1}_{j+j+1} = (n-1) \hat{V}(\tilde{v}_j) d_{jj}$$

SKIP § 9.7.3, 9.7.4

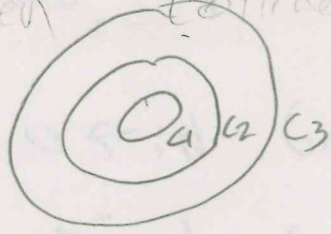
2) ^{p261} Collinearity has little effect on prediction.

$$\hat{y}_e = \underline{x}_e' \hat{\beta}, \quad V(\hat{y}_e) = \sigma^2 \frac{\underline{x}_e' (\underline{X}' \underline{X})^{-1} \underline{x}_e}{n}$$

$$= \sigma^2 \underline{x}_e' (\underline{X}' \underline{X})^{-1} \underline{x}_e = \sigma^2 h_e. \quad \text{Still have}$$

$\max h_i \rightarrow 0$. So $h_e \rightarrow 0$ if $h_e \leq \max(h_i)$.

When collinearity is low, $h_i = \text{constant}$



$h_i = c$
 when collinearity is low

$h_i = c$ when
 collinearity is high

2 predictors $\underline{x}_1, \underline{x}_2$

Ch 10 More Departures From Assumptions

residual $e_i = y_i - \hat{y}_i = y_i - \underline{x}_i^T \hat{\underline{\beta}} =$

$$\underline{x}_i^T \underline{\beta} + \varepsilon_i - \underline{x}_i^T \hat{\underline{\beta}} = \varepsilon_i + \underline{x}_i^T (\underline{\beta} - \hat{\underline{\beta}}) = \varepsilon_i + N_i$$

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

So $\hat{\underline{\beta}} - \underline{\beta} \sim AN_p(\underline{0}, \sigma^2 \frac{W}{n}) \sim AN_p(\underline{0}, \text{MSE}(\underline{X}'\underline{X})^{-1})$

So $\underline{\beta} - \hat{\underline{\beta}} \sim AN_p(\underline{0}, \text{MSE}(\underline{X}'\underline{X})^{-1})$

So $N_i = \underline{x}_i^T (\underline{\beta} - \hat{\underline{\beta}}) \sim AN_1(0, \text{MSE} \underline{x}_i^T (\underline{X}'\underline{X})^{-1} \underline{x}_i)$
 $= AN_1(0, \text{MSE} h_i)$

$\frac{\underline{X}'\underline{X}}{n} \rightarrow W^{-1}, n(\underline{X}'\underline{X})^{-1} \underline{x}_i \underline{x}_i^T \rightarrow W$ so $\frac{W}{n} \approx (\underline{X}'\underline{X})^{-1}$

2) So the residual $e_i = \varepsilon_i + N_i$ where

$$N_i \sim AN(0, \text{MSE} \cdot h_i) \quad \text{and} \quad h_i \rightarrow 0$$

as $n \rightarrow \infty$. But for moderate n ,

the N_i term could dominate, and

the e_i tend to be more normal

than the ε_i if the ε_i are not

$$i.i.d. N(0, \sigma^2).$$

3) The residual plot \hat{y} vs. e should always be used. The plot

of e vs. y is not used

$$\text{since } \text{corr}(\underline{e}, \underline{y}) = \sqrt{1 - R^2}.$$

want plots where $\text{corr}(\underline{y}, \underline{e}) = 0$.

If the predictors are no good, $R^2 \approx 0$,

$$\hat{y}_i \approx \bar{y}, \quad e_i \approx y_i - \underbrace{\bar{y}}_{\text{constant}} \quad \text{and} \quad \text{corr}(\underline{e}, \underline{y}) \approx 1.$$

Proof

LM 56

$$\text{Proof } \text{corr}(\underline{e}, \underline{y}) = 1 - R^2 !$$

$$\text{corr}(\underline{e}, \underline{y}) = \frac{\sum (e_i - \bar{e})(y_i - \bar{y})}{\sqrt{\sum e_i^2 \sum (y_i - \bar{y})^2}}$$

$$\text{But } \bar{e} = 0 \text{ so top} = \sum e_i (y_i - \bar{y}) =$$

$$\sum e_i y_i + \bar{y} \sum e_i = \underline{e}^T \underline{y}$$

$$\text{Now } \sum e_i^2 = \underline{e}^T \underline{e} = \underline{y}^T (\underline{I} - \underline{P}) \underline{y} = \underline{y}^T \underline{e} = \underline{e}^T \underline{y}$$

$$\text{So } \text{corr}(\underline{e}, \underline{y}) = \frac{\sum e_i^2}{\sqrt{\sum e_i^2 \sum (y_i - \bar{y})^2}} = \sqrt{\frac{\sum e_i^2}{\sum (y_i - \bar{y})^2}}$$

$$= \sqrt{\frac{\text{SSE}}{\text{SSTO}}} \quad \therefore R^2 = 1 - \frac{\text{SSE}}{\text{SSTO}}$$

$$\text{So } R^2 = 1 - \sqrt{\text{corr}(\underline{e}, \underline{y})} \quad \text{or}$$

$$\text{corr}(\underline{e}, \underline{y}) = \sqrt{1 - R^2}$$

4) p266 → Use $R = H$ in this chapter.

$$E(e) = (I - H)Y, \quad \hat{Y} = HY$$

$$E(e) = 0$$

$$\text{cov}(e) = \sigma^2(I - H)$$

$$\text{cov}(e, \hat{Y}) = 0 \text{ scalar}$$

$$\text{cov}(\hat{Y}) = \sigma^2 H$$

5) Let $h_i = H_{ii}$ be the i th leverage

The internally studentized residual

$$r_i = \frac{e_i}{\sqrt{\text{MSE}} \sqrt{h_i}} \quad \text{handwritten}$$

externally studentized residual

$$t_i = \frac{e_i}{\sqrt{\text{MSE}(i)} \sqrt{1 - h_i}} \quad \text{have } U(r_i | \alpha), V(t_i | \alpha).$$

Here $\text{MSE}(i)$ is the MSE from the regression model without the i th case.

6] P 269

$$\sum_{i=1}^n h_i = P = \text{tr}(H)$$

$\frac{1}{n} \leq h_i \leq 1$ if X has a column of 1's.

The average h_i value is $\frac{P}{n}$.

check $E(Y|X) = \beta'X$ with

response and residual plots.

Ship § 10.3.1

8]

P 276

power transformation

$$f_\lambda(x) = \begin{cases} \sqrt[\lambda]{x} & \lambda \neq 0 \\ \log(x) & \lambda = 0 \end{cases}$$

are useful for removing strong nonlinearities from the predictors.

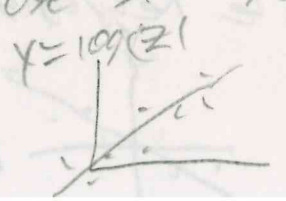
9]

P 297

If Z is the "response"

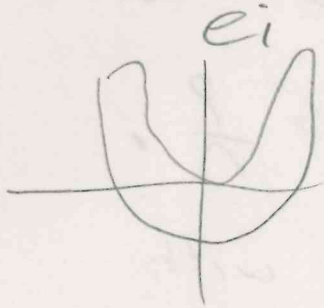
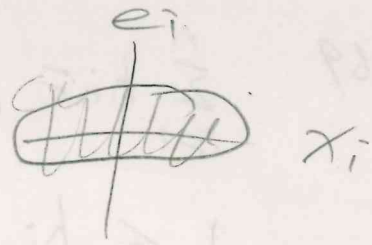
Plot $Y = f_\lambda(Z)$ vs \hat{Y}_λ for $\lambda = -1, -\frac{1}{2}, -\frac{1}{3}, 0, \frac{1}{3}, \frac{1}{2}$

Choose λ so response plot is linear.



log transformation is good

10) plot x_i vs e_i



x_i \leftarrow add x_i^2 to model

113 P276 In a (generalized) additive model

$$(GAM), y_i = \alpha + \sum_{j=1}^{p-1} S_j(x_{ij}) + \epsilon_i$$

$$= \alpha + \epsilon_i$$

where $S_j(x_{ij}) = B_j x_{ij}$ or S_j is

unspecified. If S_j is unspecified

the software will show a plot of \hat{S}_j .

If \hat{S}_j is linear, $B_j x_{ij}$ may be good.

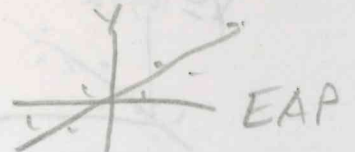
If \hat{S}_j is a quadratic, add x_{ij}^2 to

the OLS model. (A GAM can also be

used when you can't find an OLS model with good response and residual plots.

Let $EAP_i = \hat{\alpha} + \sum_{j=1}^{p-1} \hat{S}_j(x_{ij})$. The response

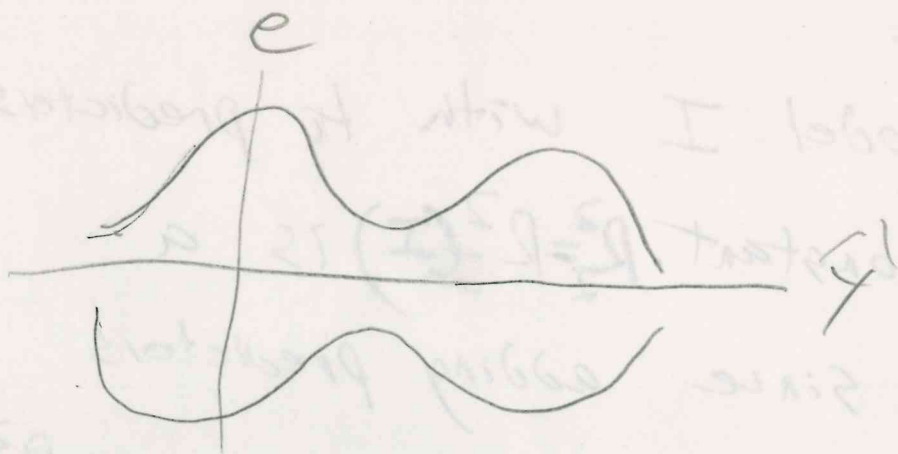
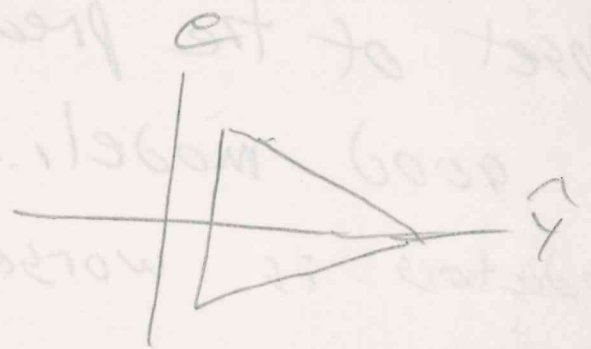
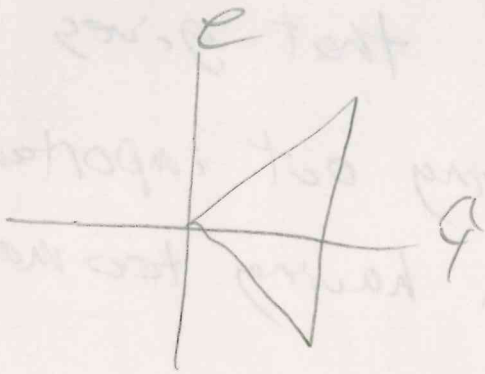
plot of EAP vs y should be linear for a good GAM.



Step 10.33

LM 58

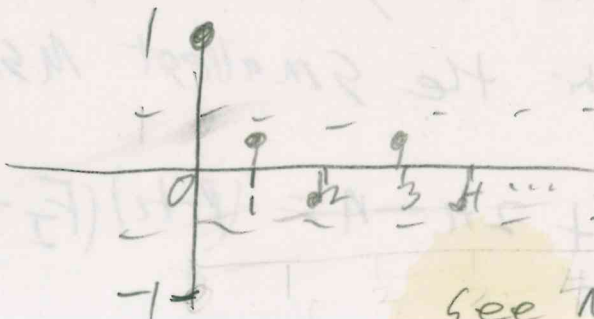
Step 10.4 12) Detect non constant variance
from residual plots



13) WLS, GLS are remedies

14) e_i could be correlated

AIC and PACF of e_i are useful



want plotted points between
SE bars except at

lag = 0

see Math 474 Time series

Step § 10.5 - 10.7

Ch 12 } Variable selection or subset selection tries to find a subset of the predictors that gives a good model. Leaving out important predictors is worse than having too many predictors.

2) ^{p400-1} Consider model I with k predictors including a constant. $R^2 = R^2(I)$ is a poor measure since adding predictors does not decrease and usually increases R^2 .

$MSE(I)$ and adjusted R^2

not everyone uses this formula

$$\bar{R}^2(I) = 1 - \frac{(1 - R^2(I))}{n - k} = 1 - \frac{MSE(I)}{SSTO}$$

are better measures.

The model with the largest $\bar{R}^2(I)$ is also the model with the smallest $MSE(I)$.

3) $C_p(I) = \frac{SSE(I)}{MSE(F)} + 2k - n = (p - k)(F_I - 1) + k$

where F_I is the partial F test statistic when the reduced model is I and $MSE(F)$ is the MSE of the full model that uses all p predictors, including a constant. LM 59

$$4) \text{corr} \left(\underset{\substack{\uparrow \\ \text{Full}}}{\sqrt{F}}, \tilde{e}_I \right) = \sqrt{\frac{SSE(F)}{SSE(I)}} = \sqrt{\frac{n-p}{C_p(I) + n - 2k}}$$

5) Models with $C_p(I) \leq \min(2k, p)$ are of interest.

6) Let I_{\min} be the model (eg from forward selection) with the smallest $C_p(I)$.

Find the submodel I_I with the fewest number of predictors such that

$C_p(I_I) \leq C_p(I_{\min}) + 1$. Then model I_I is the initial model to be examined, it is possible that $I_I = I_{\min}$ or $I_I = I_{(p)}$ the full model.

7) a) Models I with fewer predictors than I_{min} such that $C_p(I) \leq C_p(I_{min}) + 4$ should also be examined.

b) Suppose model I has k predictors including a constant. Models I with fewer predictors than I_{min} such that

$C_p(I_{min}) + 4 < C_p(I) \leq \min(2k, p)$ should be checked but often underfit; important predictors are omitted.

8) ^{PHH-418} Forward selection and backward elimination; see E3 rev 113).

9) C_p needs a full model, with $n \geq 10p$ or so. Other criteria such as $\bar{R}^2(I)$ ($MSE(I)$) can be used.

If there are J potential predictors

where $J > n$ (or $J > \frac{n}{3}$), use

forward selection with $\bar{R}^2(I)$ until there are $\approx \frac{n}{3}$ (or $\frac{n}{10}$) predictors in the model.

10) The full model should be good,
 Make the usual checks on the
 selected submodel. The response
 and residual plots of the submodel
 should be similar to those of the
 full model.

ex) backward elimination
 current terms x1 x2 x3 x4 x5

	k	Cp
delete x1	5	3.132
x3	5	3.143
x4	5	4.378
x2	5	6.698
x5	5	297.743

Current terms x2 x3 x4 x5

	k	Cp
delete x3	4	2.027
x4	4	2.536
x2	4	4.834
x5	4	306.913

Current terms x2 x4 x5

	k	Cp
delete x4	3	1.617
x2	3	4.258
x5	3	305.047

← $I_{min} = I_1$

Current terms x2 x5

	k	Cp
delete x2	2	4.456
x5	2	337.149

← Fewer predictors than I_{min}

what are the terms, including an intercept,
in II ?

Soln intercept, x_2 , x_5
 $k=3$

since x_4 was deleted

ex) forward selection
if model current terms

	x_3, x_4, x_5		
add	x_6	k	CP
		5	6.357 ← best
	x_7	5	16.748
	x_1	5	17.891
	x_2	5	19.763

"best model" has

intercept, x_3, x_4, x_5 , and x_6 ,
 $k=5$

since x_6 was added to the model

see HW 9 (5), (6).