

LM 61
11) Variable selection is a form of data snooping, so "inference is not valid after selecting a model using variable selection." Can use the model for description.

If possible do a pilot study with variable selection, then do a bigger study using the model selected in the pilot study as the candidate full model. Inference from the bigger study can be done.

12) Often building a good full model also uses data snooping, then inference is no longer valid.

13) ^{p. 312-3} For a given data set, may have to use data snooping to find a model that fits the data reasonably well. The fitted model tends to fit the training data $(Y_1, X_1), \dots, (Y_n, X_n)$ better than future test data, and "inference can not be justified by large sample theory."

14) ^(don't worry about several candidate models) If the variable selection is done by computer, eg use I_1 or I_{min} and data snooping was not used to build the full model, it is possible to develop valid inference for the selected model. It is known that I_{min} overfits, asymptotically, see p 448.

Skim § 12.2

Ignore cross validation, on p 403-6) § 12.3.3-4

Skim § 12.4.3 on stepwise selection?
Terms that were deleted can come back in the model.

§ 12.5 15) Suppose $\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$ and $p > \frac{n}{5}$ or there is collinearity, Shrinkage methods shrink some $|\hat{\beta}_i|$ values towards 0.

ex) Forward selection and backwards elimination make $\hat{\beta}_i = 0$ for the omitted variables.

16) PH23

The ridge regression estimator

$$\text{is } \hat{\beta}(\lambda) = \underbrace{(X'X + \lambda I)}_{p \times p}^{-1} X'Y, \lambda \geq 0$$

non-singular even if $X'X$ is singular for $\lambda > 0$

$\hat{\beta}(\lambda)$ minimizes the criterion

$$RSS(\alpha) = (Y - X\beta)^T (Y - X\beta) + \lambda \beta_s^T \beta_s$$

$$= RSS(\beta) + \lambda \beta_s^T \beta_s = RSS(\beta) + \lambda \sum_{i=1}^{p-1} \beta_i^2$$

where $\beta = \begin{pmatrix} \beta_0 \\ \beta_s \end{pmatrix}$. This estimator can

be fit even if $p \geq n$. Typically use

$\hat{\beta}(\hat{\lambda})$ where λ is estimated by $\hat{\lambda}$,
(GCV CV etc)

(Let $w = X'X$ then $X'X + \lambda I = w + \lambda I = \frac{1}{n} \sum_{i=1}^n w_i^2 \geq 0$ so $X'X \geq 0$)

so $X'(X + \lambda I)X = X'X + \lambda X'X \geq 0$ for $\lambda > 0$ unless $X = 0$.

Step §12.5.3, 12.6

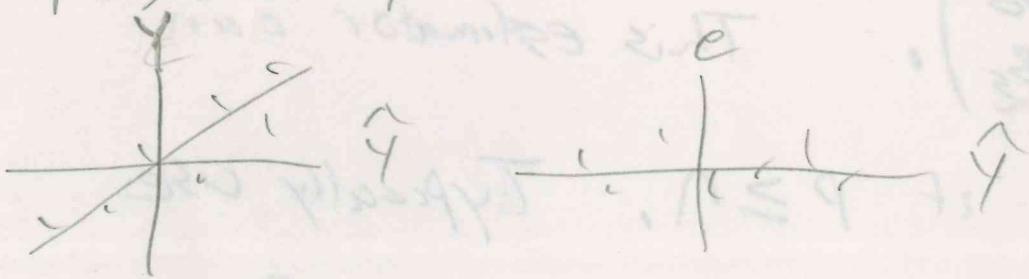
$\sum_{i=1}^p w_i^2 > 0$ unless $X = 0$. so $X'X + \lambda I > 0$ for $\lambda > 0$

Step §12.7,

Step §12.8, 12.9 Step §12.9 except Step §12.9.1

ch 8 } p187 In experimental design models or design of experiments 'DOE', the entries of \underline{X} are coded, often $-1, 0, 1$. Often the \underline{X} matrix is not a full rank matrix.

2) Some DOE models, like the 2^k factorial model, have one Y_i per \underline{x}_i , lots of \underline{x}_i 's, and the response and residual plots look like those for MLR.



3) Some DOE models have n_i Y_i 's for each distinct \underline{x}_i . Then the response and residual plots no longer look like those for MLR.

4) Suppose there are p distinct \underline{x}_i , called treatments, $n = n_1 + \dots + n_p$, and $n_i \equiv m = \frac{n}{p}$. A dot plot of z_1, \dots, z_m consists of

an axis and m points corresponding to z_1, \dots, z_m .

LM 63

If $m \geq 5$ and p is small,

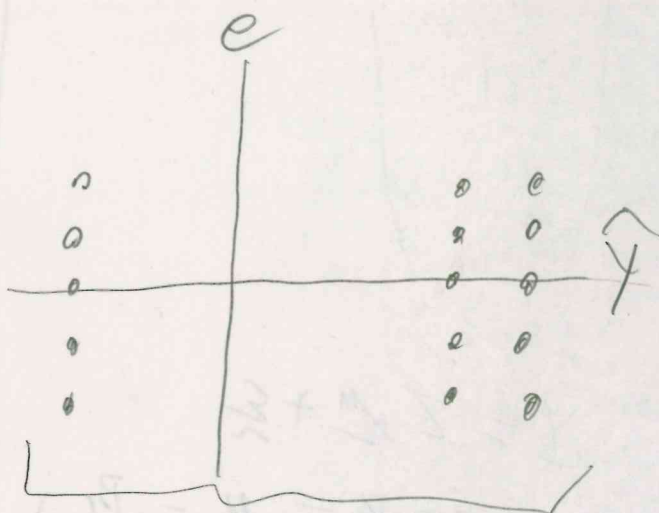
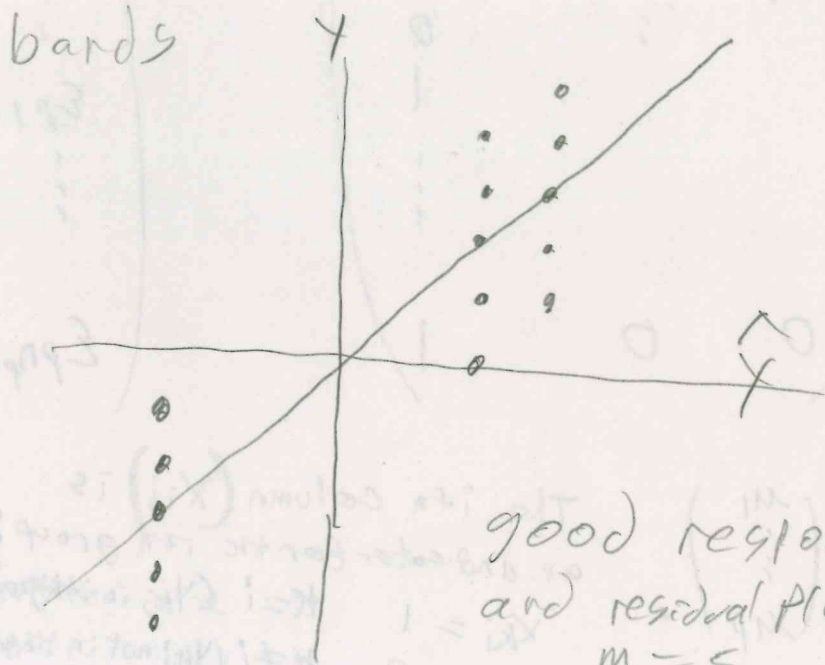
the response plot consists of

p dot plots, one for each treatment.

If $m \geq 10$, the p dot plots should

have roughly the same shape and spread.

The residual plot also consists of p dot plots. The points in the response and residual plots should scatter about the identity and $e=0$ lines, but the scatter need not be in evenly populated bands



good response and residual plots with $m=5$. 3 dot plots

6) P189 For the cell means model,

LM 64

X is full rank, $\underline{1}$ is not a column of X , but $\underline{1} \in C(X)$

Since if $X = (\underline{v}_1, \dots, \underline{v}_p)$, then $\underline{1} = \sum_{i=1}^p \underline{v}_i$.

$$X'X = \begin{pmatrix} 1 & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 & 0 & 1 & \dots & 0 \\ & & & & & & & 1 \end{pmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}$$

$$= \begin{pmatrix} n_1 & 0 & \dots & 0 \\ 0 & n_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & n_p \end{pmatrix}$$

$$(X'X)^{-1} = \text{diag}\left(\frac{1}{n_1}, \dots, \frac{1}{n_p}\right)$$

$$X'\underline{y} = \begin{pmatrix} \sum_{j=1}^{n_1} y_{1j} & \sum_{j=1}^{n_2} y_{2j} & \dots & \sum_{j=1}^{n_p} y_{pj} \end{pmatrix} = (y_{10} \ y_{20} \ \dots \ y_{p0})^T$$

$$\text{So } \hat{\underline{\mu}} = \hat{\underline{\beta}} = (X'X)^{-1} X'\underline{y} = \text{diag}\left(\frac{1}{n_1}, \dots, \frac{1}{n_p}\right) \begin{pmatrix} y_{10} \\ \vdots \\ y_{p0} \end{pmatrix} = \begin{pmatrix} \bar{y}_{10} \\ \vdots \\ \bar{y}_{p0} \end{pmatrix}$$

$$\hat{y} = X(X'X)^{-1}X'y = X\hat{\mu} =$$

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & & & \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & & & \\ 0 & 1 & \dots & 0 \\ \vdots & & & \\ 0 & \dots & 0 & 1 \\ \vdots & & & \\ 0 & \dots & 0 & 1 \end{pmatrix} \begin{pmatrix} \bar{y}_{10} \\ \bar{y}_{20} \\ \vdots \\ \bar{y}_{p0} \end{pmatrix} = \begin{pmatrix} \bar{y}_{10} \\ \vdots \\ \bar{y}_{10} \\ \bar{y}_{20} \\ \vdots \\ \bar{y}_{20} \\ \vdots \\ \bar{y}_{p0} \\ \vdots \\ \bar{y}_{p0} \end{pmatrix}$$

So $\hat{y}_{ij} = \bar{y}_{i0}$, $i=1, \dots, p$, $j=1, \dots, n_i$.

→ Hence the dot plot for the j th treatment crosses the identity line at \bar{y}_{j0} in the response plot.

