6) p189 For the cell means model,

$\underset{\sim}{X}$ is full rank, $\underset{\sim}{1}$ is not a

column of $\underline{X}$, but $\underset{\sim}{1} \in C(\underline{X})$

Since if $\underline{X} = (\underset{\sim}{v_1}, \ldots, \underset{\sim}{v_p})$, then $\underset{\sim}{1} = \sum_{j=1}^{p} \underset{\sim}{v_j}$.

$$\underset{p \times n}{\underline{X}'} \underset{n \times p}{\underline{X}} = \begin{pmatrix} 1 & \cdots & 1 & 0 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 1 & \cdots & 1 & \vdots & \vdots & & \vdots \\ \vdots & & \vdots & & & & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & 0 & 0 & 0 & & 1 & \cdots & 1 \end{pmatrix} \begin{bmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & & 0 \\ 0 & 1 & & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 1 & & 0 \\ 0 & 0 & & 1 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & & 1 \end{bmatrix}$$

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

$$(\underline{X}'\underline{X})^{-1} = \operatorname{diag}\left(\frac{1}{n_1}, \ldots, \frac{1}{n_p}\right).$$

$$\underline{X}'\underset{\sim}{Y} = \left( \sum_{j=1}^{n_1} y_{1j} \quad \sum_{j=1}^{n_2} y_{2j} \quad \cdots \quad \sum_{j=1}^{n_p} y_{pj} \right)' = \left( y_{1 \bullet} \quad y_{2 \bullet} \cdots y_{p \bullet} \right)^T$$

So $\underset{\sim}{\hat{\mu}} = \underset{\sim}{\hat{\beta}} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underset{\sim}{Y} = \operatorname{diag}\left(\frac{1}{n_1}, \ldots, \frac{1}{n_p}\right) \begin{pmatrix} y_{1\bullet} \\ \vdots \\ y_{p\bullet} \end{pmatrix} = \begin{pmatrix} \bar{y}_{1\bullet} \\ \vdots \\ \bar{y}_{p\bullet} \end{pmatrix}$

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

$$\begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & & & \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & & & \\ 0 & 1 & \cdots & 0 \\ \vdots & & & \\ 0 & \cdots & 0 & 1 \\ \vdots & & & \\ 0 & \cdots & 0 & 1 \end{pmatrix} \begin{pmatrix} \bar{Y}_{10} \\ \bar{Y}_{20} \\ \vdots \\ \bar{Y}_{P0} \end{pmatrix} = \begin{pmatrix} \bar{Y}_{10} \\ \vdots \\ \bar{Y}_{10} \\ \bar{Y}_{20} \\ \vdots \\ \bar{Y}_{20} \\ \vdots \\ \bar{Y}_{P0} \\ \vdots \\ \bar{Y}_{P0} \end{pmatrix}$$

So $\hat{Y}_{ij} = \bar{Y}_{i0}$, $i=1,...,P$, $j=1,...,n_i$.

7) Hence the dot plot for the $i$th treatment crosses the identity line at $\bar{Y}_{i0}$ in the response plot.

8} One way Anova F test

$H_0$ $\mu_1 = \cdots = \mu_p$    $H_A$ not $H_0$
(not all of the $p$ means are equal)

9} For $p = 2$ this is the pooled $t$ test

$H_0$ $\mu_1 = \mu_2$    $H_A$ $\mu_1 \neq \mu_2$.

P189

10} If $H_0$ is true, let $\mu_1 = \cdots = \mu_p \equiv \mu$

Then the $y_{ij}$ are iid,

$y_{ij} = \mu + \varepsilon_{ij}$    and    $\hat{\mu} = \overline{Y_{00}} = \dfrac{\sum\sum y_{ij}}{n}$.

$\varepsilon_{ij}$ iid $N(0, \sigma^2)$

From P100, $F = \dfrac{(RSS(H) - RSS)/q}{RSS/(n-p)} \sim F_{q, n-p}$

$q = p-1$

To show $q = p-1$, need full rank

Need $A$ $\ni$ $A\underline{\mu} = \underline{0}$ is equivalent to $H_0$.

$q \times p$    $q \times p$

$A = \begin{pmatrix} 1 & -1 & 0 & 0 & 0 \\ 1 & 0 & -1 & 0 & 0 \\ \vdots & & & -1 & 0 \\ 1 & 0 & \cdots & 0 & -1 \end{pmatrix}$ works since $A\begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix} = \begin{pmatrix} \mu_1 - \mu_2 \\ \vdots \\ \mu_1 - \mu_p \end{pmatrix}$.

$(p-1) \times p$

$$RSS = \sum_i \sum_j (Y_{ij} - \overline{Y}_{io})^2, \quad RSS(H) = \sum_i \sum_j (Y_{ij} - \overline{Y}_{oo})^2,$$
$$= SSTO$$

11) p192 total sum of squares $\quad SSTO = \sum_{i=1}^{p} \sum_{v=1}^{n_i} (Y_{ij} - \overline{Y}_{oo})^2$

treatment sum of squares $\quad SSTR = \sum_{i=1}^{p} n_i (\overline{Y}_{io} - \overline{Y}_{oo})^2$

residual or error sum of squares $SSE = \sum_i \sum_j (Y_{ij} - \overline{Y}_{io})^2$

$SSTO = SSTR + SSE, \quad RSS(H) - RSS = SSTO - SSE = SSTR,$

one way ANOVA table $\quad MS = SS/df$

| Source | df | SS | MS | F |
|---|---|---|---|---|
| between or treatment | p-1 | SSTR | MSTR | $\dfrac{MSTR}{MSE}$ |
| error | n-p | SSE | MSE | |

$H_0: \mu_1 = \cdots = \mu_p \quad H_A: \text{not } H_0$

$$Pval = P(F_{p-1, n-p} > F).$$

reject $H_0$ if $pval \le \delta$

fail to reject $H_0$ if $pval \ge \delta$

12) Rule of thumb: Let $R_1, \ldots, R_p$ be the ranges of the $p$ dot plots. If

$max(R_1, \ldots, R_p) \le 2 \, min(R_1, \ldots, R_p)$ then the one way Anova F test $pval \approx$ correct.

if the response and residual plots suggest that the remaining one way Anova model assumptions are reasonable.

So the test has some robustness to the assumption $V(\varepsilon_{ij}) \equiv \sigma^2$.

Note: P195 CIs are less robust to this assumption,

13} $Y_{ij} = \mu_i + \varepsilon_{ij}$, $\varepsilon_{ij}$ iid, $E(\varepsilon_{ij})=0$, $V(\varepsilon_{ij})=\sigma^2$
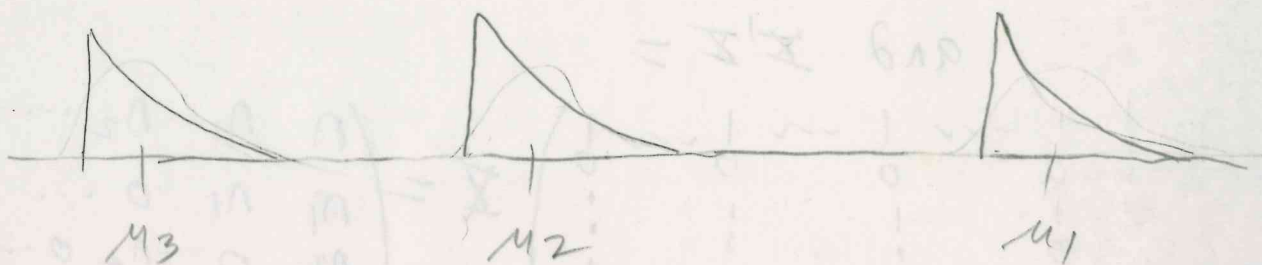
implies that the $\varepsilon_{ij}$ have a pdf $f_\varepsilon(z)$. So the $Y_{ij}$ $(j=1,\ldots,n_i)$ have pdf

$f_\varepsilon(z-\mu_i)$ for $i=1,\ldots,P$ with $\sigma_i^2 \equiv \sigma^2$

**strong assumption**

location family, same shape, but different means $\mu_i = E(Y_{ij})$.



$\mu_3$  $\mu_2$  $\mu_1$

The $V(\varepsilon_{ij}) \equiv \sigma^2$ assumption is much stronger for the One way Anova model than for the MLR model. The F test is a large sample test if $V(\varepsilon_{ij}) \equiv \sigma^2$

14) Another $X$ matrix for the one way Anova model adds a constant and deletes the last column of the $X_c$ for the cell means model.

$$Y = X\underline{B} + \underline{\varepsilon}, \qquad \underline{B} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{P-1} \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & 0 & & 0 \\ 1 & 1 & 0 & & 0 \\ 1 & 0 & 1 & & 0 \\ \vdots & 0 & 1 & & 0 \\ 1 & 0 & 1 & & 0 \\ & & & & \vdots \\ & & & & 0 \\ 1 & 0 & 0 & & 0 \end{pmatrix}, \qquad \text{so } X'Y = \begin{pmatrix} Y_{00} \\ Y_{10} \\ \vdots \\ Y_{P-1,0} \end{pmatrix}$$

and $X'X =$

$$\begin{pmatrix} 1 & \cdots & 1 & 1 & \cdots & 1 & \cdots & 1 & \cdots & 1 \\ 1 & \cdots & 1 & 0 & \cdots & 0 & & & & 0 \\ 0 & \cdots & 0 & 1 & \cdots & 1 & & & & 0 \\ \vdots & & \vdots & \vdots & & & & \ddots & & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 & 1 & \cdots & 1 & 0 \end{pmatrix} \qquad X = \begin{pmatrix} n & n_1 & n_2 & \cdots & n_{P-1} \\ n_1 & n_1 & 0 & \cdots & 0 \\ n_2 & 0 & n_2 & 0 \cdots & 0 \\ \vdots & & & & \\ n_{P-1} & 0 & \cdots & 0 & n_{P-1} \end{pmatrix}$$

$$X'X = \begin{pmatrix} n & (n_1 \ n_2 \ \cdots \ n_{p-1}) \\[2em] \begin{pmatrix} n_1 \\ n_2 \\ \vdots \\ n_{p-1} \end{pmatrix} & diag(n_1, \cdots, n_{p-1}) \end{pmatrix}$$

$$(X'X)^{-1} = \frac{1}{np}\begin{bmatrix} 1 & -1 & \cdots & & -1 \\ -1 & 1+\frac{np}{n_1} & 1 & \cdots & 1 \\ \vdots & 1 & 1+\frac{np}{n_2} & \cdots & 1 \\ & \vdots & & \ddots & \vdots \\ -1 & 1 & \cdots & 1 & 1+\frac{np}{n_{p-1}} \end{bmatrix}$$

$$= \frac{1}{np}\begin{pmatrix} 1 & -\underline{1}' \\[1em] -\underline{1} & \underline{1}\,\underline{1}' + diag\left(1+\frac{np}{n_1}, \cdots, 1+\frac{np}{n_{p-1}}\right) \end{pmatrix}.$$

<span style="color:red">typo in 134)</span>

<span style="color:red">$\frac{np}{n_i}$ not $1+\frac{np}{n_i}$</span>

$$\hat{\underline{\beta}} = (X'X)^{-1}X'\underline{Y} = \begin{pmatrix} \bar{Y}_{po} \\ \bar{Y}_{1o} - \bar{Y}_{po} \\ \vdots \\ \bar{Y}_{p-1,o} - \bar{Y}_{po} \end{pmatrix}, \qquad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix} = \begin{pmatrix} \mu_p \\ \mu_1 - \mu_p \\ \vdots \\ \mu_{p-1} - \mu_p \end{pmatrix}$$

15) This model is interesting since the
one-way Anova F test
$H_0$ $\mu_1 = \cdots = \mu_p$ vs $H_A$ not $H_0$

corresponds to the MLR Anova F test

$H_0$ $\beta_1 = \cdots = \beta_{p-1} = 0$ vs $H_A$ not $H_0$.

16) p192 A contrast $\theta = \sum_{i=1}^{p} c_i \mu_i = \underline{c}' \underline{\mu}$.
where $\sum_{i=1}^{p} c_i = 0$.
A $100(1-\delta)\%$ CI for $\theta$ is

$$\sum_{i=1}^{p} c_i \bar{Y}_{i0} \pm t_{n-1, 1-\frac{\delta}{2}} \sqrt{MSE} \sqrt{\sum_{i=1}^{p} \frac{c_i^2}{n_i}} .$$

Assuming groups are independent,

$$V\left(\sum_{i=1}^{p} c_i \bar{Y}_{i0}\right) = \sum_{i=1}^{n} c_i^2 V(\bar{Y}_{i0}) = \sum_{i=1}^{n} c_i^2 \frac{\sigma^2}{n_i} .$$

so $SE\left(\sum_{i=1}^{p} c_i \bar{Y}_{i0}\right) = \sqrt{MSE} \sqrt{\sum_{i=1}^{p} \frac{c_i^2}{n_i}} .$

p193
17) can use Bonferroni CIs and Scheffe CIs
for $k$ contrasts where want $P($ all $k$ CIs
contain $\theta_j$, $j = 1, \ldots, k) \geq 1 - \delta$.

Inference After Variable Selection and Lasso,
see ch4 of online notes

1) One simple method for inference
after variable selection is <u>data splitting</u>
with 2 sets: let the training set have
$n_T \leq \frac{n}{2}$ cases and the validation set
have $n_V = n - n_T \geq \frac{n}{2}$ cases. Select the $n_T$
cases without replacement from the $n$ cases.
Assume the cases are independent and
follow a statistical model, eg MLR,

I) Build the model I with the training set,
possibly using variable selection and
using the response to select predictors
and predictor transformation.
Let model I have $t_I$ predictors.

II) Act as if I is the full model
for the validation set. want
$n \geq 5k$ and preferably $n \geq 10t_i$.
Need model I to be a good model for the data.

*2) Variant: use, for example, $\frac{n}{10}$ cases
for the training set. If you can not
get a good model, select $\frac{n}{10}$ cases from the

validation set for the new training set. Cases that remain are the new validation set.

Repeat until $I$ is a good model or $n_T \leq \frac{n}{2}$ with $n_T \approx \frac{n}{2}$.

$$\frac{n}{10} \quad \frac{2n}{10} \quad \frac{3n}{10} \quad \frac{4n}{10} \quad \frac{5n}{10} \qquad \text{training}$$

$$\frac{9n}{10} \quad \frac{8n}{10} \quad \frac{7n}{10} \quad \frac{6n}{10} \quad \frac{5n}{10} \qquad \text{validation,}$$

3) Data splitting can work for $n \geq 10p$ and $n << 10p$. Inefficient inference is much better than invalid inference. Efficiency $\approx \frac{n_V}{n} = 1 - \frac{n_T}{n}$.

4) The bootstrap is useful if $n \geq 10p$. The bootstrap is used for tests, CIs and confidence regions.

5) Suppose $z_1, \ldots, z_n$ are iid and there is a statistic $T = T(z_1, \ldots, z_n)$,

$$\underset{p \times m}{T} \quad \underset{p \times 1}{} \quad \underset{1 \times 1}{}$$

Suppose we could gather $B$ iid samples

$$T^{(1)} = T\left(\underline{z}_i^{(1)}, \ldots, \underline{z}_n^{(1)}\right)$$
$$\vdots$$

$$T^{(B)} = T\left(\underline{z}_i^{(B)}, \ldots, \underline{z}_n^{(B)}\right), \qquad \text{If } B \text{ is large}$$

could examine $T^{(1)}, \ldots, T^{(B)}$ for inference.

6) The $\underline{\text{empirical distribution}}$ gives probability $\frac{1}{n}$

to the iid data $\quad \underline{z}_1, \ldots, \quad \underline{z}_n$
$$\frac{1}{n} \qquad \qquad \frac{1}{n}$$

$\left( \text{So } E\underline{w} = \sum_{i=1}^{n} \underline{z}_i \frac{1}{n} = \bar{\underline{z}} \text{ and } E\underline{w}\underline{w}^T = \sum_{i=1}^{n} \underline{z}_i \underline{z}_i^T \frac{1}{n} \right)$, $\text{cov } \underline{w} = \frac{1}{n} \sum (\underline{z}_i - \bar{\underline{z}})(\underline{z}_i - \bar{\underline{z}})^T$

Let $\underline{w}$ be from the emp. dist.
$$E(\underline{w} - E\underline{w})(\underline{w} - E\underline{w})^T =$$
$$E\underline{w}\underline{w}^T - E(\underline{w})(E\underline{w})^T$$

So draw $\underline{w}$ from the empirical distribution

and $P(\underline{w} = \underline{z}_i) = \frac{1}{n} \qquad i = 1, \ldots, n.$

Sample with replacement to get

$\underline{w}_1, \ldots, \underline{w}_n \qquad$ which are $n$ iid observations

from the empirical distribution.

7) Suppose iid sample is $Y_1, \ldots, Y_n$.

Then the empirical c.d.f. is

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^{n} \underbrace{I\{Y_i \leq x\}}_{}$$

$$\begin{cases} 1 & \text{if } Y_i \leq x \\ 0 & \text{if } Y_i > x \end{cases}$$

So $I\{Y_i \leq x\}$ are iid binomial $(1, p)$

where $p = P(Y_i \leq x) = F_Y(x)$.

So $E[I\{Y_i \leq x\}] = F_Y(x)$ and $V[I\{Y_i \leq x\}] = F_Y(x)(1 - F_Y(x))$,

For fixed $x$, $\hat{F}(x)$ is a sample mean, so

by the CLT $\sqrt{n}\left(\hat{F}(x) - F(x)\right) \xrightarrow{D} N\left(0, \underbrace{F_Y(x)(1 - F_Y(x))}_{\in [0, \frac{1}{4}]}\right)$

So $\hat{F}(x)$ is a $\sqrt{n}$ consistent estimator of $F_Y(x)$.

8) For many statistics $T(\underline{Z}_1, \dots, \underline{Z}_n) - T(\underline{w}_1, \dots \underline{w}_n) \xrightarrow{D} 0$

where the $\underline{w}_i$ are iid drawn with replacement from
the empirical distribution of $\underline{z}_1, \dots, \underline{z}_n$,

9) Suppose $T$ is a vector (or $T$ is $p \times m$ and look at $T_{ij}$). Let

$$T = T(\underline{Z}_1, \dots, \underline{Z}_n) \quad \text{where } \underline{Z}_1, \dots, \underline{Z}_n \text{ are iid.}$$

Let $\underline{Z}_{11}^*, \underline{Z}_{12}^*, \dots, \underline{Z}_{1n}^*, \quad T_1^* = T(\underline{Z}_{11}^*, \dots, \underline{Z}_{1n}^*)$

$\underline{Z}_{21}^*, \underline{Z}_{22}^*, \dots, \underline{Z}_{2n}^*, \quad T_2^* = T(\underline{Z}_{21}^*, \dots, \underline{Z}_{2n}^*)$

$\vdots$

$\underline{Z}_{B1}^* \ \underline{Z}_{B2}^* \dots \ \underline{Z}_{Bn}^*, \quad T_B^* = T(\underline{Z}_{B1}^*, \dots, \underline{Z}_{Bn}^*)$

$B$th bootstrap sample

ex) data $Z_1, \dots, Z_7 = 1, 2, 3, 4, 5, 6, 7$

$$T = \text{median}(Z_1, \dots, Z_7) = 4$$

Let $B = 2$

$(2, 2, 2, 3, 3, 5, 6 \leftarrow \text{ordered})$

1st: $3, 2, 3, 2, 5, 2, 6 \quad T_1^* = 3$

2nd: $3, 5, 3, 4, 3, 5, 7 \quad T_2^* = 4$

$(3 \ 3 \ 3 \ 4 \ 5 \ 5 \ 7 \quad \leftarrow \text{ordered})$

10) Let $T_{1n}^*, \dots, T_{Bn}^*$ be iid with same dist as statistic $T_n$. $\sqrt{n}(T_1^* - T_n), \dots, \sqrt{n}(T_B^* - T_n)$ are often pseudodata for $\sqrt{n}(T_{1n} - \theta), \dots, \sqrt{n}(T_{Bn} - \theta)$

$\bar{T}^* = \frac{1}{B} \sum_{i=1}^{B} T_i^* \qquad S_{T^*}^2 = \frac{1}{B-1} \sum (T_i^* - \bar{T}^*)(T_i^* - \bar{T}^*)'$

11] If $P=1$) $T_{(1)}^*, T_{(2)}^*, \ldots, T_{(B)}^{(*)}$

a) $100(1-\delta)\%$ percentile CI for $E(T)=\theta$

discards smallest and largest $100\frac{\delta}{2}\%$ of $T_{(i)}^*$

or b) compute the shorth $(c=\lceil B(1-\delta)\rceil)$ interval of the $T_{(i)}^*$

$$\left[ T_{\lceil B\frac{\delta}{2}\rceil}^*, \; T_{\lceil B(1-\frac{\delta}{2})\rceil}^* \right] \qquad eg \; 1-\delta=0.9, \; \delta=0.1, \frac{\delta}{2}=0.05$$

$$\left[ T_{\lceil B(0.05)\rceil}^*, \; T_{\lceil B \, 0.975\rceil}^* \right] \quad is \; 95\%\, CI \; for \; \theta$$

reject $H_0 \; \theta=\theta_0$ if $\theta_0$ is outside the

CI.

12] Sample indices $1, \ldots, n$ with replacement

get $i_1, \ldots, i_n$ $eg !!$

nonparametric __bootstrap__ (empirical, naive, rowwise).

ex $\quad n=6 \quad \underset{\sim}{i_1} = 3, 2, 3, 2, 5, 6$

so use $(Y_3, X_3), (Y_2, X_2) (Y_3, X_3), (Y_2, X_2) (Y_5, X_5) (Y_6, X_6)$

is the 1st bootstrap sample.

For MLR and $T_i^* = \underset{\sim}{\hat{B}}_i^*$, this should work well

if $x_i = (1, u_i^T)^T$ and the

$z_i = (Y_i, u_i^T)^T$ are iid from some distribution

with nonsingular covariance matrix $\Sigma_{z}$.
This regularity condition is <u>strong</u> eg $z_i \sim N_p(\mu, \sigma^2 I)$.

13) $Y_i = E(Y_i) + \varepsilon_i = \hat{Y}_i + r_i$

The <u>residual bootstrap</u> samples the residuals

$r_1, \ldots, r_n$ with replacement giving

$r_{11}^*, \ldots, r_{1n}^*$  $Y_{1i}^* = \hat{Y}_i + r_{1i}^*$ $i=1,\ldots,n$  $\hat{\beta}_1^* = \hat{\beta}$ from

regressing $Y_1^*$ on $X$

$\vdots$

$r_{B1}^*, \ldots, r_{Bn}^*$  $Y_{Bi}^* = \hat{Y}_i + r_{Bi}^*$ $i=1,\ldots,n$  $\hat{\beta}_B^*$ from

regressing $Y_B^*$ on $X$,

$$\overbrace{\widehat{cov}(\hat{\beta})}^{\text{bootstrap est of } cov(\hat{\beta}_{OLS})} = \frac{1}{B-1} \sum_{i=1}^{B} (\hat{\beta}_i^* - \bar{\beta}^*)(\hat{\beta}_i^* - \bar{\beta}^*)^T$$

As $B \to \infty$

$$\to \frac{\sum_{i=1}^{n} r_i^2}{n} (X'X)^{-1}. \quad \text{Usual estimator is } \underbrace{\frac{\frac{\sum_{i=1}^{n} r_i^2}{n-p}}{}}_{MSE} (X'X)^{-1}$$

14) Idea for variable selection

(regulitation methods, lasso, ridge regression

etc.) where $n \geq 10p$. Let $\underline{\beta} = (\beta_1, \ldots, \beta_p)'$.

Consider automated variable selection

eg model $J = I$, $J = I_{min}$ with forward selection,

often fit $\underset{\sim}{Y} = X_{I_{min}} \underset{\sim}{\beta}_{I_{min}} + \underset{\sim}{\varepsilon}$ with MLR
$k \times 1$

and use $COV(\hat{\beta}_{I_{min}}) = \overbrace{\left(\frac{1}{n-k} \sum_{i=1}^{n} r_{i_{I_{min}}}^2\right)}^{MSE(I_{min})} (X_{I_{min}}' X_{I_{min}})^{-1}$,

which is an incorrect estimator,

Instead get bootstrap data as in $\boxed{B}$

do variable selection to get model $I_{min,i}^*$, $i = 1,...,B$.

Then $\underset{\sim}{\hat{\beta}}_{I_{min,i}}^*$ has the $k_i$ values of $\underset{\sim}{\hat{\beta}}_{I_{min,i}}$ and $p-k_i$ 0's
$p \times 1$

eg $p = 5$ $\underset{\sim}{\hat{\beta}}_{I_{min}}^* = \begin{pmatrix} 1 \\ 3.7 \\ 4.3 \end{pmatrix}$ uses $x_1 \equiv 1$, $x_3$ and $x_5$

$\underset{\sim}{\hat{\beta}}_i^* = \underset{\sim}{\hat{\beta}}_{I_{min,i}}^* = \begin{pmatrix} 1 \\ 0 \\ 3.7 \\ 0 \\ 4.3 \end{pmatrix}$, Compute $\hat{\beta}_1^*, ..., \hat{\beta}_B^*$
$p \times 1$

$\overline{\hat{\beta}}^* = \frac{1}{B} \sum_{i=1}^{B} \hat{\beta}_i^*$ and $\hat{COV}(\underset{\sim}{\hat{\beta}}) = \frac{1}{B-1} \sum_{i=1}^{B} (\hat{\beta}_i^* - \overline{\hat{\beta}}^*)(\hat{\beta}_i^* - \overline{\hat{\beta}}^*)^T$

Display CI's