

$$36) \sqrt{n}(\hat{\beta}_{MIX} - \beta) \text{ and } \sqrt{n}(\hat{\beta}_{VS} - \beta)$$

are selecting from $\underline{U}_{In} = \sqrt{n}(\hat{\beta}_{I_{H_0}} - \beta)$ and asymptotically from \underline{U}_H (where $\underline{U}_H \sim NP(0, U_{H0})$ if $S \subseteq I_H$). Random selection does not change the dist of \underline{U}_{In} and \underline{U}_H , but variable selection changes the dist of selected \underline{U}_{In} to \underline{W}_{In} and the dist of selected \underline{U}_j to \underline{W}_j .

37) Variable selection CLT: Assume

$P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$ and let

$\hat{\beta}_{VS} = \hat{\beta}_{I_{H_0}}$ with probs π_{In} where

$\pi_{In} \rightarrow \pi_H$ as $n \rightarrow \infty$. Denote the $\pi_H > 0$ by π_j . Assume

$\underline{W}_{In} = \sqrt{n}(\hat{\beta}_{I_{H_0}} - \beta) \xrightarrow{D} \underline{W}_j$. Then

$\underline{W}_n = \sqrt{n}(\hat{\beta}_{VS} - \beta) \xrightarrow{D} \underline{W}$ where the cdf of \underline{W}

is $F_{\underline{W}}(\underline{z}) = \sum_j \pi_j F_{\underline{W}_j}(\underline{z})$. Thus \underline{W} is a mixture dist of the \underline{W}_j with prob's π_j .

38) $\hat{\beta}_{S, MIX}$ seems to be a good approx
for $\hat{\beta}_{S, VS}$ unless the predictors are
highly correlated. Most of the selection
bias is due to predictors in E
which make selection of S almost random.
 $\hat{\beta}_{E, MIX}$ and $\hat{\beta}_{E, VS}$ tend to differ, but
both use zero padding.

Let $\beta = (\beta_1, \beta_2, 0, 0)^T = (\beta_1, \beta_2, \beta_3, \beta_4)^T$

ex) Suppose $\sqrt{n}(\hat{\beta}_S - \beta_S) \xrightarrow{D} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 \begin{pmatrix} v_{11} & v_{12} \\ v_{21} & v_{22} \end{pmatrix} \right)$

OLS CLT

$= \sqrt{n} \left(\begin{pmatrix} \hat{\beta}_{1S} \\ \hat{\beta}_{2S} \end{pmatrix} - \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} \right)$

Then $\sqrt{n}(\hat{\beta}_{S,0} - \beta) \xrightarrow{D} N_4 \left(\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \sigma^2 \begin{pmatrix} v_{11} & v_{12} & 0 & 0 \\ v_{21} & v_{22} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \right)$

$= \sqrt{n} \left(\begin{pmatrix} \hat{\beta}_{1S} \\ \hat{\beta}_{2S} \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} \beta_1 \\ \beta_2 \\ 0 \\ 0 \end{pmatrix} \right)$

$\sigma^2 \begin{pmatrix} v_{11} & v_{12} & 0 & 0 \\ v_{21} & v_{22} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$

$\hat{\beta}_{MIX}$ has a mixture dist of $\hat{\beta}_{S,0}$ $j=1, \dots, 4$ asymptotically

if $P(S \subseteq I_{min}) \rightarrow 1$. $\beta_{I_1} = \beta_S$, $\beta_{I_2} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$, $\beta_{I_3} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_4 \end{pmatrix}$, $\beta_{I_4} = \beta$

$a_1=2$ $a_2=3$ $a_3=3$ $a_4=4$

39) Geometric Argument. Suppose there

is an iid sample T_1, \dots, T_n
 Problem usually $\beta = 1$

of size B of the statistic. A large sample 100 (1- δ)% prediction region for

T_n is $\left\{ \underline{w} : D_{\underline{w}}^2(\bar{T}, \hat{\underline{w}}_T) \leq D_{(B)}^2 \right\}$

Need $\bar{T} \rightarrow \underline{\theta}$ faster than $T_i \rightarrow \underline{\theta}$: eg $E(T_i) = \underline{\theta}$ with MCLT or $\sqrt{n}(\bar{T} - \underline{\theta}) \rightarrow N_2(\underline{0}, \Sigma)$

Now $D_{T_i}^2(\bar{T}, \hat{\underline{w}}_T) = D_{\bar{T}}^2(T_i, \hat{\underline{w}}_T)$

\uparrow center of hyperellipsoid

See Figure on back of HW 10.



$\frac{95}{100}$ T_i are in the prediction region

contains $\approx 95\%$ of T_i if n, B are large

So $\bar{T} \in \{ \underline{W} : D_{\underline{W}}^2(T_i, \hat{\mu}_T) \leq D_{(\underline{W})}^2 \}$ iff

$T_i \in \{ \underline{W} : D_{\underline{W}}^2(\bar{T}, \hat{\mu}_T) \leq D_{(\underline{W})}^2 \}$ which occurs

with $\approx 100(1-\delta)\%$ e.g. 95%. If \bar{T} goes to $\underline{\theta}$

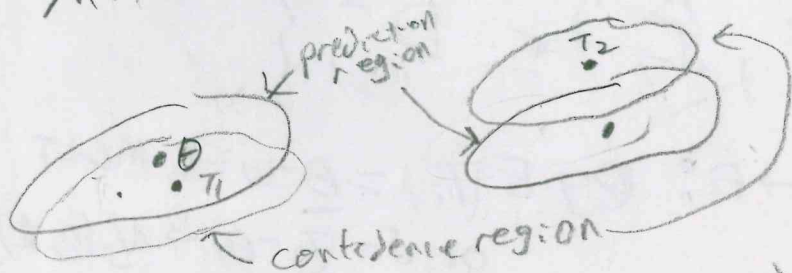
(fast enough, then $P[\underline{\theta} \in \{ \underline{W} : D_{\underline{W}}^2(T_i, \hat{\mu}_T) \leq D_{(\underline{W})}^2 \}] \rightarrow 1-\delta$.

e.g. $E \bar{T} = \underline{\theta}$ $\text{cov}(\bar{T}) = \frac{\Sigma}{B}$ $\text{cov}(T_i) = \Sigma$

or $\sqrt{n}(\bar{T} - \underline{\theta}) \xrightarrow{D} N_g(\underline{\theta}, \Sigma)$ then $\sqrt{n}(\bar{T} - \underline{\theta}) \xrightarrow{D} N_g(\underline{\theta}, \frac{\Sigma}{B})$

40) If $\sqrt{n}(\bar{T}_n - \underline{\theta}) \xrightarrow{D} \underline{U}$ and $\sqrt{n}(\bar{T}_i^* - \bar{T}_n) \xrightarrow{D} \underline{U}$

Then the bootstrap sample takes the iid cloud of the T_1, \dots, T_B centered at $\underline{\theta}$ (or \bar{T}) and shifts the cloud to be centered at \bar{T}_n .



So by the Geometric Argument, the prediction region

Bootstrap and Ron and Hybrid methods give large sample confidence regions. Need $B \geq 50P$ if $\underline{\theta}$ is $p \times 1$

and need n large, sometimes LM 79 extremely large. If n is not large enough, undercoverage ($<$ nominal 100 (1- δ)% coverage) can occur. For OLS var set want $n \geq 20p$, error dist unimodal, not highly skewed.

4) For variable selection with the residual and parametric bootstrap (and the nonparametric bootstrap under strong regularity conditions) it can be shown that for large n , the bootstrap sample cloud T_n^* , $(\hat{\beta}_I^*, \hat{\beta}_{I,0}^*)$ tends to be slightly more variable than the iid data cloud T_n , $(\hat{\beta}_I, \hat{\beta}_{I,0})$.

Suppose $S \subseteq I$, then

$$E(\hat{\beta}_I) = \beta_I, \quad E(\hat{\beta}_{I,0}) = \beta_{I,0}$$

$$\text{COV}(\hat{\beta}_I) \rightarrow \sigma^2 V_I, \quad \text{COV}(\hat{\beta}_{I,0}) \rightarrow \sigma^2 V_{I,0}$$

$$\text{COV}(\hat{\beta}_I^*) \rightarrow \sigma^2 V_I, \quad \text{COV}(\hat{\beta}_{I,0}^*) \rightarrow \sigma^2 V_{I,0}$$

$$E(\hat{\beta}_I^*) = \hat{\beta}_I \quad \text{and} \quad E(\hat{\beta}_{I,0}^*) = \hat{\beta}_{I,0}$$

For the residual bootstrap $\text{COV}(\hat{\beta}^*) \approx \frac{1-p}{n} \text{COV}(\hat{\beta})$.

$$\text{Let } T_n = A \hat{\beta}_{I_{\min},0} \quad \text{and} \quad \theta = A \beta$$

Let the j -th component of the bootstrap cloud be

$$A \hat{\beta}_{I_j,0}^* \quad (\text{I}_j \text{ selected})$$

B_jn times

while the j th component of the iid cloud is $A \hat{\beta}_{I_j,0}, \dots, A \hat{\beta}_{I_j,0}^{*}$

(I_j selected n_j times),

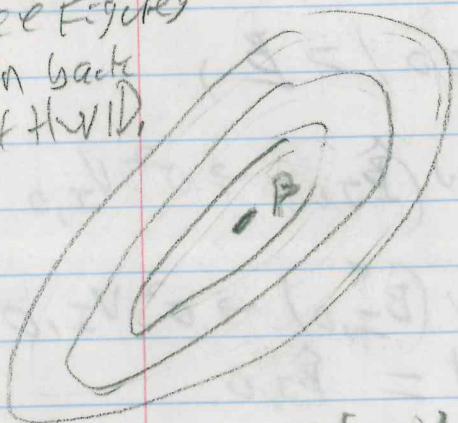
Asymptotically, the bootstrap and iid sample component clouds have the same variability and shape, but the iid sample component clouds are centered at $E(\hat{\beta}_{I_j,0}) = \beta$

while the bootstrap sample component clouds are centered at

$$E(\hat{\beta}_{I_j,0}^*) = \hat{\beta}_{I_j,0}$$

Different

see figures on back of HW11



iid data cloud

separate so clouds are centered at $\hat{\beta}_{I_j,0}$

Makes the overall cloud have greater variability.

Hence by the Geometric Argument expect coverage \geq nominal.

Go to MREG

Undercoverage can occur if the bootstrap data cloud is less variable than the iid cloud eg $\frac{n-p}{n}$ NOT

Multivariate Linear Regression

LM 88
OK 8 in course notes

1) Multivariate linear regression (mreg)

is a special case of the multivariate linear model where at least one predictor variable takes on many values. MANOVA is another special case where the values of X are coded often $\{-1, 0, 1\}$.

See handouts 164) - 190) and references on HW 11 ^{update}

2) The mreg model is $\underset{\sim}{y}_i = B^T \underset{\sim}{x}_i + \underset{\sim}{\epsilon}_i$

for $i=1, \dots, n$ with $m \geq 2$ response variables y_1, \dots, y_m and p predictor variables

$x_1 \equiv 1, x_2, x_3, \dots, x_p$. The i th case
non trivial predictors

$= (\underset{\sim}{x}_i^T, \underset{\sim}{y}_i^T)^T$ where $x_{i1} = 1$ is often omitted.

data matrix $\begin{pmatrix} \underline{x}_1^T & \underline{y}_1^T \\ \vdots & \vdots \\ \underline{x}_n^T & \underline{y}_n^T \end{pmatrix}$ but column of 1s
is usually omitted

2) In matrix form,

$$\underline{Z} = \underline{X} B + E \quad \text{where}$$

$n \times m$ $n \times p$ $p \times m$ $n \times m$

$$\underline{Z}_{n \times m} = \left(\underline{y}_1, \underline{y}_2, \dots, \underline{y}_m \right) = \begin{pmatrix} \underline{y}_1^T \\ \vdots \\ \underline{y}_n^T \end{pmatrix}$$

$$B = \left[\underline{B}_1 \quad \underline{B}_2 \quad \dots \quad \underline{B}_m \right], \quad E = \left[\underline{e}_1 \quad \dots \quad \underline{e}_m \right] = \begin{pmatrix} \underline{\varepsilon}_1^T \\ \vdots \\ \underline{\varepsilon}_n^T \end{pmatrix}$$

$$E(\underline{\varepsilon}_k) = \underline{0}, \quad \text{COV}(\underline{\varepsilon}_k) = \underline{\Sigma}_{\underline{\varepsilon}_k} = (\sigma_{ij}), \quad k=1, \dots, m$$

$$E(\underline{e}_i) = \underline{0}, \quad \text{COV}(\underline{e}_i, \underline{e}_j) = \sigma_{ij} I_n \quad \text{for}$$

$i, j = 1, \dots, m,$

3) The \underline{e} 's are no longer residuals. Use $\underline{\hat{\varepsilon}}, \underline{\Sigma}$

and $\hat{\epsilon}_{ij}$ for residuals.

LM 87

4) Each response variable Y_j in the m reg model follows a univariate multiple linear regression (MLR) model

$$Y_j = \underset{n \times 1}{\tilde{Y}_j} = \underset{n \times p}{\tilde{X}} \underset{p \times 1}{\beta_j} + \underset{n \times 1}{e_j} \quad j=1, \dots, m$$

$$Y_{ij} = \underset{1 \times p}{x_i^T} \underset{p \times 1}{\beta_j} + e_i = \beta_{j0} + x_{i2} \beta_{j2} + \dots + x_{ip} \beta_{jp} + e_i$$

\uparrow
error not residual

with $x_{i1} \equiv 1$, \tilde{X} does not depend on j so the same predictors are used for all m response variables.

5) $m=1 \rightarrow$ MLR model, so for m reg we usually assume $m \geq 2$.

6) Suppose $p=1$ and $\tilde{X} = \underline{1}$, $\underline{Y} = \tilde{X} B + E$

$$\text{Then } \underset{n \times m}{\begin{pmatrix} \underline{Y}_1 & \dots & \underline{Y}_m \end{pmatrix}} = \underset{n \times 1}{\underline{1}} \overbrace{\begin{pmatrix} \beta_1 & \dots & \beta_m \end{pmatrix}}^{\underline{\mu}^T} + \begin{pmatrix} e_1 & \dots & e_m \end{pmatrix}$$

$1 \times m$

So

$$\begin{pmatrix} y_1^T \\ \vdots \\ y_n^T \end{pmatrix} = \begin{pmatrix} \underline{\mu}^T \\ \vdots \\ \underline{\mu}^T \end{pmatrix} + \begin{pmatrix} \underline{\varepsilon}_1^T \\ \vdots \\ \underline{\varepsilon}_n^T \end{pmatrix}$$

and $y_i = \underline{\mu} + \underline{\varepsilon}_i \quad i=1, \dots, n$

$m \times 1$ $m \times 1$ $m \times 1$

m location models $Y_{ij} = \mu_j + \varepsilon_{ij}$

$i=1, \dots, m, j=1, \dots, n$, Here y_1, \dots, y_n are iid,
 $E(\varepsilon_{ij}) = \mu_j$ $\text{cov}(y_i) = \Sigma = \Sigma_j$.

This model is the multivariate location and dispersion model and the (multivariate) least squares (LS) estimator is

$$\hat{\underline{\mu}} = \overline{y} = \frac{1}{n} \sum_{k=1}^n y_k = \begin{pmatrix} \overline{Y}_1 \\ \overline{Y}_2 \\ \vdots \\ \overline{Y}_m \end{pmatrix} \text{ where } \overline{Y}_j = \frac{1}{n} \sum_{i=1}^n Y_{ij} = \overline{Y}_{0j} \text{ for } j=1, \dots, m.$$

Note $\hat{\mu}_j = \overline{Y}_j = \overline{Y}_{0j}$ as in the univariate location model.

For each variable Y_k , make a response plot of \hat{Y}_{ik} vs Y_{ik} and a

0's except for a 1 in the position corresponding to the i th variable to be deleted. To delete the last k variables

use $L = \begin{bmatrix} \bar{0} & I_k \end{bmatrix}$. The MANOVA F

test and F_j test are special cases.

(0) The MANOVA F test uses

$L = \begin{bmatrix} \bar{0} & I_{p-1} \end{bmatrix}$, H_0 : the nontrivial predictors are not needed in the m reg model

H_1 : at least one of the nontrivial predictors is needed

This test is the analog of the MLR Anova F_j test

(1) The F_j test uses $L_j = \begin{bmatrix} \bar{0}, \dots, 0, 1, 0, \dots, 0 \end{bmatrix}$ ^{j th position}
 $1 \times p$

H_0 $\underline{b}_j^T = \underline{0}^T \Leftrightarrow$ predictor X_j is not needed in the m reg model given the other predictors are in the model

H_1 $\underline{b}_j^T \neq \underline{0}^T \Leftrightarrow$ predictor X_j is needed

The F_j test is the analog of the MLR

F tests for $H_0 \beta_j = 0$.

MLR was $\hat{\beta}$

model in matrix form

$$Y = X\beta + e$$

$n \times 1$ $n \times p$ $p \times 1$ $n \times 1$

$$E(Y) = X\beta$$

$$E(e) = 0$$

$$H = P = X(X'X)^{-1}X'$$

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

$$\hat{Y} = PY$$

$$\hat{e} = (I - P)Y$$

$$E(\hat{\beta}) = \beta$$

$$E(\hat{Y}) = E(Y) = X\beta$$

$$\hat{\sigma}^2 = MSE = \frac{e'e}{n-p} = \frac{\sum e_i^2}{n-p}$$

response and residual plots

LS CLT

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

"F-stat" \xrightarrow{D}

MREG LM 89.5 and Q11

$$Z = XB + E$$

$n \times m$ $n \times p$ $p \times m$ $n \times m$

$$E(Z) = XB$$

$$E(E) = 0, E(E_j) = 0$$

$$H = P = X(X'X)^{-1}X'$$

$$\hat{\beta} = (X'X)^{-1}X'Z$$

$\hat{\beta}_j = (X'X)^{-1}X'z_j$

$$\hat{Z} = PZ, \hat{z}_j = PY_j$$

$$\hat{E} = (I - P)Z, \hat{e}_j = (I - P)z_j$$

$$E(\hat{\beta}) = \beta, E(\hat{\beta}_j) = \beta_j$$

$$E(\hat{Z}) = E(Z) = XB$$

$E(z_j) = X\beta_j$

$$\hat{\sigma}^2 = \frac{\hat{E}'\hat{E}}{n-p}, \hat{\sigma}_j^2 = \frac{\hat{e}_j'\hat{e}_j}{n-p}$$

m pairs of response and residual plots

MLS CLT

$$\sqrt{n} \text{vec}(\hat{\beta} - \beta) \xrightarrow{D} N(0, \sigma^2 \otimes W)$$

$$\sqrt{n}(\hat{\beta}_j - \beta_j) \xrightarrow{D} N_p(0, \sigma_j^2 W)$$

MLR
 $H_0: LB = 0$ L $r \times p$

$r F_R \xrightarrow{D} \chi^2_r$

$F_R = \underline{L\hat{\beta}}^T \left(L(X^T X)^{-1} L^T \right)^{-1} \underline{L\hat{\beta}}$

$r \hat{\sigma}_{rest}^2$
 $\approx F_{r, n-p}$

MREG
 $H_0: LB = 0$ L $r \times p$
 $(n-p) U(L) \xrightarrow{D} \chi^2_{rm}$

$\frac{(n-p) U(L)}{r} =$

$\frac{\text{vec}(\underline{L\hat{\beta}})^T \left[\hat{\Sigma}^{-1} \otimes (L(X^T X)^{-1} L^T)^{-1} \right] \text{vec}(\underline{L\hat{\beta}})}{r}$

$\approx F_{rm, n-mp}$

$E[y_i] = x_i^T \underline{\beta} = \underline{\beta}^T x_i$

$E(\underline{y}_i) = \underline{B}^T \underline{x}_i$

$E[y_{ij}] = \underline{x}_i^T \underline{\beta}_j = \underline{\beta}_j^T \underline{x}_i$

$V(e_i) = \sigma^2$

$\text{cov}(\underline{\varepsilon}_i) = \underline{\Sigma}_{\varepsilon}$

$y_i = \underline{x}_i^T \underline{\beta} + e_i$
 $= \underline{\beta}^T \underline{x}_i + e_i$

$\underline{y}_i = \underline{B}^T \underline{x}_i + \underline{\varepsilon}_i$

$\underline{y}_i^T = \underline{x}_i^T \underline{B} + \underline{\varepsilon}_i^T$