

16) The prediction region method large sample $100(1-\delta)\%$ confidence region for $\underline{\theta}$ is (33')

$$\left\{ \underline{w} : \underbrace{(\underline{w} - \bar{T}^*)^T [\underline{S}_T^*]^{-1} (\underline{w} - \bar{T}^*)}_{D_{\underline{w}}^2(\bar{T}^*, S_T^*)} \leq D_{(U_B)}^2 \right\}$$

where $D_i^2 = (\bar{T}_i^* - \bar{T}^*)^T [\underline{S}_T^*]^{-1} (\bar{T}_i^* - \bar{T}^*)$.

Reject $H_0: \underline{\theta} = \underline{\theta}_0$ if $(\bar{T}^* - \underline{\theta}_0)^T [\underline{S}_T^*]^{-1} (\bar{T}^* - \underline{\theta}_0) > D_{(U_B)}^2$.

17) The (modified) Bickel and Ren large sample $100(1-\delta)\%$ confidence region is

$$\left\{ \underline{w} : (\underline{w} - T_n)^T (S_T^*)^{-1} (\underline{w} - T_n) \leq D_{(U_B, T)}^2 \right\}$$

\downarrow 100th Percentile

where $D_i^2 = (\bar{T}_i^* - T_n)^T (S_T^*)^{-1} (\bar{T}_i^* - T_n)$.

18) The hybrid large sample $100(1-\delta)\%$ confidence region is

$$\left\{ \underline{w} : (\underline{w} - T_n)^T (S_T^*)^{-1} (\underline{w} - T_n) \leq D_{(U_B)}^2 \right\}$$

19) 16) is the nonparametric prediction region applied to the bootstrap sample.

20) Ch 2 applied 2 PIs to the bootstrap sample to get large sample CIs.

21) The confidence regions tend to work 29.2

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

where $E(\underline{U}) = \underline{0}$ and $\text{cov}(\underline{U}) = \underline{\Sigma}_U \succ 0$.

Positive definite

WRT the bootstrap distr, \bar{T}_n is a constant,

$$\text{and } \sqrt{n} \begin{pmatrix} \bar{T}_1^* - \bar{T}_n \\ \vdots \\ \bar{T}_B^* - \bar{T}_n \end{pmatrix} \xrightarrow{D} \begin{pmatrix} \underline{V}_1 \\ \vdots \\ \underline{V}_B \end{pmatrix}, \quad \underline{V}_i \stackrel{iid}{\sim} \underline{U}, \text{ multiply}$$

both sides by $A = (\frac{1}{B} \dots \frac{1}{B})$ to take the sample mean of both sides.

$$\text{For fixed } B, \text{ get } \sqrt{n}(\bar{T}^* - \bar{T}_n) \xrightarrow{D} \frac{1}{B} \sum_{i=1}^B \underline{V}_i$$

As $B \rightarrow \infty$, $\sqrt{n}(\bar{T}^* - \bar{T}_n) \xrightarrow{D} \underline{0}$ so \bar{T}^* and \bar{T}_n are asymptotically equivalent and $\sqrt{n}(\bar{T}^* - \theta) \xrightarrow{D} \underline{U}$.

$$\underline{U}' = \sqrt{n}(\bar{T}^* - \bar{T}_n + \bar{T}_n - \theta) = \sqrt{n}(\bar{T}^* - \bar{T}_n) + \sqrt{n}(\bar{T}_n - \theta) \xrightarrow{D} \underline{0} + \underline{U} = \underline{U}$$

22) Comparing the iid data cloud T_1, \dots, T_B with the bootstrap data cloud T_1^*, \dots, T_B^* is useful where $T_i = \bar{T}_n$. If $\sqrt{n}(\bar{T}_n - \theta) \xrightarrow{D} \underline{U}$

and $\sqrt{n}(\bar{T}_i^* - \bar{T}_n) \xrightarrow{D} \underline{U}$, then the bootstrap data cloud is like the iid data cloud (which is centered at θ) shifted to be centered at \bar{T}_n .

23} Geometric argument: suppose

$$\sqrt{n}(\bar{T}_n - \underline{\theta}) \xrightarrow{D} \underline{0} \quad \text{with } E(\underline{0}) = \underline{0} \text{ and } \text{cov}(\underline{0}) = \underline{\Sigma}_0$$

Assume T_1, \dots, T_B are iid with nonsingular covariance matrix $\underline{\Sigma}_{T_n}$. Assume $(\underline{1}, \underline{S}_T)^T \perp \underline{\Sigma}_A^{-1}$.

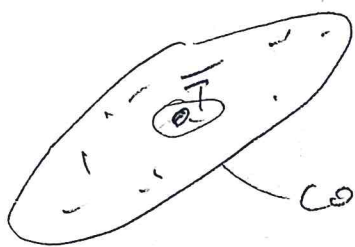
The nonparametric prediction region applied to T_1, \dots, T_B is $R_p = \left\{ \underline{\omega} : D_{\underline{\omega}}^2(\bar{T}, \underline{S}_T) \leq D_{(\underline{0}, \underline{B})}^2 \right\}$ is centered at \bar{T} and contains a future value of the statistic T_B with prob $1 - \delta_B \rightarrow t_\delta$ as $B \rightarrow \infty$ (eventually bounded below).

The region $R_c = \left\{ \underline{\omega} : D_{\underline{\omega}}^2(T_n, \underline{S}_T) \leq D_{(\underline{0}, \underline{B})}^2 \right\}$ is a large sample confidence region for $\underline{\theta}$ where T_n is a randomly selected T_i .

Proof} R_c contains \bar{T} with prob $1 - \delta_B$

$$D_{T_i}(\bar{T}, \underline{S}_T) = D_{\bar{T}}(T_i, \underline{S}_T)$$

\swarrow centered at \bar{T} \swarrow centered at T_i



contains $100(1 - \delta_B)\%$ of the T_i

Any T_i in the hyperellipsoid R_p has R_c that contains \bar{T} . Any T_i outside of R_p has R_c that does not contain \bar{T} .

Ex B, $\sqrt{n}(\bar{T} - \underline{\theta}) \xrightarrow{D} \frac{1}{B} \sum_{i=1}^B \underline{v}_i$, $\underline{v}_i \stackrel{iid}{\sim} \underline{0}$. So \bar{T} gets arbitrarily close to $\underline{\theta}$ as $B \rightarrow \infty$ compared to the T_i .

So R_c contains $\underline{\theta}$ with prob $\rightarrow 1 - \delta_B \rightarrow t_\delta$.

contains θ with 99% prob is not a point. (H.S)

So the coverage tends to be less than 1- α if

$\frac{2}{3} = \sqrt{B} \pm \delta$. Increasing the coverage using $\frac{2}{3}$ makes the observed coverage closer to the nominal 1- α .

25) The bootstrap data cloud is like the iid data cloud skitted to be centered at T_n . So applying the nonparametric prediction region to the bootstrap data cloud gives a confidence region.

26) Suppose $\sqrt{n}(T_n - \theta) \xrightarrow{D} U$, $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} U$, and $\sqrt{n}(T_i^* - \bar{T}^*) \xrightarrow{D} U$ and $\sqrt{n}(\bar{T}^* - \theta) \xrightarrow{D} U$.

see 21). Then

$$D_1^2 = D_{T_i^*}^2(\bar{T}^*, S_T^*) = \sqrt{n}(T_i^* - \bar{T}^*) (n S_T^*)^{-1} \sqrt{n}(T_i^* - \bar{T}^*)$$

$$D_2^2 = D_{\theta}^2(T_n, S_T^*) = \sqrt{n}(T_n - \theta) (n S_T^*)^{-1} \sqrt{n}(T_n - \theta)$$

$$D_3^2 = D_{\theta}^2(\bar{T}^*, S_T^*) = \sqrt{n}(\bar{T}^* - \theta) (n S_T^*)^{-1} \sqrt{n}(\bar{T}^* - \theta)$$

$$D_4^2 = D_{T_i^*}^2(T_n, S_T^*) = \sqrt{n}(T_i^* - T_n) (n S_T^*)^{-1} \sqrt{n}(T_i^* - T_n)$$

If $(n S_T^*)^{-1} \xrightarrow{P} \Sigma_A^{-1}$, then $D_j^2 \xrightarrow{D} U^T \Sigma_A^{-1} U$. If

$(n S_T^*)^{-1}$ is "not too ill conditioned," then

$D_j^2 \approx U^T (n S_T^*)^{-1} U$ for large n . Then the 3 bootstrap

confidence regions will have coverage near 1- α .

27) Data splitting divides the training data $\underline{x}_1, \dots, \underline{x}_n$ into two sets H and V where H has $n_H \geq \frac{1}{2}$ of the cases and V has $n - n_H = n_V$ of the cases. The estimator (\bar{T}_H, C_H) is computed using the data set H . Then the squared validation distances

$$D_j^2 = D_{\underline{x}_j}^2(\bar{T}_H, C_H) = (\underline{x}_j - \bar{T}_H)^T C_H^{-1} (\underline{x}_j - \bar{T}_H)$$

are computed for the $j=1, \dots, n_V$ cases in the validation set V . Let $D_{(U_V)}^2$ be the U th order statistic of the D_j^2 where

$$U_V = \text{med}(n_V, \lceil (n_V+1)(1-\delta) \rceil).$$

28) The large sample $100(1-\delta)\%$ data splitting prediction region for \underline{x} is

$$\left\{ \underline{z} : D_{\underline{z}}^2(\bar{T}_H, C_H) \leq D_{(U_V)}^2 \right\}.$$

29) Can use $(\bar{T}_H, C_H) = (\bar{T}_{RMVN}, C_{RMVN})$, or,

$$\text{if } p > n, (\bar{T}_H, C_H) = (\bar{X}, I) \text{ or } (\text{MED}(W), I).$$

30) $P_{(1)}, \dots, P_{(n)}$ have rank $1, \dots, n$.

Consider D_K^2 for $K=1, \dots, N_V, N_V+1$ where $D_{N_V+1}^2$ is for X_t . Since the N_V+1 cases are iid, the prob that D_X^2 has rank j for $j=1, \dots, N_V+1$ is $\frac{1}{N_V+1}$ for each j . So the ranks follow the discrete uniform distribution.

Let $D_{(j)}^2$ be the order statistics using $j=1, \dots, N_V$ without using the unknown $D_{N_V+1}^2$. Then

$D_{(j)}^2$ has rank j if $D_{(j)}^2 < D_{N_V+1}^2$ but rank $j+1$ if $D_{(j)}^2 > D_{N_V+1}^2$. Thus $D_{(U_V)}^2$ has rank U_V+1 if $D_{X_t}^2 < D_{(U_V)}^2$ and

$$P\left(X_t \in \left\{z : D_z^2(T_H, C_H) \leq D_{(U_V)}^2\right\}\right) = P\left(D_{X_t}^2 \leq D_{(U_V)}^2\right)$$

$$\geq \frac{U_V}{N_V+1} \rightarrow 1-\delta \text{ as } N_V \rightarrow \infty.$$

\uparrow
 \equiv if there are no ties.

$$\underbrace{\left(P(X_t \text{ has rank } \leq U_V) \right)}_{= P(X_t \text{ has rank } < U_V+1)}$$

32) N_V 1 2 19 20 39 59 79 99

$\alpha = .05$ $\frac{U_V}{N_V+1}$ $\frac{1}{2}$ $\frac{2}{3}$ $\frac{19}{20} = .95$ $\frac{20}{21}$ $\frac{38}{40} = .95 = \frac{57}{60}$ $\frac{76}{80} = .95 = \frac{95}{100}$

\approx so chance

33) Need $N_V \geq 100$ to estimate the $100(1-\delta)$ percentile of the D_j^2 , $.05 \leq \delta \leq .5$. Volume can be huge for small N_V .

			Zemen, Samuel G.
--	--	--	------------------

or if (T_H, C_H) is bad

1) The multiple linear regression (MLR) model (36)

$$y_i = x_{i1}\beta_1 + \dots + x_{ip}\beta_p + e_i = \underline{x}_i^T \underline{\beta} + e_i \quad i=1, \dots, n$$

In matrix form $\underline{y} = \underline{X} \underline{\beta} + \underline{e}$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} \quad \text{where}$$

$$\underline{X} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix}.$$

2) MLR is important for Math 484, 584 (linear model theory) 586 (statistical learning) and robust statistics.

- 3) a) The constant variance MLR model has the e_i iid, $E(e_i) = 0$, $V(e_i) = \sigma^2$.
- b) The unimodal MLR model is as in a) but adds the assumption that the e_i are from a unimodal distribution that is not highly skewed.
- c) The normal MLR model has $e_i \stackrel{iid}{\sim} N(0, \sigma^2)$.

§5.4

To get an asymptotically optimal PI for Y_e given \underline{x}_e and $\underbrace{(Y_1, \underline{x}_1), \dots, (Y_n, \underline{x}_n)}_{\text{training data}}$, $\underbrace{\quad}_{\text{test data}}$

Find the shortest estimator $\{\hat{\Gamma}(e), \hat{\Gamma}(e+1)\}$ of the residuals and use $\{\hat{Y}_e + \hat{\Gamma}(e), \hat{Y}_e + \hat{\Gamma}(e+1)\}$.

$\hat{Y}_e \xrightarrow{P} E(Y_e | \underline{x}_e)$ and $\{\hat{\Gamma}(e), \hat{\Gamma}(e+1)\}$ estimates the pop shortk of the iid e_i .

Here $a_{ii} = (1 + \frac{15}{n}) \sqrt{\frac{n}{n-p}} \sqrt{1+h_e}$ where

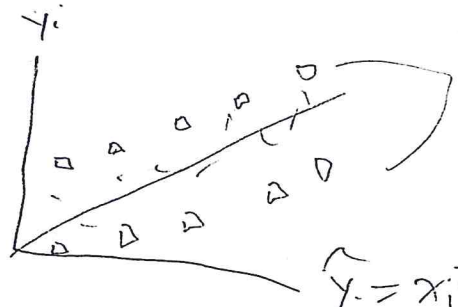
$$r_i \approx \sqrt{\frac{n}{n-p}} e_i \quad \text{in that } \sigma^2 = \frac{1}{n-p} \sum_{i=1}^n r_i^2 \approx \frac{1}{n} \sum_{i=1}^n e_i^2$$

and $h_e = \underline{x}_e^T (X^T X)^{-1} \underline{x}_e$ is the leverage of \underline{x}_e .

want $h_e \leq \max \underline{x}_i^T (X^T X)^{-1} \underline{x}_i = h_i$ where

h_i is the i th diagonal element of $H = X(X^T X)^{-1} X^T$.

response plot



PI limits not quite parallel to identity line

$$\hat{Y}_i = \underline{x}_i^T \hat{\beta} = ESP_i$$

stem §1.5