

$[L_n(\underline{y}), U_n(\underline{y})]$  is a large sample

$100(1-\delta)\%$  CI for  $\theta$  if  $P(L_n \leq \theta \leq U_n) = P(L_n(\underline{Y}) \leq \theta \leq U_n(\underline{Y}))$  is eventually bounded below by  $1-\delta$  as  $n \rightarrow \infty$ .

2} A large sample  $100(1-\delta)\%$  confidence region for  $\theta$  is a set  $A_n$  such that  $P(\theta \in A_n)$  is eventually bounded below by  $1-\delta$  as  $n \rightarrow \infty$ .

CIs are a special case.

3} A statistic is a function of the data that does not depend on any unknown parameters  $\theta$ .

4} <sup>p231</sup> The quantity  $R(\underline{Y} | \theta)$  is a pivot or pivotal quantity if the distribution of  $R(\underline{Y} | \theta)$  does not depend on  $\theta$ .  $R(\underline{Y}, \theta)$  is an asymptotic pivot or asymptotic pivotal quantity if the limiting dist of  $R(\underline{Y}, \theta)$  does not depend on  $\theta$ .

$$\text{ex) a) } \sqrt{n} (T_n - \theta) \xrightarrow{D} N_p(0, \Sigma)$$

67.5

does not depend on  $\theta$   
if  $\Sigma$  doesn't depend  
on  $\theta$   
(rare)

$$\sqrt{n} \hat{\Sigma}^{-1/2} (T_n - \theta) \xrightarrow{D} N_p(0, I)$$

useful asymptotic pivot

$$\text{then } n (T_n - \theta)^T \hat{\Sigma}^{-1} (T_n - \theta) \xrightarrow{D} \chi_p^2$$

$$\text{b) } Y_1, \dots, Y_n \text{ iid } E(Y_i) = \mu, \quad V(Y_i) = \sigma^2$$

$$\sqrt{n} \frac{(\bar{Y} - \mu)}{\sigma} \xrightarrow{D} N(0, 1)$$

$$\sqrt{n} (\bar{Y} - \mu) \xrightarrow{D} N(0, \sigma^2)$$

need  
 $\sigma$  known  
to be a useful  
pivot

$$\sqrt{n} \left( \frac{\bar{Y} - \mu}{\sigma} \right) \sim N(0, 1), \quad Y_i \text{ iid } N(\mu, \sigma^2)$$

$$\sqrt{n} \left( \frac{\bar{Y} - \mu}{s} \right) \sim t_{n-1}, \quad Y_i \text{ iid } N(\mu, \sigma^2)$$

$$\sqrt{n} \left( \frac{\bar{Y} - \mu}{s} \right) \xrightarrow{D} N(0, 1), \quad Y_i \text{ iid as above}$$

useful asymptotic pivot

## C) Linear Models

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

$$E(e_i) = 0, \quad V(e_i) = \sigma^2$$

$\underline{x}_i = (1, x_{2i}, \dots, x_{pi})^T$  so a constant is in the model,  $i=1, \dots, n$ .

$$H_0 \underline{A} \underline{\beta} = \underline{0} \quad H_A \underline{A} \underline{\beta} \neq \underline{0}$$

$r \times p$   $A$  full rank  $r$

$$F_0 = \frac{1}{r \text{MSE}} (\underline{A} \hat{\underline{\beta}})^T (\underline{A} (\underline{X}^T \underline{X})^{-1} \underline{A}^T)^{-1} (\underline{A} \hat{\underline{\beta}})$$

$\underline{X}$   $n \times p$  with  $i$ th row  $\underline{x}_i^T$ .

If  $e_i \sim N(0, \sigma^2)$ ,  $H_0$  true  $F_0 \sim F_{r, n-p}$

and  $r F_0 \xrightarrow{D} \chi_r^2$  under mild conditions if  $e_i$  iid.

$r F_{r, n-p} \xrightarrow{D} \chi_r^2$ , so the  $F$  tests are large sample tests.

5) Consider a statistic  $T_n$   
 where  $\sqrt{n}(T_n - \theta) \xrightarrow{D} N_p(0, \Sigma_T)$ .

A bootstrap data set,

$$Y_1^*, \dots, Y_n^* \quad \text{or} \quad (Y_1^*, X_1^*), \dots, (Y_n^*, X_n^*)$$

etc, is used to compute  
 a bootstrap statistic  $T_n^*$ .

Repeat  $B$  times to get a  
bootstrap sample  $T_1^*, \dots, T_B^*$ .

6) A case consists of the  
 measurements taken on a person,  
 object or thing to get a data set.

eg  $Y_1, \dots, Y_n$  height  $Y_i = \text{ith case}$

eg  $(Y_1, X_1^T)^T, \dots, (Y_n, X_n^T)^T$   $Y_i = \text{height}$   $X_{i1} = \text{weight}$   
 $X_{i2} = \text{gender of ith person}$   
 $(Y_i, X_i^T)^T$  is the  $i$ th case.

7) Let the data set have  $n$  cases,

The nonparametric bootstrap

creates a bootstrap data set by drawing a sample of size  $n$  with replacement from the  $n$  cases. Then  $T^*$  is computed from the bootstrap data set.

If  $\underline{x}_1, \dots, \underline{x}_n$  is the data set, then the empirical distribution is a discrete distribution where the  $\underline{x}_i$  are equally likely. If  $\underline{w}$  is a random vector from the empirical dist,

then  $P(\underline{w} = \underline{x}_i) = \frac{1}{n} \quad i = 1, \dots, n.$

The cdf of the empirical dist is  $F_n.$



ex) data 1, 2, 3, 4, 5, 6, 7

LS 70

with  $n=7$  and sample median  $T_n = 4$ .

R output gives  $B=2$  bootstrap data sets drawn with replacement from the data.

3 2 3 2 5 2 6

1st bootstrap dataset

(2 2 2 3 3 5 6)

$T_1^* = 3$

3 5 3 4 3 5 7

2nd bootstrap dataset

(3 3 3 4 5 5 7)  $T_2^* = 4$

8) Often  $\sqrt{n}(T_n - \theta) \xrightarrow{D} \underline{U}$  and

large sample theory can be used to

show  $\sqrt{n}(T_n^* - T_n) \xrightarrow{D} \underline{\mu}$ .

An iid sample  $T_{1n}, \dots, T_{Bn}$  of the statistic would be useful, but we only have  $T_n = T_{1n}$ .

$$\sqrt{n} (T_n^* - T_n), \dots, \sqrt{n} (T_{Bn}^* - T_n),$$

70.5

where  $T_i^* = T_n^*$ , is pseudodata

$$\text{for } \sqrt{n} (T_n - \underline{\theta}), \dots, \sqrt{n} (T_{Bn} - \underline{\theta})$$

if  $\sqrt{n} (T_n - \underline{\theta}) \xrightarrow{D} \underline{U}$  and  $\sqrt{n} (T_i^* - T_n) \xrightarrow{D} \underline{U}$ ,

9) For the bootstrap sample

$$T_1^*, \dots, T_B^*, \text{ let}$$

$$\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^*, \quad S_T^* = \frac{1}{B-1} \sum_{i=1}^B (T_i^* - \bar{T}^*) (T_i^* - \bar{T}^*)^T$$

the usual sample mean and sample covariance matrix of  $T_1^*, \dots, T_B^*$ .

$\bar{T}^*$  is known as the bagging estimator  
or smoothed bootstrap estimator.

Let  $T_i$  be a random variable.

LS 71

10) The bootstrap percentile method large sample  $100(1-\delta)\%$  CI for  $\theta$

is an interval  $[T_{(k_L)}^* , T_{(k_U)}^*]$  containing

$\approx [B(1-\delta)]$  of the  $T_i^*$ . Let

$$k_1 = [B \frac{\delta}{2}] \text{ and } k_2 = [B(1 - \frac{\delta}{2})].$$

The usual percentile method CI is

$$[T_{(k_1)}^* , T_{(k_2)}^*].$$

11) The large sample  $100(1-\delta)\%$

Shortcut CI for  $\theta$  is

$$[T_{(s)}^* , T_{(s+c-1)}^*] \text{ where}$$

$$c = \min \left( B, \left\lceil B \left[ 1 - \delta + 1.12 \sqrt{\frac{\delta}{B}} \right] \right\rceil \right).$$

12) Let  $T^*$  be a  $g \times 1$  random vector.

The large sample  $100(1-\delta)\%$  standard bootstrap confidence region is

$$\left\{ \underline{w} : (\underline{w} - T_n)^T [S_T^*]^{-1} (\underline{w} - T_n) \leq \chi_{p, 1-\delta}^2 \right\}$$

$$= \left\{ \underline{w} : D_{\underline{w}}^2(T_n, S_T^*) \leq \chi_{p, 1-\delta}^2 \right\}$$

13) The Prediction region method

large sample 100(1-δ)% conf reg for Θ is

$$\left\{ \underline{w} : (\underline{w} - \bar{T}^*)^T [S_T^*]^{-1} (\underline{w} - \bar{T}^*) \leq D_{(UB)}^2 \right\}$$

$$= \left\{ \underline{w} : D_{\underline{w}}^2(\bar{T}^*, S_T^*) \leq D_{(UB)}^2 \right\} \text{ where}$$

$D_{(UB)}^2$  is computed from

$$D_i^2 = (T_i^* - \bar{T}^*)^T (S_T^*)^{-1} (T_i^* - \bar{T}^*) \text{ for}$$

$i=1, \dots, n$

14) Let  $g_B = \min(1-\delta+0.05, 1-\delta+\frac{g}{B})$ ,  $\delta > 0.1$

$g_B = \min(1-\frac{\delta}{2}, 1-\delta+10\frac{g}{B})$ ,  $\delta \leq 0.01$ .

If  $1-\delta < 0.999$  and  $g_B < 1-\delta+0.001$ ,

Set  $g_B = 1-\delta$ . Let  $D_{(UB)}^2$  be the 100  $g_B$ th