

## INFERENCE ON MEANS USING THE BOOTSTRAP

BY G. JOGESH BABU<sup>1</sup> AND KESAR SINGH<sup>2</sup>

Rutgers University

We study the asymptotic accuracy of the bootstrap approximation to the distribution of a  $k$ -sample studentized mean.

**1. Introduction and main results.** Let  $F_1, F_2, \dots, F_k$  be the distributions of  $k$  populations with means  $\mu_1, \mu_2, \dots, \mu_k$ . Let  $\theta = \sum_{i=1}^k l_i \mu_i$  where  $l_1, l_2, \dots, l_k$  are non-zero constants. Let  $(X_{i1}, X_{i2}, \dots, X_{in_i}), i = 1, 2, \dots, k$ , be independent random samples of sizes  $n_1, n_2, \dots, n_k$  from  $F_1, F_2, \dots, F_k$ . Let  $n$  denote the vector  $(n_1, n_2, \dots, n_k)$  and  $N = \sum_{i=1}^k n_i$ . A natural estimator for  $\theta$  is  $\hat{\theta}_n = \sum_{i=1}^k l_i \bar{X}_i$  and a consistent estimator for its variance is  $\hat{\sigma}_n^2 = \sum_{i=1}^k l_i^2 s_i^2 / n_i$  where  $\bar{X}_i = n_i^{-1} \sum_{j=1}^{n_i} X_{ij}$  and  $s_i^2 = (1/n_i) \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2$ . Here we study the accuracy of the bootstrap approximation to the distribution of the studentized random variable  $t_n = (\hat{\theta}_n - \theta) / \hat{\sigma}_n$ . This approximation is discussed in the next paragraph. Although one could base an inference about  $\theta$  on the difference  $\hat{\theta}_n - \theta$  itself, it turns out that the bootstrap approximation is asymptotically more accurate for  $t_n$  than for  $(\hat{\theta}_n - \theta)$ .

Let  $G_i$  denote the empirical distribution function based on  $(X_{i1}, X_{i2}, \dots, X_{i,n_i}), i = 1, 2, \dots, k$ . The dependence of  $G_i$ 's on the sample sizes is suppressed in the notation. Now let  $(Y_{i1}, Y_{i2}, \dots, Y_{in_i}), i = 1, 2, \dots, k$ , denote independent random samples from the populations  $G_1, G_2, \dots, G_k$ ;  $\bar{Y}_i = n_i^{-1} \sum_{j=1}^{n_i} Y_{ij}$  and  $\gamma_i^2 = n_i^{-1} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$ . Then, by definition, the distribution of  $t_n^* = \sum_{i=1}^k l_i (\bar{Y}_i - \bar{X}_i) / \sum_{i=1}^k l_i^2 \gamma_i^2$  under  $G_1, G_2, \dots, G_k$  is the bootstrap distribution of  $t_n^*$ . Under the conditions given below, the bootstrap distribution of  $t_n^*$  is shown to be asymptotically close to the actual distribution of  $t_n$  up to  $o(N^{-1/2})$ . In applications the bootstrap distribution is approximated by drawing samples of sizes  $n_1, n_2, \dots, n_k$  from  $G_1, G_2, \dots, G_k$  a large number of times, say  $M$  times, calculating  $t_n^*$  each time and finally forming an empirical histogram. It is shown here that this second stage approximation is good up to  $o(N^{-1/2})$  provided  $M/(N \log N) \rightarrow \infty$ .

We now state the main results proved in this note. Throughout, we make the following assumptions, to be referred to as A in the sequel.

A.  $F_i$  has finite 6th moment for all  $1 \leq i \leq k$ . For at least one  $i$ ,  $F_i$  is continuous. Without loss of generality we shall assume that  $F_1$  is continuous. The  $n_i$ 's tend to infinity at the same rate. In other words, the  $N/n_i \leq \lambda < \infty$  for all  $i = 1, 2, \dots, k$ . In practice this last condition means that the  $n_i$ 's are of comparable size.

In what follows, for any distribution  $F$ , let  $F^{-1}(t) = \inf\{x: F(x) \geq t\}$ , where  $0 < t < 1$ .

**THEOREM.** If  $H_n$  denotes the d.f. of  $t_n$  and  $H_n^*$  denotes the d.f. of  $t_n^*$  then, under A, as  $N \rightarrow \infty$

$$(1) \quad N^{1/2} \sup_{x \in \mathbb{R}} |H_n(x) - H_n^*(x)| \rightarrow 0$$

and

$$(2) \quad N^{1/2} |H_n^{-1}(t) - H_n^{*-1}(t)| \rightarrow 0$$

and for all  $t \in (0, 1)$ . Further let  $H_{n,M}$  denote the approximation to  $H_n^*$  described in the second paragraph above with  $M$  samples from  $G_i$ 's. If  $M/(N \log N) \rightarrow \infty$  as  $N \rightarrow \infty$ , then

Received November 1981; revised August 1982.

<sup>1</sup> On leave from the Indian Statistical Institute.

<sup>2</sup> Research supported by NSF Grant MCS-81-02341.

AMS 1970 subject classifications. 62G05; 62G15.

Key words and phrases. Bootstrap, Edgeworth expansions, empirical distribution function, confidence interval.

for almost all sample sequences  $(X_{ij})$

$$(3) \quad N^{1/2} \sup_{x \in R} |H_{n,M}(x) - H_n^*(x)| \rightarrow 0$$

and

$$(4) \quad N^{1/2} |H_{n,M}^{-1}(t) - H_n^{-1}(t)| \rightarrow 0$$

a.s. for all  $t \in (0, 1)$  as  $N \rightarrow \infty$ . The a.s. here refers to the random mechanism generating the samples from  $G$ 's. (We assume that all the second stage sample sequences are defined on the same space.)

It may be mentioned here that (1) in the above theorem is an extension of (1.5) in [8] which is a result involving  $(\bar{X} - \mu)/\sigma$ . For constructing a confidence interval for  $\theta$ , one may replace an actual quantile  $H_n^{-1}(\alpha)$  of  $t_n$  by its bootstrap approximation  $H_{n,M}^{-1}(\alpha)$ . This approximation in the one sample case has been investigated by Efron [6] on simulated data from an asymmetric population. The procedure performed quite well (see Table 5 of [6]).

**2. Proof of the theorem.** We first develop some notation. Let  $\phi_x, \Phi_x$  denote the density and the d.f. of a normal variable with mean zero and dispersion matrix  $\Sigma$ ; let  $\phi, \Phi$  denote the density and the d.f. of a standard normal variable in  $R$ ; let  $c$  denote a constant, the later may denote different constants at different places. For non-negative integral vectors  $\beta = (\beta_1, \dots, \beta_r)$  and  $\mathbf{x} \in R^r$  let  $x^\beta = \prod_{j=1}^r x_j^{\beta_j}$ ,  $\beta! = \beta_1! \beta_2! \dots \beta_r!$ ,  $|\beta| = \beta_1 + \dots + \beta_r$  and  $D^\beta = D_1^{\beta_1} \dots D_r^{\beta_r}$ , where  $D_i^{\beta_i}$  denotes the  $\beta_i$ th order derivative with respect to the  $i$ th variable. Finally let  $\|\mathbf{x}\|^2 = x_1^2 + \dots + x_r^2$  and  $\mathbf{x} \cdot \mathbf{y} = x_1 y_1 + \dots + x_r y_r$ , where  $\mathbf{x} = (x_1, \dots, x_r)$  and  $\mathbf{y} = (y_1, \dots, y_r)$ .

We shall show that

$$(5) \quad P(t_n^* \leq x) = \Phi(x) + N^{-1/2} \int_{-\infty}^x d(y) \phi(y) dy + o(N^{-1/2}) \quad \text{a.s.}$$

where  $d$  is a polynomial whose coefficients depend upon  $F_i$ 's. The same steps will also yield

$$(6) \quad P(t_n \leq x) = \Phi(x) + N^{-1/2} \int_{-\infty}^x d(y) \phi(y) dy + o(N^{-1/2}).$$

Clearly (5) and (6) imply (1). Before proving (5) we shall deduce (2), (3) and (4) from (5) and (6).

To prove (3), first note that in distribution (given the original sample)  $H_{n,M}$  is the same as the empirical d.f. of  $H_n^{-1}(U_i)$  where  $U_1, U_2, \dots, U_M$  are i.i.d.  $U[0, 1]$  random variables. If  $E_M$  denotes the empirical d.f. of  $U_1, \dots, U_M$  then  $\#(H_n^{-1}(U_i) \leq x) = ME_M(H_n(x))$ . Hence, using a well known bound on  $E_n$ , we have

$$\begin{aligned} P(\sup_{x \in R} |H_{n,M}(x) - H_n(x)| \geq 4M^{-1/2}(\log M)^{1/2}) \\ \leq P(\sup_{t \in [0,1]} |E_M(t) - t| \geq 4M^{-1/2}(\log M)^{1/2}) = O(M^{-2}). \end{aligned}$$

Consequently, in view of Borel-Cantelli lemma,

$$\limsup_{M \rightarrow \infty} M^{1/2}(\log M)^{-1/2} \sup_{x \in R} |H_{n,M}(x) - H_n(x)| \leq 4 \quad \text{a.s.}$$

Here dependence of  $M$  on  $N$  is suppressed. The claim (3) in the theorem follows from this, since  $M/(N \log N) \rightarrow \infty$ .

The claims (2) and (4) on quantiles follow using Lemma 1 given below which is an easy consequence of Taylor's expansion.

**LEMMA 1.** Let  $L_N$  be a sequence of d.f.'s on the real line such that, for a polynomial

$a_N$  with its coefficients bounded in  $N$ ,

$$L_N(x) = \int_{-\infty}^x [1 + N^{-1/2} a_N(y)] \phi(y) dy + o(N^{-1/2})$$

uniformly in  $x$ . Then for each  $a \in (0, 1)$ ,

$$L_N^{-1}(a) = z - (\phi(z)\sqrt{N})^{-1} \int_{-\infty}^a a_N(y)\phi(y) dy + o(N^{-1/2})$$

where  $z = \phi^{-1}(a)$ .

The proof of (5) is based on Lemmas 2-5 that follow. In the proofs below we assume w.l.g. that  $i_1 = i_2 = \dots = i_k = 1$ .

All proofs are given for a single sequence of realizations of  $\{X_{ij}\}$  for which  $G_j$  converges weakly to  $F_j$  and  $\int x^k dG_j \rightarrow \int x^k dF_j$  for  $j = 1, 2, \dots, k$ . Thus in view of A the results hold a.s.

**LEMMA 2.** Let  $Y$  be a random vector in  $R^2$  with mean zero and dispersion matrix  $V = ((v_{ij}))$ . Suppose for some  $b > 1$ ,  $\max\{|v_{11}|, |v_{22}|, |v_{21}|\}$ ,  $E\|Y\|^2 < b$ . Let  $a > 2$  be such that  $\Delta(a) < 1/10$ , where  $\Delta(a) = (1/a) + E(\|Y\|^2 I(\|Y\| > a))$ . Then for all  $\|t\| \leq a^{-2}\sqrt{N}$  and all non-negative integral vectors  $\alpha$ , with  $|\alpha| \leq 3$ ,

$$\begin{aligned} & |D^\alpha(g^N(t/\sqrt{N}) - (1 - (i/6\sqrt{N})E(t \cdot Y)^2) \exp(-t' V t/2))| \\ & \leq cb^2(\Delta(a) + N^{-1/2})N^{-1/2}(\|t\|^2 + 1) \exp(-t' V t/2 + ca^{-1}b^2\|t\|^2) \end{aligned}$$

where for  $t \in R^2$ ,  $g(t) = E(\exp(it \cdot Y))$ .

The proof is similar to the proof of Theorem 9.9 of [3].

**LEMMA 3.** Suppose A holds. Let  $\lambda_j = (N/n_j)^{1/2}$ ,  $Z_j = [\lambda_j(Y_{j1} - \bar{X}_j), \lambda_j^2(Y_{j1} - \bar{X}_j)^2 - s_j^2]$ ,  $g_j$  denote the characteristic function of  $Z_j/\sqrt{n_j}$ ,  $B_j$  denote the dispersion matrix of  $Z_j$  and  $B = \sum_{j=1}^k B_j$ . Then for any  $\eta > 0$ ,

$$\begin{aligned} \max_{j|\beta| \leq 3} \int_{\|t\| \leq a\sqrt{N}} \left| D^\beta \left[ \prod_{j=1}^k g_j^\beta(t) - e^{-t' B t/2} \left( 1 - \frac{i}{6\sqrt{N}} \sum_{j=1}^k \lambda_j E(t \cdot Z_j)^2 \right) \right] \right| dt \\ = o(N^{-1/2}). \end{aligned}$$

**PROOF.** Define

$$f_j(t) = (1 - (i/6\sqrt{n_j})E(t \cdot Z_j)^2) \exp(-t' B_j t/2).$$

First note that

$$\begin{aligned} (7) \quad \max_{j|\beta| \leq 3} |D^\beta (e^{-t' B t/2} (1 - (i/6\sqrt{N}) \sum_{j=1}^k \lambda_j E(t \cdot Z_j)^2) - \prod_{j=1}^k f_j(t))| \\ = O(N^{-1}(\|t\|^{2k+1} + 1) \exp(-t' B t/2)), \end{aligned}$$

and for any non-empty subset  $J$  of  $\{1, 2, \dots, k\}$ ,

$$\begin{aligned} (8) \quad \max_{j|\beta| \leq 3} |D^\beta \prod_{j \in J} g_j^\beta(t)| \leq \max_{j|\beta| \leq 3} E \left\| \sum_{j \in J} n_j^{-1/2} \sum_{i=1}^2 Z_{ji} \right\|^{|\beta|} \\ \leq 1 + k^3 \max_{j|\beta| \leq 3} E \left\| \sum_{i=1}^2 Z_{ji} \right\|^2 = O(1), \end{aligned}$$

where  $Z_{ji}$  are independent copies of  $Z_j$ . The last inequality above follows from the proof of Lemma 14.7 of [3] as  $\sup E \|Z_j\|^2 \leq b < \infty$  from some  $b > 0$ . Also for  $1 \leq j \leq k$

$$(9) \quad \max_{j|\beta| \leq 3} |D^\beta \prod_{j \in J} f_j(t)| = O((1 + \|t\|^{2k}) \exp(-t' B t/2)).$$

By (8), (9) and Lemma 2, we have for any  $\alpha > 2$ ,  $|\beta| \leq 3$  and  $\|t\| \leq \lambda^{-1}\alpha^{-\sqrt{N}}$ ,

$$\begin{aligned} & |D^{\beta}(\prod_{i=1}^k g_i^{\beta}(t) - \prod_{i=1}^k f_i(t))| \\ (10) \quad & \leq \sum_{i=1}^k |D^{\beta}(\prod_{i < j} f_i(t) \prod_{j > i} g_j^{\beta}(t))(g_i^{\beta}(t) - f_i(t))| \\ & = O(N^{-1/2}(r(\alpha) + N^{-1/2}))(1 + \|t\|^{s+\alpha}) \exp(-t' B_1 t/2 + c\alpha^{-1} \|t\|^{\beta}), \end{aligned}$$

where  $r(\alpha) = 1/\alpha + \sup_y E(\|Z\|^2 I(\|Z\| > \alpha))$ . It now follows from (7) and (10) that, for a  $|\beta| \leq 3$ , the integral in Lemma 3 over  $\|t\| \leq \lambda^{-1}\alpha^{-\sqrt{N}}$  is  $O(r(\alpha)N^{-1/2}) + O(N^{-1})$ .

Since  $F_1$  is continuous, the dispersion matrix of  $X = (X_{11}, (X_{11} - \mu_1)^2)$  is positive definite and the c.f.  $h$  of  $X$  satisfies the condition  $|h(t)| < 1$  for all  $t \neq 0$ . As a result of this and the fact that weak convergence implies convergence of c.f.'s. uniformly over compact sets, it follows that

$$\sup \|g_i(t)\|; \|t\| \in [\lambda^{-1}\alpha^{-\sqrt{N}}, \eta\sqrt{N}] \leq \delta < 1$$

for all large  $N$ . Also,

$$\inf((t' B_1 t) / \|t\|^2; t \neq 0) \geq b > 0$$

for all large  $N$  under A. Finally for any  $|\beta| \leq 3$

$$|D^{\beta}(g_i^{\beta}(t))| \leq n^{|\beta|} E \|Z_i n^{-1/2}\|^{|\beta|} |g_i(t)|^{n-1} \leq c N^2 |g_i(t)|^{n-2}.$$

Thus for  $|\beta| \leq 3$ , the integral in the lemma over  $\lambda^{-1}\alpha^{-\sqrt{N}} \leq \|t\| \leq \eta\sqrt{N}$  is  $O(N^{-1})$ . The claim now follows by letting  $\alpha \rightarrow \infty$ .

Next, an inversion theorem is obtained by combining a modification of Lemma 5 in [9] with Lemma 11.6 of [3]. The proof is deleted.

**LEMMA 4.** Let  $P$  be a probability on  $R^k$  and  $Q$  denote a measure with density  $\{1 + N^{-1/2}p(y)\}\phi_{\Sigma}(y)$  where  $p(y)$  is a polynomial and  $\Sigma$  is a positive definite matrix of order  $k \times k$ . Let the coefficients of  $p(y)$ ,  $\lambda_{\max}$  and  $\lambda_{\min}$  be bounded by  $M > 0$  where  $\lambda_{\max}$  and  $\lambda_{\min}$  denote the maximum and minimum eigen values of  $\Sigma$ . Then for any  $\varepsilon > 0$

$$\begin{aligned} |P(C) - Q(C)| & \leq c(k) \max_{|t| \leq k+1} \int_{\|t\| \leq \varepsilon^{-1}\sqrt{N}} |D^{\beta}(\hat{P}(t) - \hat{Q}(t))| dt \dots \\ & + c(M)[\Phi_{\Sigma}(\partial C)^{1/\sqrt{N}}] + O(N^{-1}). \end{aligned}$$

Here  $\hat{P}$  and  $\hat{Q}$  stand for c.f.'s of  $P$  and  $Q$ ;  $(\partial C)^{1/\sqrt{N}}$  is the  $\varepsilon/\sqrt{N}$  neighborhood of the boundary of  $C$ .

Finally Lemma 5 justifies converting a multivariate one-term Edgeworth expansion into an univariate one. This result is a modification of Lemma 2.1 of [2]. A proof for the present version is contained in [1].

**LEMMA 5.** Let  $t = (t_1, t_2, \dots, t_r)$  be a vector,  $L = \{L_{ij}\}$  be a  $r \times r$  matrix and  $q$  be a polynomial in  $r$  variables. Let  $M \geq \max(|v_{ij}|, |u_{ij}|, |t_i|, |L_{ij}|, |c_i|)$ , where  $V = ((v_{ij}))$  is a positive definite matrix  $(u_{ij}) = V^{-1}$  and  $c_i$  are the coefficients of  $q$ . Let  $|t_i| > t_0 > 0$ . Then there exists a polynomial  $p$  in one variable, whose coefficients are continuous functions of  $t_i, L_{ij}, v_{ij}, u_{ij}$  and  $c_i$ , such that

$$\int_{\|z\| \leq N^{-1/2}q(z)} (1 + N^{-1/2}q(z))\phi_V(z) dz = \int_{-\infty}^{\infty} (1 + N^{-1/2}p(y))\phi(y) dy + o(N^{-1/2})$$

where the  $o(\cdot)$  term depends on  $M$  and  $t_0$  only.

We now briefly sketch the proof of (5) using the lemmas. Define,  $\xi_j = n_j^{-1/2}\Sigma_j^{-1}(Y_j - \bar{X}_j)^2 - s_j^2$  and  $s^2 = \sum_1^r s_j^2/n_j$ . From Lemmas 3 and 4 it follows that for a measurable  $C$  and

$\epsilon > 0$ ,

$$(11) \quad P\{(\sqrt{N} \sum_1^k (\bar{Y}_j - \bar{X}_j), N^{1/2} \sum_1^k (\xi_j/n_j)) \in C\} \\ = \int_C \phi_B(x) [1 + N^{-1/2} \alpha_N(x)] dx + o(N^{-1/2}) + O(\Phi_B(\beta B)^{1/\sqrt{N}})$$

where  $\alpha_N$  is a polynomial whose coefficients are polynomials in  $\{\lambda_j\}$ , and the moments of  $G_j$  of order 6 or less. Note that  $B = \{b_{ij}\}_{2 \times 2}$  with  $b_{11} = N\sigma^2$ , the variance of  $\sqrt{N} \sum_1^k \bar{Y}_j$ . Now (11) combined with Lemma 5 entails

$$(12) \quad P\{s^{-1} \sum_1^k (\bar{Y}_j - \bar{X}_j) [1 - (1/4)s^{-2} \sum_1^k (\xi_j/n_j)] \leq x\} \\ = \Phi(x) + N^{-1/2} \int_{-\infty}^x b(y) \phi(y) dy + o(N^{-1/2}),$$

where  $b$  is a polynomial whose coefficients are continuous functions of  $B^{-1}$ ,  $\lambda_j$  and the moments of  $G_j$  of order 6 or less.

Define  $C_N = \{\sqrt{N} \sum_1^k |\bar{Y}_j - \bar{X}_j| < \log N\}$  and  $D_N = \{N^{1/2} \sum_1^k (\xi_j/n_j) < \log N\}$ . On  $C_N \cap D_N$  one has

$$(13) \quad t_N^2 = s^{-1} \sum_1^k (\bar{Y}_j - \bar{X}_j) [1 - (1/4)s^{-2} \sum_1^k (\xi_j/n_j)] + O(N^{-1}(\log N)^2)$$

(taking  $l_1 = \dots = l_4 = 1$ ). Since the 6th moments of  $\{Y_{1j}\}$  are bounded, it follows from the proof of Theorem 2 of [7] that

$$(14) \quad [1 - P(C_N)] + [1 - P(D_N)] = o(N^{-1/2}).$$

Thus (12), (13) and (14) yield (5)

## REFERENCES

- [1] BANU, G. J. and SINGH, K. (1981). On one term Edgeworth correction by Efron's bootstrap. Unpublished manuscript.
- [2] BHATTACHARYA, R. N. and GHOSH, J. K. (1978). On the validity of formal Edgeworth expansion. *Ann. Statist.* 6 435-451.
- [3] BHATTACHARYA, R. N. and RANGA RAO, R. (1976). *Normal Approximation and Asymptotic Expansions*. Wiley, New York.
- [4] BICKEL, P. J. and FREEDMAN, D. (1981). Some asymptotics on the bootstrap. *Ann. Statist.* 9 1196-1217.
- [5] EFRON, B. (1979). Bootstrap—another look at Jackknife. *Ann. Statist.* 7 1-26.
- [6] EFRON, B. (1981). Nonparametric standard errors and confidence intervals. Stanford Technical Report No. 67, April 1981.
- [7] MICHEL, R. (1976). Non-uniform central limit bounds with applications to probabilities of deviations. *Ann. Probab.* 4 102-106.
- [8] SINGH, K. (1981). On asymptotic accuracy of Efron's bootstrap. *Ann. Statist.* 9 1187-1195.
- [9] SWEETING, T. J. (1977). Speed of convergence for the multidimensional central limit theorem. *Ann. Probab.* 5 28-41.

MATH. STAT. DIVISION  
INDIAN STATISTICAL INSTITUTE  
203 B. T. ROAD  
CALCUTTA-700035  
INDIA

DEPARTMENT OF STATISTICS  
HILL CENTER FOR THE MATHEMATICAL SCIENCES  
RUTGERS UNIVERSITY  
NEW BRUNSWICK, NEW JERSEY 08903