Departamento de Estadística y Econometría

Universidad Carlos III de Madrid

Calle Madrid, 126

28903 Getafe (Spain)

Fax 34 - 91- 624.9849

# FINITE SAMPLE NONPARAMETRIC INFERENCE AND LARGE SAMPLE EFFICIENCY

Joseph P. Romano and Michael Wolf*

Abstract

Given a sample $X_1,...,X_n$ from a distribution $F$, the problem of constructing nonparametric confidence intervals for the mean $\mu(F)$ is considered. Unlike bootstrap procedures or those based on normal approximations, we insist on any procedure being truly nonparametric in the sense that the probability that the confidence interval contains $\mu(F)$ based on a sample of size $n$ from $F$ be at least $1-\alpha$ for all $F$ and all $n$. Bahadur and Savage (1956) proved it is impossible to find an effective confidence interval for $\mu(F)$ without some restrictions. Thus, we assume that $F$ is supported on a compact set, which we take to be [0,1]. In this setting, and asymptotic efficiency result is obtained that gives a lower bound on the size of any conservative interval. We then provide a construction of an interval that meets our finite sample requirement on level, yet has an asymptotic efficiency property. Thus, the price to be paid for using fully nonparametric procedures when considering the tradeoff between exact inference statements and asymptotic efficiency is negligible. Much of what is accomplished for the mean generalizes to other settings as well.

Key Words

Confidence Intevals; Efficiency; Nonparametrics.

*Department of Statistics, Stanford University, e-mail: romano@stat.stanford.edu; Department of Statistics y Econometrics, Universidad Carlos III de Madrid, e-mail: mwolf@est-econ.uc3m.es.

# 1. Introduction

Suppose $X_1, \ldots, X_n$ are i.i.d. according to a distribution $F$ on the line. Consider the problem of constructing a level $1 - \alpha$ confidence interval for $\mu(F)$, the mean of $F$. The distribution $F$ is assumed to belong to a large class $\mathbf{F}$ of distributions. Clearly, $\mathbf{F}$ must be restricted somewhat since we are assuming $\mu(F)$ exists. If $I_n$ is a random interval (or set), define the coverage level over $\mathbf{F}$ to be

$$inf\{Pr_F\{\mu(F) \in I_n\} : \quad F \in \mathbf{F}\}.$$

In fact, even if we assume $\mathbf{F}$ consists of all distributions $F$ having finite moments of all orders, Bahadur and Savage (1956) proved the negative result that it is impossible to construct an effective confidence interval for $\mu(F)$. In particular, if $I_n$ is a random interval (depending on $X_1, \ldots, X_n$) such that, even for one $F$, the probability under $F$ that $I_n$ is a bounded set is one, then the coverage level over $\mathbf{F}$ is zero!

Nevertheless, the aforementioned result has not deterred the search for valid inference procedures, especially since Efron's (1979) discovery of the bootstrap, as well as methods based on Edgeworth expansions, likelihood, and other resampling refinements. Indeed, there are several methods yielding intervals $I_n$ of nominal level $1 - \alpha$ satisfying, for fixed $F$,

$$|Pr_F\{\mu(F) \in I_n\} - (1 - \alpha)| = O(n^{-p}), \tag{1.1}$$

for some $p > 0$. In fact, $p = 1$ for intervals whose coverage error is of the same order as that provided by the normal approximation, $p = 2$ for second-order accurate intervals such as the bootstrap $t$-interval, and $p$ can even be larger by bootstrap iteration (under assumptions to ensure the validity of Edgeworth expansions); these properties are derived in Hall (1992). Unfortunately, all these intervals have the property that their coverage level over $\mathbf{F}$ is zero.

The technical reason why these methods can misbehave so badly yet still satify (1.1) is that the convergence result in (1.1) holds for each fixed $F$, and is not uniform over $F$. Therefore, a question worthy of investigation is to find methods that are appropriately uniform, and indeed, some results are obtained in Hall and Jing (1995), by imposing restrictions on $\mathbf{F}$. Here, however, we insist on considering only procedures whose coverage level over a large $\mathbf{F}$ is the nominal level. Uniform convergence to the nominal level over $\mathbf{F}$,

while preferable to pointwise convergence for each fixed $F$, is not strong enough for our inferential purposes; we aim for a coverage level $1 - \alpha$ for any finite sample size rather than in the limit only. Because of the Bahadur and Savage result, we do need to make some restriction. The assumption imposed then is that the unknown $F$ has support in a fixed compact set, which we take to be $[0, 1]$; otherwise, $F$ is arbitrary (and need not be continuous, for example).

The main problem considered in this paper is then the following. Let $\mathbf{F_0}$ be the class of all distributions on $[0, 1]$. The goal is to construct a confidence interval $I_n$ for $\mu(F)$ that is effective. First, among confidence intervals $I_n$ whose coverage level over $\mathbf{F_0}$ is $1 - \alpha$, can we find an optimality result giving a lower bound on how large the length of $I_n$ can be? Second, can we actually construct an optimal interval? The answer to the first question is not that surprising, as intervals based on the normal approximation serve as our gold standard. Thus, if $I_n = [L_n, U_n]$ and $D_n = n^{1/2}(U_n - L_n)/2$, then the standard asymptotic interval $\bar{X}_n \pm z_{1-\frac{\alpha}{2}} s_n/n^{1/2}$ (where $z_{1-\frac{\alpha}{2}}$ is the $1 - \frac{\alpha}{2}$ quantile of the standard normal distribution and $s_n$ is the sample standard deviation) satisfies $D_n \to z_{1-\frac{\alpha}{2}} \sigma(F)$ in probability under $F$, where $\sigma^2(F)$ is the variance of $F$. Not surprisingly, in a nonparametric setting, this asymptotic constant $z_{1-\frac{\alpha}{2}} \sigma(F)$ is in some sense the best attainable.

Of course, the normal theory intervals are not fully nonparametric as their coverage level over $\mathbf{F_0}$ is also zero. Such is the case for bootstrap intervals as well, as made clear in Romano (1989). To appreciate why, consider $F_n \in \mathbf{F_0}$ with $F_n$ the distribution that assigns mass $p_n$ to 1 and $1 - p_n$ to 0. Fix any $\theta$ in $(\alpha, 1)$ and choose $p_n$ small enough (but not 0) so that $(1 - p_n)^n \geq \theta$. Then, under $F_n$, a sample of size $n$ will be a sample of all zeroes with probability at least $\theta$. Since the resampled or bootstrap data sets will all be degenerate as well, the resulting bootstrap confidence interval will not include $\mu(F_n) = p_n$. Thus, with probability at least $\theta$, the resulting bootstrap confidence interval will not cover the true mean, and hence the coverage probability is no bigger than $1 - \theta < 1 - \alpha$. In fact, since $\theta$ is arbitrary in the argument, the coverage level over $\mathbf{F_0}$ is zero.

A variety of conservative methods are presented in Bickel (1992), including some justification for Stringer's (1963) widely used proposal. However, none of the methods presented there are both conservative in level and efficient. We believe our proposal is the first one that is both nonparametric (or conservative) and efficient.

The paper is organized as follows. In Section 2, we derive a result which can be viewed

as an asymptotic efficiency result for conservative confidence intervals for the mean. In essence, the result says that we can not do better, in terms of length of the interval, than an interval based on a normal approximation. The novel aspect of the result is that its proof draws upon local asymptotic minimax estimation theory by looking at an appropriate least favorable parametric submodel, but intervals are not compared in the usual minimax way of worst case behavior over shrinking neighborhoods. In Section 3, we review a proposal of Anderson (1967), which leads to genuinely nonparametric (conservative) intervals whose length is of the right order, but the constant is too big. His proposal is generalized, though efficiency is not obtained. The reason for considering Anderson's approach is because it makes use of the fact that we can construct genuinely nonparametric confidence bands for the c.d.f. $F$ by the usual Kolmogorov-Smirnov bands. Anderson's procedure uses these bands in a primary way. The efficient construction we present in Section 4 also makes use of the Kolmogorov-Smirnov bands, but somehow the construction relies on these bands in a more secondary way so that efficiency can ensue. In Section 5, we present some conclusions, variations, generalizations, and directions for future work.

## 2. An Asymptotic Lower Bound

In this section, it is more convenient to index the probability distribution generating the data by the measure $P$ (as opposed to the c.d.f. $F$). Suppose $X_1, \ldots, X_n$ are i.i.d. according to a probability $P$, concentrated on [0,1]; otherwise, nothing else is assumed about $P$. Inference focuses on the mean $\mu(P)$. The question considered in this section is the following. Among confidence intervals for the mean with guaranteed coverage, what is an asymptotic lower bound for the length of the interval? The following theorem answers the question. In fact, much more is shown. Specifically, the interval must be centered (to an appropriate order) at the sample mean, in order to be efficient. Moreover, an efficient interval, to order $o_P(n^{-1/2})$ behaves like an interval provided by the usual normal approximation, namely $\bar{X}_n \pm \sigma(P) z_{1-\frac{\alpha}{2}}/n^{1/2}$. However, the normal approximation interval is not conservative, and so can not qualify to be efficient. In the statement of the theorem and in the proof, $P^n$ refers to the product measure of $n$ i.i.d. observations from $P$.

**Theorem 2.1.** Let $I_n = [L_n, U_n]$ be a sequence of intervals (being measurable functions of $X_1, \ldots, X_n$) satisfying

$$P^n\{\mu(P) \in I_n\} \geq 1 - \alpha$$

for all $n$ and all $P$. Define $D_n = D_n(I_n) = n^{1/2}(U_n - L_n)/2$. Assume $D_n$ is asymptotically concentrated on $[0, a]$; that is, for every $\epsilon > 0$, $P^n\{D_n < a + \epsilon\} \to 1$ as $n \to \infty$. Then,

(i) $a \geq z_{1-\frac{\alpha}{2}}\sigma(P)$.

(ii) If $I_n$ is an interval such that the lower bound $z_{1-\frac{\alpha}{2}}\sigma(P)$ is attained, then $I_n$ is centered at $\bar{X}_n$ in the sense

$$n^{1/2}[\frac{(U_n + L_n)}{2} - \bar{X}_n] \to 0$$

in $P^n$-probability.

(iii) If $I_n$ is an interval such that the lower bound $z_{1-\frac{\alpha}{2}}\sigma(P)$ is attained, then $D_n \to z_{1-\frac{\alpha}{2}}\sigma(P)$ in $P^n$-probability, and so

$$I_n = \bar{X}_n \pm \frac{\sigma(P)z_{1-\frac{\alpha}{2}}}{n^{1/2}} + o_P(n^{-1/2}).$$

**Proof.** Fix $P = P_0$ having mean $\mu_0 = \mu(P_0)$, standard deviation $\sigma_0 = \sigma(P_0)$, and density $f$ (with respect to some $\sigma$-finite measure $dx$). Introduce the parametric submodel $\{P_{n,\theta}^n\}$,

where $P_{n,\theta}^n$ is the n-fold product measure of $P_{n,\theta}$, and $P_{n,\theta}$ has density

$$f_{n,\theta} = f(x)[1 + \frac{\theta(x - \mu_0)}{n^{1/2}}]$$

for $0 \leq x \leq 1$. Note that for any fixed real number $\theta$, $f_{n,\theta}$ defines a density as soon as $n \geq \theta^2$. Further note that $f_{n,\theta}$ has mean

$$\mu(P_{n,\theta}) = \mu_0 + \theta n^{-1/2} \int_0^1 x(x - \mu_0)f(x)dx = \mu_0 + \frac{\theta\sigma_0^2}{n^{1/2}}.$$

By the usual examination of loglikelihood ratios (see Proposition 2.3 of Millar, 1983), it is immediate that the experiments $\{P_{n,\theta}^n\}$ converge to $\{P_\theta\}$, where $P_\theta$ is the Gaussian measure on the line with mean $\theta\sigma_0$ and variance one. Indeed, the loglikelihood ratio at $\theta$ versus $\theta = 0$ is simply

$$log[\prod_{i=1}^n f_{n,\theta}(X_i) / \prod_{i=1}^n f_{n,0}(X_i)] = \sum_{i=1}^n log[1 + \frac{\theta(X_i - \mu_0)}{n^{1/2}}].$$

Under $P_0^n$ (by a Taylor expansion and the law of large numbers), this loglikelihood ratio is asymptotic to $\theta n^{1/2}(\bar{X}_n - \mu_0) - \theta^2\sigma_0^2/2$, which of course is asymptotically normal with mean $-\theta^2\sigma_0^2/2$ and variance $\theta^2\sigma_0^2$ (implying contiguity). Moreoever, for this parametric submodel, it is then clear that the sample mean $\bar{X}_n$ is then locally asymptotically minimax for the mean. Now, let $\ell$ be any bounded subconvex loss function. Then, by the asymptotic minimax theorem (see Millar, 1983, p.146),

$$lim_{c \uparrow \infty} lim_{n \to \infty} inf_{T_n} sup_{|\theta| \leq c} \int \ell[n^{1/2}(T_n - \mu(P_{n,\theta}))]dP_{n,\theta}^n$$

$$= lim_{c \uparrow \infty} lim_{n \to \infty} inf_{T_n} sup_{|\theta| \leq c} \int \ell[\sigma_0^2(T_n - \theta)]dP_{n,\theta}^n = inf_T sup_\theta \int \ell[\sigma_0^2(T - \theta)]dP_\theta$$

where the infimum over $T_n$ and $T$ refer to the infimum over all estimators available for experiments $P_{n,\theta}^n$ and $P_\theta$, respectively. But, if $Q_\theta$ denotes the Gaussian distribution with mean $\theta$ and variance one, the last expression can be evaluated as

$$inf_T sup_\theta \int \ell[\sigma_0(\sigma_0 T - \sigma_0\theta)]dP_\theta = inf_T sup_\theta \int \ell[\sigma_0(T - \theta)]dQ_\theta = E\ell(\sigma_0 Z),$$

where $Z$ is a standard normal variable. In particular, let $\ell = \ell_d$ be the (subconvex) loss function satisfying $\ell_d(T - \theta) = 1$ if $|T - \theta| > d$ and is zero otherwise. The previous result then says

$$lim_{c \uparrow \infty} lim_{n \to \infty} inf_{T_n} sup_{|\theta| \leq c} P_{n,\theta}^n\{n^{1/2}|T_n - \mu(P_{\theta,n})| > d\} = P\{|Z| \geq \frac{d}{\sigma_0}\} \qquad (2.1)$$

for any $d$. Now, given the interval $I_n = [L_n, U_n]$, let $\hat{\mu}_n$ be the estimator $(L_n + U_n)/2$. Then, by the hypothesis on $I_n$, we get

$$P^n\{n^{1/2}|\hat{\mu}_n - \mu(P)| \leq D_n\} \geq 1 - \alpha.$$

We now claim that, for any $P$ and any $\epsilon > 0$,

$$P^n\{n^{1/2}|\hat{\mu}_n - \mu(P)| \leq a + \epsilon\} \geq (1 - \alpha) \cdot P^n\{D_n \leq a + \epsilon\}. \tag{2.2}$$

To prove (2.2),

$$P^n\{n^{1/2}|\hat{\mu}_n - \mu(P)| \leq a + \epsilon\} \geq P^n\{n^{1/2}|\hat{\mu}_n - \mu(P)| \leq a + \epsilon | D_n \leq a + \epsilon\} P^n\{D_n \leq a + \epsilon\}$$

$$\geq P^n\{n^{1/2}|\hat{\mu}_n - \mu(P)| \leq D_n\} P^n\{D_n \leq a + \epsilon\} \geq (1 - \alpha) \cdot P^n\{D_n \leq a + \epsilon\},$$

and (2.2) follows. Now, under $P_0$, $P_0^n\{D_n \leq a + \epsilon\} \to 1$, by assumption. But, for any sequence $\{\theta_n\}$ with $|\theta_n| \leq c$, $\{P_{n,\theta_n}^n\}$ is contiguous to $\{P_0^n\}$. It then follows that $P_{n,\theta_n}^n\{D_n \leq a + \epsilon\} \to 1$ as well. Hence, by (2.2),

$$\liminf_{n \to \infty} P_{n,\theta_n}^n\{n^{1/2}|\hat{\mu}_n - \mu(P_{n,\theta_n})| \leq a + \epsilon\} \geq 1 - \alpha.$$

Therefore (by a subsequence argument), for any $\epsilon > 0$ and any fixed $c$, we have

$$\limsup_{n \to \infty} \sup_{|\theta| \leq c} P_{n,\theta}^n\{n^{1/2}|\hat{\mu}_n - \mu(P_{n,\theta})| > a + \epsilon\} \leq \alpha. \tag{2.3}$$

Now, to prove the result $a \geq z_{1-\frac{\alpha}{2}}\sigma_0$, assume the opposite and let $\Delta = z_{1-\frac{\alpha}{2}}\sigma_0 - a$. Choose $\epsilon$ so that $(\Delta - \epsilon)/\sigma_0 \equiv y > 0$. Then, choose $\delta$ small enough so that

$$P\{|Z| > z_{1-\frac{\alpha}{2}} - y\} - \delta > \alpha.$$

Finally, choose $c$ large enough so that the result (2.1) with $d = a + \epsilon$ gives

$$\lim_{n \to \infty} \inf_{T_n} \sup_{|\theta| \leq c} P_{n,\theta}^n\{n^{1/2}|T_n - \mu(P_{n,\theta})| > a + \epsilon\} \geq P\{|Z| \geq \frac{a + \epsilon}{\sigma_0}\} - \delta. \tag{2.4}$$

But the right hand side of (2.4) is greater than or equal to

$$P\{|Z| \geq \frac{z_{1-\frac{\alpha}{2}}\sigma_0 - \Delta + \epsilon}{\sigma_0}\} - \delta = P\{|Z| \geq z_{1-\frac{\alpha}{2}} - y\} - \delta > \alpha,$$

by choice of $\delta$. This yields a contradiction because by (2.3), the estimator $\hat{\mu}_n$ would have a smaller asymptotic minimax risk than the bound of $\alpha$ which is valid for any estimator $T_n$ given by (2.4).

To prove (ii), the result (2.3) implies that $\hat{\mu}_n$ is almost locally asymptotic minimax, which would be immediate if (2.3) were true with $\epsilon$ replaced by 0. The result would then follow by Theorem 4.1 of Hajek (1972). Since we can not just replace $\epsilon$ by 0 in (2.3), we sketch an argument based on the proof of (4.3) in Hajek (1972). For notational simplicity, assume $\sigma_0 = 1$. Assume $n^{1/2}(\hat{\mu}_n - \bar{X}_n)$ fails to converge to 0 in $P$-probability. Then, there exists an $\bar{\epsilon} > 0$ such that for all large $m$,

$$P_m^m(m^{1/2}|\hat{\mu}_m - \bar{X}_m| > \bar{\epsilon}) > \bar{\epsilon}.$$

If $Z$ denotes an observation from the model $\{P_\theta\}$, then by Lemmas 3.1 and 3.2 of Hajek (1972), this would imply the existence of an estimator sequence $\xi_m(Z, U)$ (possibly depending on an auxiliary randomization $U$, independent of $Z$), satisfying

$$P_m^m\{|\xi_m(Z, U) - Z| > \bar{\epsilon}\} > \bar{\epsilon}$$

for all $m$ large enough, and $\xi_m(Z, U) - \theta$ is constructed to have the same asymptotic properties as $m^{1/2}(\hat{\mu}_m - \mu(P_{\theta,m}))$; see (4.8) of Hajek. But by Lemma 2.1 of Hajek, there exists $\beta > 0$ depending only on $\bar{\epsilon}$ so that the maximum risk of $\xi_m(Z, U)$ for the loss function $\ell_{z_{1-\frac{\alpha}{2}}+\epsilon}$ is at least $\alpha + \beta$. But, the maximum risk of $\xi_m(Z, U)$ for the loss function $\ell_{z_{1-\frac{\alpha}{2}}+\epsilon}$ is asymptotically (by our (2.3)) less than or equal to $\alpha$, which is within $(2/\pi)^{1/2}\epsilon$ of the best obtainable rate, since

$$|\alpha - P\{|Z| \geq z_{1-\frac{\alpha}{2}} + \epsilon\}| \leq (2/\pi)^{1/2}\epsilon.$$

Hence, by choosing $\epsilon$ small enough so that $(2/\pi)^{1/2}\epsilon < \beta$, we get a contradiction.

Finally, to prove (iii), writing $\mu$ for $\mu(P)$,

$$1 - \alpha \geq P\{L_n \leq \mu \leq U_n\} = P\{n^{1/2}(\bar{X}_n - U_n) \leq n^{1/2}(\bar{X}_n - \mu) \leq n^{1/2}(\bar{X}_n - L_n)\}. \quad (2.5)$$

But, by result (ii), $n^{1/2}(\bar{X}_n - L_n) = D_n + o_P(1)$. Therefore, invoking the asymptotic normality of $n^{1/2}(\bar{X}_n - \mu)$, the probability on the right side of (2.5) is

$$P\{-D_n \leq n^{1/2}(\bar{X}_n - \mu) \leq D_n\} + o_P(1).$$

But, the fact that $D_n$ is asymptotically concentrated on $[0, z_{1-\frac{\alpha}{2}} \sigma(P)]$ and the asymptotic normality of $n^{1/2}(\bar{X}_n - \mu)$ with asymptotic variance $\sigma^2(P)$ forces $D_n \to z_{1-\frac{\alpha}{2}} \sigma(P)$ in $P^n$-probability.

**Remark 2.1.** If a conservative interval $I_n$ is constructed such that the lower bound $z_{1-\frac{\alpha}{2}} \sigma(P)$ is obtained for every fixed $P$, then we will call $I_n$ efficient; otherwise, we will call the interval inefficient.

**Remark 2.2.** Surprisingly, there is not a vast literature on efficiency theory in the construction of confidence intervals. A noteable exception are Beran and Millar (1985), who develop an asymptotic theory for confidence sets in a decision theoretic local asymptotic minimax framework. More specifically, a loss function is introduced to measure the performance of an interval $I_n$ for a real-valued parameter $\theta$ (and more generally for abstract parameter sets) as follows. If $\mu$ is the true parameter and interval $I_n$ is used, the loss is $g(n^{1/2} sup_{y \in I_n} |y - \mu|)$; here, $g$ is an increasing function on the positive reals. Evidently, the loss penalizes if the interval is too wide or if it is miscentered. Now, once a loss function is introduced, different procedures (satisfying the constraint on level) can be compared by comparing their risk functions, by looking at the maximum risk over shrinking neighborhoods as in the local asymptotic minimax theory for estimation problems. Note that our simple result (while indeed making use of local asymptotic minimax ideas) does not compare procedures by looking at the maximum risk over shrinking neighborhoods. The assumption that $D_n$ is asymptotically concentrated on $[0, a]$ need only hold under a single law $P$, and then the conclusion is the asymptotic constant $a$ must then be no less than $z_{1-\frac{\alpha}{2}} \sigma(P)$, under $P$. Issues like superefficiency do not come into play once we impose the coverage constraint, which must hold for all $P$ for $I_n$ to even be considered. Moreover, the optimality result we achieve is based directly on the length of the interval. In essence, however, our result is an asymptotic admissibility and asymptotic minimax result of sorts. Indeed, one can construct an interval $I_n$ whose standardized half-length $D_n$ is strictly less than the bound $z_{1-\frac{\alpha}{2}} \sigma(P)$ with probability tending to some number $p \in (0, 1)$. The number $p$ can be positive but it cannot be one (by the theorem), so that shorter intervals are possible but only at the expense of intervals that must then be bigger some of the time. That is, our result can be viewed as a minimax result because it says that no interval (sequence) can have a half-length that is asymptotically concentrated on a *smaller* set. Furthermore, if an interval achieves the lower bound on the length in the sense of part (i)

of Theorem 2.1, then it *already* must be appropriately centered in the sense of (ii) of the theorem.

**Remark 2.3.** The theorem generalizes to parameters $\theta(P)$ other than the mean. If the functional $\theta(\cdot)$ is appropriately differentiable, then, in a nonparametric setting, efficient intervals must behave like

$$\theta(\hat{P}_n) \pm z_{1-\frac{\alpha}{2}} \tau(P)/n^{1/2} + o_P(n^{-1/2}),$$

where $\hat{P}_n$ is the empirical measure and $\tau^2(P)$ is the asymptotic variance of $n^{1/2}\theta(\hat{P}_n)$ under $P$.

# 3. Inefficient Methods with Guaranteed Coverage

The possibility of finding confidence intervals for the mean with guaranteed coverage, but which are not too big, seems plausible given the following construction, due to Anderson (1967) (and later rediscovered by Breth, Maritz and Williams, 1978). At this point, we switch notation and index the probability distribution generating the sample $X_1, \ldots, X_n$ by the cumulative distribution function (c.d.f.) $F$. Then, $\mu(F) \equiv E_F(X_i)$. Here, all we are assuming is that $F \in \mathbf{F_0}$, the class of all c.d.f.s supported on $[0,1]$. Let $\hat{F}_n$ be the empirical c.d.f.. For c.d.f.s $F$ and $G$, let the sup (Kolmogorov) distance $d_K$ be defined by

$$d_K(F, G) = sup_t|F(t) - G(t)|.$$

Let $\hat{R}_{n,1-\alpha}$ be the Kolmogorov-Smirnov uniform confidence band for $F$ of nominal level $1 - \alpha$ defined by

$$\hat{R}_{n,1-\alpha} = \{F \in \mathbf{F_0} : n^{1/2}d_K(\hat{F}_n, F) \le c_n(1 - \alpha)\}, \tag{3.1}$$

where $c_n(1 - \alpha)$ is the $1 - \alpha$ quantile of the distribution of $n^{1/2}d_K(\hat{F}_n, F)$ under $F$ when $F$ is any continuous distribution. Note that, for any $F$ (discrete or otherwise),

$$P_F\{F \in \hat{R}_{n,1-\alpha}\} \ge 1 - \alpha;$$

the inequality is an equality iff $F$ is continuous. This leads to a nominal level $1 - \alpha$ confidence interval $I_{n,0}$ for $\mu(F)$ as follows. In words, the value $\mu$ is included in $I_{n,0}$ if there is some distribution $F$ in $\hat{R}_{n,1-\alpha}$ that has mean $\mu$. Then, the event $\{F \in \hat{R}_{n,1-\alpha}\}$ implies $\{\mu(F) \in I_{n,0}\}$ and so

$$P_F\{\mu(F) \in I_{n,0}\} \ge 1 - \alpha.$$

In fact, $I_{n,0} = \bar{X}_n \pm O_P(n^{-1/2})$, which follows from the following simple proposition.

**Proposition 3.1.** Let $m_k(F) = E_F(X_i^k)$. If $F$ and $G$ are in $\mathbf{F_0}$ and $d_K(F, G) \le \epsilon$, then $|m_k(F) - m_k(G)| \le \epsilon$.

**Proof.** By integration by parts,

$$|m_k(F) - m_k(G)| = |\int_0^1 x^k d(F - G)(x)| = |\int_0^1 (F(x) - G(x))kx^{k-1}dx|$$

$$\leq \epsilon k \int_0^1 x^{k-1} dx = \epsilon.$$

It follows immediately from Proposition 3.1 that the interval $I_{n,0}$ is contained in the interval $\bar{X}_n \pm c_n(1-\alpha)/n^{1/2}$. The claim that $I_{n,0} = \bar{X}_n \pm O_P(n^{-1/2})$ follows since $c_n(1-\alpha) \to c(1-\alpha)$, where $c(1-\alpha)$ is the upper $1-\alpha$ quantile of the distribution of $sup_{0<t<1}|B(t)|$, where $B(\cdot)$ is a Brownian Bridge process.

In fact, it can be argued that

$$I_{n,0} = \bar{X}_n \pm c(1-\alpha)/n^{1/2} + o_P(n^{-1/2}).$$

The constant $c(1-\alpha)$ is too big and should be compared with $z_{1-\frac{\alpha}{2}}\sigma(F)$. For example, with $\alpha = .05$, $c(.95) \doteq 1.36$, while $z_{1-\frac{\alpha}{2}}\sigma(F) \doteq 1.96\sigma(F) \leq .98$ for all $F$. In particular, when $F$ is the uniform distribution on $(0,1)$, the asymptotically best constant is $1.96/(12)^{1/2} \doteq .57$. The ratio $1.36/.57 \doteq 2.4$ measures the inefficiency of Anderson's procedure. Indeed, a sample of approximately $(2.4)^2 n = 5.76n$ is needed when using Anderson's procedure to be as efficient as an efficient procedure based on a sample of size $n$.

Anderson's procedure, with the hope of improved efficiency, is generalized in Gasko (1991)▮ (who treats the one-sided case) as follows. If $\underline{x}_n = (X_1, \ldots, X_n)$ is a sample of size $n$ from $F$, let $X_{i:n}$ denote the $i$-th order statistic. Define $\pi_i = \pi_i(\underline{x}_n, F) = F(X_{i+1:n}) - F(X_{i:n})$ for $i = 1, \ldots, n-1$; also, let $\pi_0 = F(X_{1:n})$ and $\pi_n = 1 - F(X_{n:n})$. Then, $\underline{\pi} = \underline{\pi}(\underline{x}_n, F)$ is a random point in the $n+1$ simplex $\Sigma_n$ in $\mathbf{R}^{n+1}$. Note that the distribution of $\underline{\pi}(\underline{x}_n, F)$ under $F$ is the same for all continuous $F$.

Now, let $S_{n,1-\alpha} \subset \Sigma_n$ be a region of probability content at least $1-\alpha$, meaning

$$P_F\{\underline{\pi}(\underline{x}_n, F) \in S_{n,1-\alpha}\} \geq 1 - \alpha \tag{3.2}$$

for all $F$. Then, the region $S_{n,1-\alpha}$ induces a confidence set $C_{n,1-\alpha}$ for $F$ by including $F$ in $C_{n,1-\alpha}$ if $\underline{\pi}(\underline{x}_n, F) \in S_{n,1-\alpha}$. Clearly,

$$P_F\{F \in C_{n,1-\alpha}\} \geq 1 - \alpha,$$

and so $C_{n,1-\alpha}$ is a level $1-\alpha$ confidence set for $F$.

We can now proceed as before to construct an interval $I_n$ for the mean. That is, a value $\mu$ belongs in the interval $I_n$ if there is some distribution $F \in C_{n,1-\alpha}$ with mean

$\mu$. Then, $P_F\{\mu(F) \in I_n\} \geq 1 - \alpha$. However, regardless of the choice of region $S_{n,\alpha}$ the resulting interval $I_n$ is not efficient.

**Proposition 3.2.** Any region $S_{n,1-\alpha}$ satisfying (3.2) leads to a conservative $1 - \alpha$ confidence interval $I_n$ for $\mu(F)$. However, no choice of $S_{n,1-\alpha}$ leads to efficiency.

**Proof.** Fix (a sequence of regions) $S_{n,1-\alpha}$ satisfying (3.2), and let $I_n$ denote the induced confidence interval for $\mu(F)$. Assume efficiency holds; then, by Theorem 2.1 (iii), we must have

$$I_n = \bar{X}_n \pm \frac{z_{1-\frac{\alpha}{2}} \cdot \sigma(F)}{n^{1/2}} + O_P(n^{-1/2})$$

for all $F$. Now, $S_{n,1-\alpha}$ also induced a level $1-\alpha$ confidence interval, $M_n$, for the parameter $m(F) = E_F(X_i^2)$ by a similar prescription: a value $m$ is included in $M_n$ if there is an $F$ in $C_{n,1-\alpha}$ with $m(F) = m$. By construction, we in fact have that $(\mu(F), m(F)) \in (I_n, M_n)$ with probability at least $1 - \alpha$ for all $F$.

We now argue that the interval for $m(F)$ based on data $\underline{x}_n$ is the same as the interval for $\mu(F)$ based on data $\underline{y}_n$, where $\underline{y}_n = (Y_1, \ldots, Y_n)$ and $Y_i = X_i^2$. To appreciate why this is so, if the interval $M_n$ for $m(F)$ based on $\underline{x}_n$ includes a value $m$, then there is an $F$ in $C_{n,1-\alpha}$ with $m(F) = m$ and $\underline{\pi}(\underline{x}_n, F) \in S_{n,1-\alpha}$. But then, $\underline{\pi}(\underline{y}_n, \tilde{F}) \in S_{n,1-\alpha}$, where $\tilde{F}(\cdot) = F(\cdot^{1/2})$ is the distribution of $Y_i = X_i^2$ if $X_i$ has c.d.f. $F$. This follows because

$$F(X_{i:n}) = F(Y_{i:n}^{1/2}) = \tilde{F}(Y_{i:n}).$$

Hence, based on data $\underline{y}_n$, the value

$$\int y d\tilde{F}(y) = \int x^2 dF(x) = m$$

is included in the confidence interval for the mean. Conversely, if the interval for the mean based on $\underline{y}_n$ includes $m$, so does the interval for the second moment based on $\underline{x}_n$. Therefore, the hypothesis that $I_n$ is efficient entails $M_n$ is efficient and so

$$M_n = \bar{Y}_n \pm \frac{z_{1-\frac{\alpha}{2}} \sigma(\tilde{F})}{n^{1/2}} + o_P(n^{-1/2}).$$

We are now in a position to arrive at a contradiction. By the bivariate Central Limit Theorem,

$$1 - \alpha \leq P_F\{\mu(F) \in I_n, \ m(F) \in M_n\}$$

$$= P\{n^{1/2}|\bar{X}_n - \mu(F)| \le z_{1-\frac{\alpha}{2}}\sigma(F), \ n^{1/2}|n^{-1}\sum_{i=1}^{n}X_i^2 - m(F)| \le z_{1-\frac{\alpha}{2}}\sigma(\tilde{F})\} + o_P(1).$$

Now, the last expression converges to

$$P\{|Z_1| \le z_{1-\frac{\alpha}{2}}\sigma(F), \ |Z_2| \le z_{1-\frac{\alpha}{2}}\sigma(\tilde{F})\},$$

where $(Z_1, Z_2)$ is bivariate normal, mean 0, $Var(Z_1) = \sigma^2(F)$, $Var(Z_2) = \sigma^2(\tilde{F})$, and covariance $Cov(X_1, X_1^2)$. Hence,

$$1 - \alpha \le P\{|Z_1| \le z_{1-\frac{\alpha}{2}}\sigma(F), \ |Z_2| \le z_{1-\frac{\alpha}{2}}\sigma(\tilde{F})\}$$

$$= P\{|Z_1| \le z_{1-\frac{\alpha}{2}}\sigma(F)\} - P\{|Z_1| \le z_{1-\frac{\alpha}{2}}\sigma(F), \ |Z_2| > z_{1-\frac{\alpha}{2}}\sigma(\tilde{F})\}$$

$$= (1 - \alpha) - P\{|Z_1| \le z_{1-\frac{\alpha}{2}}\sigma(F), \ |Z_2| \ge z_{1-\frac{\alpha}{2}}\sigma(\tilde{F})\}.$$

This last probability can only be 0 if $Z_1$ and $Z_2$ are perfectly correlated, which means that $X_1$ and $X_1^2$ must be linearly dependent with probability one. Clearly, this is not the case when $X_1$ is uniform on $(0,1)$, and so a contradiction is obtained.

**Remark 3.1.** An alternative way to generalize Anderson's idea is to just change the metric. The same basic argument in the proof of Proposition 3.2 would show inefficiency.

**Remark 3.2.** While no particular region $S_{n,1-\alpha}$ leads to efficiency, Gasko (1991) argues that improvements in efficiency at a particular $F_0 \in \mathbf{F_0}$ are possible, and that perhaps $F_0$ can be chosen adaptively. Some simulations show some improvement over Anderson's method, but no proof of efficiency is supplied (and it seems doubtful that one exists).

# 4. An Efficient Construction for the Mean.

Let $X_1, \ldots, X_n$ be i.i.d. with c.d.f. $F$, mean $\mu(F)$, variance $\sigma^2(F)$, and set $\rho(F) \equiv E_F[|X_i - \mu(F)|^3]$. The unknown c.d.f. $F$ is assumed to be in $\mathbf{F_0}$. Again, the goal is to construct a confidence $I_n$ for $\mu(F)$ that contains $\mu(F)$ with probability at least $1 - \alpha$ for every $F$; in addition, the interval $I_n$ must be efficient. In particular, the interval must satisfy the square root of the sample size multiplied by the length of the interval tends in probability to $2z_{1-\frac{\alpha}{2}}\sigma(F)$, for every $F \in \mathbf{F_0}$.

Let $\bar{X}_n = n^{-1}\sum_{i=1}^{n} X_i$ and $J_n(F)$ be the distribution of $n^{1/2}(\bar{X}_n - \mu(F))$ under $F$, with corresponding c.d.f.

$$J_n(x, F) = P_F\{n^{1/2}(\bar{X}_n - \mu(F)) \le x\}.$$

Also, let

$$d_n(\alpha, F) \equiv \inf\{x: \ J_n(x, F) \ge \alpha\}.$$

Note that $J_n(d_n(\alpha, F), F) \ge \alpha$ and $J_n(d_n(\alpha, F)^-, F) \le \alpha$. Suppose $\beta_n$ is a sequence of numbers in $[0, 1]$ converging to 0. Let $\hat{R}_{1, 1-\beta_n}$ be defined by (3.1) with $\alpha$ there replaced by $\beta_n$. Define

$$\hat{d}_{n,U}(1 - \frac{\alpha}{2}) \equiv \sup\{d_n(1 - \frac{\alpha}{2} + \beta_n, F): \ F \in \hat{R}_{n, 1-\beta_n}\} \tag{4.1}$$

and

$$\hat{d}_{n,L}(\frac{\alpha}{2}) \equiv \inf\{d_n(\frac{\alpha}{2} - \beta_n, F): \ F \in \hat{R}_{n, 1-\beta_n}\}. \tag{4.2}$$

The interval we first propose is defined by

$$I_{n,1} = \{\mu: \hat{d}_{n,L}(\frac{\alpha}{2}) \le n^{1/2}(\bar{X}_n - \mu) \le \hat{d}_{n,U}(1 - \frac{\alpha}{2})\}. \tag{4.3}$$

This construction leads to an interval which, in fact, is determined by two one-sided conservative level $1 - \frac{\alpha}{2}$ intervals; as such, we are tacitly assuming $\alpha \le 0.5$. At this point, it is certainly not clear how to compute $I_{n,1}$ because of the sup in the definition (4.1), which is a sup over an infinite-dimensional set. For now, we postpone this computational issue and present the following result. Note, however, in the course of analyzing $I_{n,1}$, we will derive further more conservative intervals (which are efficient too), but which are directly computable.

**Theorem 4.1.** For each $n$ and all $F \in \mathbf{F_0}$,

$$P_F\{\mu(F) \in I_{n,1}\} \geq 1 - \alpha.$$

**Proof.**

$$P_F\{\mu(F) \in I_{n,1}\} \geq P_F\{\hat{d}_{n,L}(\frac{\alpha}{2}) \leq n^{1/2}(\bar{X}_n - \mu(F)) \leq \hat{d}_{n,U}(1 - \frac{\alpha}{2}) \bigcap F \in \hat{R}_{n,1-\beta_n}\}$$

$$\geq P_F\{d_n(\frac{\alpha}{2} - \beta_n, F) \leq n^{1/2}(\bar{X}_n - \mu(F)) \leq d_n(1 - \frac{\alpha}{2} + \beta_n, F) \bigcap F \in \hat{R}_{n,1-\beta_n}\}$$

$$\geq P_F\{d_n(\frac{\alpha}{2} - \beta_n, F) \leq n^{1/2}(\bar{X}_n - \mu(F)) \leq d_n(1 - \frac{\alpha}{2} + \beta_n, F)\} - P_F\{F \notin \hat{R}_{n,1-\beta_n}\}, \quad (4.4)$$

where we have used the trivial inequality $P(A \cap B) \geq P(A) - P(B^c)$. But, (4.4) is bounded below by

$$(1 - \alpha + \beta_n) - P_F\{F \notin \hat{R}_{n,1-\beta_n}\} \geq 1 - \alpha,$$

and the proof is complete.

**Remark 4.1.** Setting $\beta_n = 0$ in the construction of $I_{n,1}$ would lead to a conservative interval, but not one which is efficient.

**Remark 4.2.** The above construction is related to the constructions presented in Loh (1985) and Silvapulle (1996), both of whom consider parametric testing problems. Note, however, that Loh's construction in essence replaces the value $d_n(1 - \frac{\alpha}{2} + \beta_n, F)$ in the definition (4.1) by $d_n(1 - \frac{\alpha}{2}, F)$, which would not lead to a conservative interval. The goal here is not just to produce conservative intervals, but efficient ones in a nonparametric setting.

**Theorem 4.2.** Suppose $\beta_n > 0$ and $\beta_n \to 0$ in such a way so that $log(\beta_n)/n \to 0$ as $n \to \infty$. Then, $\hat{d}_{n,U}(1 - \frac{\alpha}{2}) \to \sigma(F)z_{1-\frac{\alpha}{2}}$ in probability under $F$; similarly, $\hat{d}_{n,L}(\frac{\alpha}{2}) \to \sigma(F)z_{\frac{\alpha}{2}}$ in probability under $F$. Therefore, $I_{n,1}$ is efficient in the sense of Theorem 2.1.

In order to prove Theorem 4.2, we seek an upper bound for $\hat{d}_{n,U}(1 - \frac{\alpha}{2})$ which is analytically tractable and which tends to $\sigma(F)z_{1-\frac{\alpha}{2}}$ in probability under $F$. At the same time, we constructively derive an upper bound which may be computed explicity.

First, let

$$\Delta_n(F) = C_{BE} n^{-1/2} \rho(F) \sigma^{-3}(F), \tag{4.5}$$

where $C_{BE}$ is the smallest known universal constant valid in the Berry-Esseen Theorem. Then, let $\delta_{n,\alpha}(F)$ be defined as follows:

$$\delta_{n,\alpha}(F) = \inf\{\delta > 0 : \Phi(z_{1-\frac{\alpha}{2}+\beta_n}(1+\delta)) - \Phi(z_{1-\frac{\alpha}{2}+\beta_n}) \geq \Delta_n(F)\}. \tag{4.6}$$

In (4.6), $\Phi(\cdot)$ denotes the standard normal c.d.f..

**Proposition 4.1.** The following bound holds:

$$d_n(1 - \frac{\alpha}{2} + \beta_n, F) \leq z_{1-\frac{\alpha}{2}+\beta_n}(1 + \delta_{n,\alpha}(F)) \cdot \sigma(F). \tag{4.7}$$

Therefore,

$$\hat{d}_{n,U}(1 - \frac{\alpha}{2}) \leq z_{1-\frac{\alpha}{2}+\beta_n} \sup\{(1 + \delta_{n,\alpha}(F)) \cdot \sigma(F) : F \in \hat{R}_{n,1-\beta_n}\} \tag{4.8}$$

and

$$\hat{d}_{n,L}(\frac{\alpha}{2}) \geq z_{\frac{\alpha}{2}-\beta_n} \inf\{(1 + \delta_{n,\alpha}(F)) \cdot \sigma(F) : F \in \hat{R}_{n,1-\beta_n}\}. \tag{4.9}$$

**Proof of Proposition 4.1.** Let $r \equiv r_{n,\alpha}(F) = 1 + \delta_{n,\alpha}(F)$. By the Berry-Esseen Theorem,

$$|J_n(z_{1-\frac{\alpha}{2}+\beta_n}\sigma(F)r, F) - \Phi(z_{1-\frac{\alpha}{2}+\beta_n}r)| \leq \Delta_n.$$

By definition of $\delta_{n,\alpha}(F)$ (and hence $r$),

$$\Phi(z_{1-\frac{\alpha}{2}+\beta_n}r) - (1 - \frac{\alpha}{2} + \beta_n) \geq \Delta_n.$$

Therefore, by the triangle inequality,

$$J_n(z_{1-\frac{\alpha}{2}+\beta_n}\sigma(F)r, F) - (1 - \frac{\alpha}{2} + \beta_n) \geq 0,$$

and (4.7) follows, as then do (4.8) and (4.9).

Now, in order to use the bounds (4.8) and (4.9), we first need to understand the behavior of $\delta_{n,\alpha}(F)$. Indeed, $\delta_{n,\alpha}(F)$ is order $n^{-1/2}$ in probability; in fact, the following is true. Of course, $\phi(\cdot)$ denotes the standard normal density.

**Lemma 4.1.** For all $n$ large enough

$$\delta_{n,\alpha}(F) \leq \Delta_n(F)[\phi^2(z_{1-\frac{\alpha}{2}+\beta_n}) - 2(2\pi e)^{-1/2}\Delta_n(F)]^{-1/2}/z_{1-\frac{\alpha}{2}+\beta_n};$$

the bound holds for $n$ satisfying $\phi^2(z_{1-\frac{\alpha}{2}+\beta_n}) > 2(2\pi e)^{-1/2}\Delta_n(F)$, which holds for all large $n$.

The proof of the lemma will be deferred to the appendix. Now, define $\hat{\Delta}_n$ and $\hat{\delta}_n$ to be

$$\hat{\Delta}_n = sup\{\Delta_n(F): \ F \in \hat{R}_{n,1-\beta_n}\} \tag{4.10}$$

and

$$\hat{\delta}_n = inf\{\delta > 0: \ \Phi(z_{1-\frac{\alpha}{2}+\beta_n}(1+\delta)) - \Phi(z_{1-\frac{\alpha}{2}+\beta_n}) \geq \hat{\Delta}_n\}. \tag{4.11}$$

Also, define

$$\hat{\sigma}_{n,U} = sup\{\sigma(F): \ F \in \hat{R}_{n,1-\beta_n}\}. \tag{4.12}$$

From (4.8) and (4.9), this leads to the following bounds.

**Proposition 4.2.**

$$\hat{d}_{n,U}(1 - \frac{\alpha}{2}) \leq z_{1-\frac{\alpha}{2}+\beta_n}(1+\hat{\delta}_n)\hat{\sigma}_{n,U} \tag{4.13}$$

$$\hat{d}_{n,L}(\frac{\alpha}{2}) \geq z_{\frac{\alpha}{2}-\beta_n}(1+\hat{\delta}_n)\hat{\sigma}_{n,U} \tag{4.14}$$

Note, from Lemma 4.1, we have

$$\hat{\delta}_n \leq \hat{\Delta}_n[\phi^2(z_{1-\frac{\alpha}{2}+\beta_n}) - 2(2\pi e)^{-1/2}\hat{\Delta}_n]^{-1/2}/z_{1-\frac{\alpha}{2}+\beta_n} \tag{4.15}$$

as soon as the term in brackets is nonnegative. So, to bound $\hat{\delta}_n$ in (4.13) and (4.14), it suffices to understand the behavior of $\hat{\Delta}_n$. But, using the (crude) inequality $\rho(F) \leq \sigma^2(F)$ (for $F \in \mathbf{F_0}$), we get

$$\hat{\Delta}_n \leq C_{BE}n^{-1/2}sup\{\sigma^{-1}(F): \ F \in \hat{R}_{n,1-\beta_n}\} = C_{BE}n^{-1/2}\hat{\sigma}_{n,L}^{-1}, \tag{4.16}$$

where $\hat{\sigma}_{n,L}$ is defined by

$$\hat{\sigma}_{n,L} = inf\{\sigma(F): \ F \in \hat{R}_{n,1-\beta_n}\}. \tag{4.17}$$

Therefore, to complete our series of successive approximations, it is only necessary to bound $\sigma(F)$ (in both directions) as $F$ varies in $\hat{R}_{n,1-\beta_n}$. We will appeal to the following lemma.

**Lemma 4.2.** Fix $F$ and $G$, c.d.f.s in $\mathbf{F_0}$. If $d_K(F, G) \le \epsilon$, then

$$|\sigma^2(F) - \sigma^2(G)| \le 3\epsilon.$$

Therefore,

$$|\sigma(F) - \sigma(G)| \le (3\epsilon)^{1/2} \tag{4.18}$$

or

$$|\sigma(F) - \sigma(G)| \le 3\epsilon/\sigma(F). \tag{4.19}$$

**Proof.** Apply Proposition 3.1.

Now, for $F \in \hat{R}_{n,1-\beta_n}$, it follows from (4.19) (noting we could have employed (4.18)) that

$$|\sigma(\hat{F}_n) - \sigma(F)| \le \frac{3c_n(1 - \beta_n)}{n^{1/2}\sigma(\hat{F}_n)} \tag{4.20}$$

and so the following is true.

**Proposition 4.3.**

$$\hat{\sigma}_{n,U} \le \sigma(\hat{F}_n) + \frac{3c_n(1 - \beta_n)}{n^{1/2}\sigma(\hat{F}_n)} \tag{4.21}$$

$$\hat{\sigma}_{n,L} \ge \sigma(\hat{F}_n) - \frac{3c_n(1 - \beta_n)}{n^{1/2}\sigma(\hat{F}_n)}. \tag{4.22}$$

**Theorem 4.3.** Suppose $\beta_n$ satisfies $log(\beta_n)/n \to 0$ as $n \to \infty$. Let

$$I_{n,2} = \bar{X}_n \pm n^{-1/2}z_{1-\frac{\alpha}{2}+\beta_n}(1 + \hat{\delta}_n)\hat{\sigma}_{n,U}. \tag{4.23}$$

Then, $I_{n,2}$ contains $I_{n,1}$ (hence is conservative) and $I_{n,2}$ (hence $I_{n,1}$) is efficient in the sense of Theorem 2.1.

**Proof of Theorems 4.2 and 4.3.** It suffices to show $\hat{\delta}_n \to 0$ in probability under $F$ and $\hat{\sigma}_{n,U} \to \sigma(F)$ in probability under $F$. First, to show $\hat{\sigma}_{n,U} \to \sigma(F)$ in probability, by the law of large numbers and (4.21), it suffices to show $c_n(1 - \beta_n)/n^{1/2} \to 0$. But, by the Dvoretsky-Kiefer-Wolfowitz inequality, there is a universal constant $C_{DKW}$ such that

$$P_F\{n^{1/2}d_K(\hat{F}_n, F) > t\} \leq C_{DKW}exp(-2t^2),$$

so that $c_n(1 - \beta_n)$ satisfies

$$C_{DKW}exp(-2c_n^2(1 - \beta_n)) \geq \beta_n,$$

which implies

$$\frac{c_n^2(1 - \beta_n)}{n} \leq \frac{1}{2n}log(\frac{C_{DKW}}{\beta_n}).$$

The right hand side tends to 0 by our hypothesis on $\beta_n$. Thus, $\hat{\sigma}_{n,U} \to \sigma(F)$ in probability. By the same reasoning, $\hat{\sigma}_{n,L} \to \sigma(F)$ in probability. Therefore, by (4.16), $\hat{\Delta}_n \to 0$ in probability, and by (4.15), $\hat{\delta}_n \to 0$ in probability. The proof is complete.

**Remark 4.3.** Actually, the bounds used in the proof show much more. Specifically,

$$|\hat{\sigma}_{n,U} - \sigma(F)| = O_P(\frac{c_n(1 - \beta_n)}{n^{1/2}}) = O_P(|\frac{log(n)}{n}|^{1/2})$$

if, for example, $\beta_n$ satisfies $\beta_n = n^{-p}$ for any $p > 0$. Furthermore, $|z_{1-\frac{\alpha}{2}+\beta_n} - z_{1-\frac{\alpha}{2}}| = O(\beta_n)$. Also, by (4.16) and (4.20), $\hat{\Delta}_n = O_P[(log(n)/n)^{1/2}]$, and so $\hat{\delta}_n$ is this same order, by Lemma 4.1. Hence, for $j = 1, 2$,

$$I_{n,j} = \bar{X}_n \pm n^{-1/2}(z_{1-\frac{\alpha}{2}} + O_P(\beta_n))(1 + O_P[(log(n)/n)^{1/2}])(\sigma(F) + O_P(|\frac{log(n)}{n}|^{1/2}))$$

or

$$I_{n,j} = \bar{X}_n \pm n^{-1/2}z_{1-\frac{\alpha}{2}}\sigma(F) + O_P(\frac{\beta_n}{n^{1/2}} + \frac{[log(\beta_n)]^{1/2}}{n}). \tag{4.24}$$

Taking $\beta_n = O(n^{-1/2})$ then yields the following corollary.

**Corollary 4.1.** If $\beta_n = O(n^{-1/2})$, then

$$I_{n,j} = \bar{X}_n \pm n^{1/2}z_{1-\frac{\alpha}{2}}\sigma(F) + O_P(\frac{log(n)}{n}). \tag{4.25}$$

Of course, the intervals $I_{n,j}$ are conservative in level, but asymptotically satisfy

$$P_F\{\mu(F) \in I_{n,j}\} \to 1 - \alpha$$

as $n \to \infty$, for any fixed nondegenerate $F$. In fact, the proof shows a much stronger statement. Let $\mathbf{F}_\tau$ be all distributions $F$ on $[0,1]$ with $\sigma(F) \geq \tau$. Then, the arguments in this section show the following.

**Corollary 4.2.** For any $\tau > 0$ and for $j = 1, 2$,

$$sup_{F \in \mathbf{F}_\tau}[P_F\{\mu(F) \in I_{n,j}\} - (1 - \alpha)] \to 0$$

as $n \to \infty$.

**Remark 4.4.** By now, the constructive approach used clearly shows efficient intervals can even be computed by hand. Indeed, by employing even cruder approximations while still retaining efficiency, use Proposition 4.3 to get the following conservative and efficient interval:

$$I_{n,3} = \bar{X}_n \pm n^{-1/2} z_{1-\frac{\alpha}{2}+\beta_n}(\sigma(\hat{F}_n) + \frac{3c_n(1 - \beta_n)}{n^{1/2}\sigma(\hat{F}_n)})(1 + \hat{\delta}_{n,U}),$$

where $\hat{\delta}_{n,U}$ is obtained by replacing $\hat{\Delta}_n$ in (4.11) by

$$\hat{\Delta}_{n,U} = C_{BE}n^{-1/2}[\sigma(\hat{F}_n) - \frac{3c_n(1 - \beta_n)}{n^{1/2}\sigma(\hat{F}_n)}]^{-1} ;$$

tacitly, we are assuming the term in brackets is positive for such an approximation to be employed (which happens with probability tending to one). The reasoning follows from inequalities (4.16) and (4.22).

# 5. Conclusions and Directions for Future Work

In this paper, we have demonstrated that it is indeed possible to construct confidence intervals for the mean that have a finite sample nonparametric validity and large sample efficiency, if it assumed the observations lie in a compact set. Our approach was constructive and made use of some crude inequalities (so there is hopefully room for improvement); nevertheless, we have proved existence of efficient intervals which can even be calculated by hand. Ultimately, we would not want to employ the most conservative of the intervals, $I_{n,3}$, and we seek better ways at approximating $I_{n,1}$.

There are obvious approaches based on simulation, but we have deliberately avoided these in this present work because we wanted to present a construction free of simulation error which might destroy our finite sample conservativeness (unless properly accounted for). Clearly, however, one can always do what amounts to many bootstraps by sampling from distributions in our Kolmogorov band for $F$. Then, $\hat{d}_{n,U}(1 - \frac{\alpha}{2})$ in (4.1) can be approximated by the largest value of $d_n(1 - \frac{\alpha}{2} + \beta_n, F)$ as $F$ ranges over a random selection of $F$ in $\hat{R}_{n,1-\beta_n}$.

Alternatively, we could consider replacing $\hat{d}_{n,U}(1 - \frac{\alpha}{2})$ by the bootstrap quantile $d_n(1 - \frac{\alpha}{2} + \beta_n, \hat{F}_n)$, as long as we can account for the variation of $J_n(F)$ as $F$ varies in $\hat{R}_{n,1-\beta_n}$; indeed, the Berry-Esseen theorem can be utilized, though better approximations are desired.

Our method can be generalized in a number of directions. First, the Kolmogorov Smirnov distance played a role but only a minor one; however, it proves to be convenient at this time. Second, instead of basing the interval on the distribution of $\bar{X}_n - \mu(F)$, it may be better to look at a studentized quantity with hopes of higher efficiency. Next, the basic interval $I_{n,1}$ defined in (4.3) applies to other parameters; simply, replace $\mu(F)$ by the parameter $\theta(F)$, $\bar{X}_n$ by $\hat{\theta}_n$ and let $J_n(F)$ be the distribution of $n^{1/2}(\hat{\theta}_n - \theta(F))$ under $F$. Extensions to two (or more) sample problems are immediate as well; individual confidence bands for the unknown laws are constructed and utilized appropriately. Finally, when the observations are no longer real-valued, we no longer have the convenience of distribution-free confidence bands for the unknown distribution. But, these bands really play a secondary role, and we can utilize exponential inequalities for the sup distance between the empirical measure and the true measure to get conservative bands which may

be good enough.

As a final, more philosophical point, we advocate that the goal of constructing procedures that have finite sample validity requirements should be a primary consideration. There are many methods that enjoy excellent properties, such as (1.1) with $p = 2$. Rather than finding methods that satisfy (1.1) with $p = 3$, for example, we should only do so when we do not have to compromise on finite sample validity. Thus, now that we have intervals that are conservative and efficient, can we retain these two properties and still have (1.1) with $p$ large?

# Appendix

**Proof of Lemma 4.1.** For purposes of the proof, set $z = z_{1-\frac{\alpha}{2}+\beta_n}$, $\Delta = \Delta_n(F)$, $\delta = \delta_{n,\alpha}(F)$ and $r = 1 + \delta$. By Taylor's theorem,

$$\Phi(zr) - \Phi(z) = z(r-1)\phi(z) + \frac{1}{2}z^2(r-1)^2\phi'(z^*),$$

where $z^*$ is between $z$ and $zr$. But, $\phi'(z^*) \leq 0$ and $|\phi'(z)|$ is maximized when $z = 1$, which leads to

$$\Phi(zr) - \Phi(z) \geq z(r-1)\phi(z) - \frac{1}{2}z^2(r-1)^2(2\pi e)^{-1/2}.$$

So, it suffices to choose $r$ large enough so that the right hand side of the last expression is greater than or equal to $\Delta$; equivalently, it suffices to choose $\delta$ large enough so that

$$-\frac{1}{2}z^2(2\pi e)^{-1/2}\delta^2 + z\phi(z)\delta - \Delta \geq 0.$$

Solving the quadratic leads to choosing $\delta$ so that

$$\delta \geq \frac{\phi(z) - [\phi^2(z) - c\Delta]^{1/2}}{\frac{zc}{2}},$$

where $c = 2(2\pi e)^{-1/2}$. (As in the statement of the lemma, we are assuming $\phi^2(z) \geq z\Delta$.) By expanding the function $f(\Delta) = (\phi^2(z) - c\Delta)^{1/2}$, so that

$$f(\Delta) = \phi(z) - \frac{1}{2}\Delta c[\phi^2(z) - c\Delta^*]^{-1/2}$$

for some $\Delta^*$ between $0$ and $\Delta$, it suffices to choose $\delta$ so that

$$\delta \geq \Delta[\phi^2(z) - c\Delta^*]^{-1/2}/z.$$

So, by monotonicity in $\Delta^*$, it suffices to choose $\delta$ so

$$\delta \geq \Delta[\phi^2(z) - c\Delta]^{-1/2}/z.$$

Hence, the smallest $\delta$ that will work is no bigger than the right hand side of this last expression.

# References

Anderson, T. (1967). Confidence limits for the expected value of an arbitrary bounded random variable with a continuous distribution function. *Bull. ISI* **43**, 249–251.

Bahadur, R. and Savage, L. (1956). The nonexistence of certain statistical procedures in nonparametric problems. *Annals of Mathematical Statistics* **25**, 1115–1122.

Beran, R. and Millar, W. (1985). Asymptotic theory of confidence sets. In *Proceedings of the Berkeley conference in honor of Jerzy Neyman and Jack Kiefer, volume II,* Wadsworth, 865–887.

Bickel, P. (1992). Inference and auditing: the Stringer bound. *International Statistical Review* **60**, 197–209.

Breth, M., Maritz, J. and Williams, E. (1978). On distribution–free lower confidence limits for the mean of a nonnegative random variable. *Biometrika* **65**, 529–534.

Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *Annals of Statistics* **7**, 1–26.

Gasko, M. (1991). A new distribution-free lower confidence bound for the mean of a positive random variable. Working Paper QS 92-006, San Jose State University College of Business.

Hajek, J. (1972). Local asymptotic minimax and admissibility in estimation. *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability* **1**, 175–194.

Hall, P. (1992). *The bootstrap and Edgeworth expansion.* Springer-Verlag, New York.

Hall, P. and Jing, B. (1995). Uniform coverage bounds for confidence intervals and Berry-Esseen theorems for Edgeworth expansion. *Annals of Statistics* **23**, 363–375.

Loh, W. (1985). A new method for testing separate families of hypotheses. *Journal of the American Statistical Association* **80**, 362–368.

Millar, W. (1983). The minimax principle in asymptotic statistical theory. *Ecole d'Ete-St. Flour XI, Springer Lecture Notes* **976**, 76–265.

Romano, J. (1989). Do bootstrap confidence procedures behave well uniformly in $P$? *Canadian Journal of Statistics* **17**, 75–80.

Silvapulle, M. (1996). A test in the presence of nuisance parameters. *Journal of the American Statistical Association* **91**, 1690–1693.

Stringer, K. (1963). Practical aspects of statistical sampling in auditing. *Proceedings of Business and Economic Statistics Section*, American Statistical Association.