
Towards Practical Mean Bounds for Small Samples

Anonymous Authors¹

Abstract

Bounds on the mean of an unknown distribution and their associated confidence intervals are widely used in the scientific literature to characterize uncertainty of mean estimates. Historically, to bound the mean for small sample sizes, practitioners have had to choose between using methods with unrealistic assumptions about the unknown distribution (e.g., Gaussianity) and methods like Hoeffding's inequality that use weaker assumptions but produce much looser (wider) intervals. In 1969, Anderson (1969a) proposed a mean confidence interval strictly better than Hoeffding's whose only assumption is that the support of the distribution falls in an interval $(-\infty, b]$. For the first time since then, we present a new family of bounds, based on the same assumptions, at least one of which *dominates* Anderson's, meaning that it is always as tight or tighter than Anderson's. We prove that our bounds have *guaranteed coverage*, i.e., they hold with probability at least $1 - \alpha$ for all distributions on an interval $(-\infty, b]$. In simulations, we show that for many distributions, the gain over Anderson's bound is substantial.

1. Introduction

In this work, we revisit the classic statistical problem of defining a confidence interval on the mean μ of an unknown distribution with CDF F from an i.i.d. sample $\mathbf{X} = X_1, X_2, \dots, X_n$, and the closely related problems of producing upper or lower confidence bounds on the mean. To produce a non-trivial upper bound, one must make assumptions about F , such as finite variance, sub-Gaussianity, or that its support is contained on a finite interval $[a, b]$ or a semi-infinite interval $(-\infty, b]$. We adopt this last assumption. For simplicity, we focus on upper bounds, but the development for lower bounds and confidence intervals is

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

similar.

An upper confidence bound has *guaranteed coverage* for a set of distributions \mathcal{F} if, for all sample sizes $1 \leq n \leq \infty$, for all confidence levels $1 - \alpha \in (0, 1)$, and for all distributions $F \in \mathcal{F}$, the bound $\mu_{\text{upper}}^{1-\alpha}$ satisfies

$$\text{Prob}_F[\mu \leq \mu_{\text{upper}}^{1-\alpha}(X_1, X_2, \dots, X_n)] \geq 1 - \alpha, \quad (1)$$

where μ is the mean of the unknown distribution F .

Among bounds with guaranteed coverage for distributions on an interval $[a, b]$, our interest is in bounds with good performance on *small sample sizes*. The reason is that, for 'large enough' sample sizes, excellent bounds and confidence intervals already exist. In particular, the confidence intervals based on Student's t -statistic (Student, 1908) are satisfactory in terms of coverage and accuracy for most practitioners, given that the sample size is greater than some threshold.¹

The validity of the Student's t method depends upon the Gaussianity of the sample mean, which, strictly speaking does not hold for any finite sample size unless the original distribution itself is Gaussian. However, for many applications, the sample mean becomes close enough to Gaussian as the sample size grows (due to the effects described by the central limit theorem), that the resulting bounds hold with probabilities close to the confidence level. Such results vary depending upon the unknown distribution, but it is generally accepted that a large enough sample size can be defined to cover any distributions that might occur in a given situation.² This question is what to do when the sample size is smaller than such a threshold. Our bounds provide a new and better option for guaranteed coverage with small sample sizes—our bounds are tighter for *every possible sample* than the bound by Anderson (Anderson, 1969a), which is arguably the best existing bound with guaranteed coverage for small sample sizes. Below we review bounds with coverage guarantees, those that do *not* exhibit guaranteed coverage, and those for which the result is unknown.

¹Adequate sample size for Student's t method depends upon the setting, but a common rule is $n > 30$.

²An example in which the sample mean is still visibly skewed (and hence inappropriate for use with Student's t) even after $n = 80$ samples is given for log-normal distributions in the supplementary material.

055 **1.1. Distribution free bounds with guaranteed coverage**

056 Several bounds exist that have guaranteed coverage. These
 057 include Hoeffding's inequality (Hoeffding, 1963), Anderson's bound
 058 (Anderson, 1969a), and the bound due to Maurer & Pontil (2009).
 059

060 **Hoeffding's inequality.** For a distribution F on $[a, b]$, Hoeffding's inequality (Hoeffding, 1963) provides a bound on
 061 the probability that the sample mean, $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, will deviate from the mean by more than some amount,
 062 $t \geq 0$:

$$\Pr(\mu - \bar{X}_n \leq t) \leq e^{-\frac{2nt^2}{(b-a)^2}}. \quad (2)$$

063 Defining α to be the right hand side of this inequality, solving
 064 for t as a function of α , and rewriting in terms of α
 065 rather than t , one obtains a $1 - \alpha$ upper confidence bound
 066 on the mean of

$$b^{\alpha, \text{Hoeffding}}(\mathbf{X}) \stackrel{\text{def}}{=} \bar{X}_n + (b-a)\sqrt{\frac{\ln(1/\alpha)}{2n}}. \quad (3)$$

067 **Maurer and Pontil.** One limitation of Hoeffding's inequality
 068 is that the amount added to the sample mean to obtain the
 069 upper confidence bound scales with the range of the random
 070 variable over \sqrt{n} , which shrinks slowly as n increases.
 071

072 Bennett's inequality (Bennett, 1962) considers both the sample
 073 mean and the sample variance and obtains a better dependence
 074 on the range of the random variable *when the variance is known*. Maurer & Pontil (2009) derived an
 075 upper confidence bound for the variance of a random variable,
 076 and suggest combining this with Bennett's inequality (via the
 077 union bound) to obtain the following $1 - \alpha$ upper confidence
 078 bound on the mean

$$b^{\alpha, \text{M&P}}(\mathbf{X}) \stackrel{\text{def}}{=} \bar{X}_n + \frac{7(b-a)\ln(2/\alpha)}{3(n-1)} + \sqrt{\frac{2\hat{\sigma}^2 \ln(2/\alpha)}{n}}.$$

079 Notice that Maurer and Pontil's upper confidence bound
 080 scales with the range, $(b-a)$, divided by n (as opposed
 081 to the \sqrt{n} of Hoeffding's). However, the \sqrt{n} dependence
 082 is unavoidable to some extent: Maurer and Pontil's upper
 083 confidence bound scales with the sample standard deviation,
 084 $\hat{\sigma}$, divided by \sqrt{n} . As a result, Maurer and Pontil's bound
 085 tends to be tighter than Hoeffding's when both n is large and
 086 the range of the random variable is large relative to the variance.
 087 Lastly, notice that Maurer and Pontil's bound requires
 088 $n \geq 2$ for the sample standard deviation to be defined.

089 **Anderson's bound.** Anderson's bound (Anderson, 1969a)³
 090 is based upon the well-known Kolmogorov-Smirnov (KS)
 091 statistic, which gives the maximum distance between a
 092 continuous distribution's CDF and the empirical CDF of a
 093 sample from that distribution. Using the distribution of the
 094

095 ³An easier to access and virtually equivalent version of Anderson's work can be found in Anderson (1969b).

096 KS-statistic (which is independent of the underlying distribution F , as long as it is continuous), one may define
 097 an 'envelope' around the empirical CDF that, with high
 098 probability, contains the true CDF. The upper and lower
 099 extremes of such an envelope define the CDFs with the minimum
 100 and maximum attainable means for distributions that
 101 fit within the envelope, and thus bound the mean with high
 102 probability.

103 In practice, Anderson's inequality tends to be significantly
 104 tighter than Maurer and Pontil's inequality unless the variance
 105 of the random variable is minuscule in comparison to the range of the random variable (and n is sufficiently
 106 large). However, neither Anderson's inequality nor Maurer
 107 and Pontil's inequality strictly dominates the other. That is,
 108 neither upper bound is strictly less than or equal to the other
 109 in all cases. However, Anderson's inequality *does* dominate
 110 Hoeffding's inequality (Learned-Miller & Thomas, 2019).

111 A variety of variations on the theme of Anderson's bound
 112 have been proposed, including those by Diouf & Dufour
 113 (2005), Learned-Miller & DeStefano (2008), and Romano &
 114 Wolf (2000). However, none of these variations are shown
 115 to dominate Anderson's original bound. That is, while they
 116 give tighter intervals for some samples, they are looser for
 117 others. Finally we mention a bound due to Fienberg et al.
 118 (1977). This bound applies to distributions on a discrete set
 119 of support points, but nothing prevents it, in theory, from
 120 being applied to an arbitrarily dense set of points on an interval
 121 such as $[0, 1]$. This bound has a number of appealing
 122 properties, and comes with a proof of guaranteed coverage.
 123 However, the main drawback is that it is currently computationally
 124 intractable, with a computation time that depends exponentially
 125 on the number of points in the support set, precluding many (if not most) practical applications.

1.2. Bounds that do not exhibit guaranteed coverage

126 Many bounds that are used in practice are known to violate
 127 Eq. (1) for certain distributions. These include the aforementioned
 128 Student's t method, and various bootstrap procedures,
 129 such as the bias-corrected and accelerated (BCa) bootstrap
 130 and the percentile bootstrap. See Efron & Tibshirani (1993)
 131 for details of these methods. A simple explanation of the
 132 failure of bootstrap methods for certain distributions is given
 133 by Romano & Wolf (2000, pages 757–758). Presumably if
 134 one wants guarantees of Eq. (1), one cannot use these methods
 135 (unless one has extra information about the unknown
 136 distribution).

1.3. Bounds conjectured to have guaranteed coverage

137 There are at least two known bounds that perform well in
 138 practice but for which no proofs of coverage are known. One
 139 of these, used in accounting procedures, is the so-called
 140 Stringer bound (Stringer, 1963). It is known to violate

Eq. (1) for confidence levels $\alpha > 0.5$ (Pap & van Zuijlen, 1995), but its coverage for $\alpha < 0.5$ is unknown.

A little known bound by Gaffke (2005) gives remarkably tight bounds on the mean, but has eluded a proof of guaranteed coverage. This bound was recently rediscovered by Learned-Miller & Thomas (2019), who do an empirical study of its performance and provide a method for computing it efficiently.

We demonstrate in Section 4 that our bound dominates those of both Hoeffding and Anderson. **To our knowledge, this is the first bound that has been shown to dominate Anderson's bound.**

2. A Family of Confidence Bounds

In this section we define our new upper confidence bound. Let n be the sample size. We use bold-faced letters to denote a vector of size n and normal letters to denote a scalar. Upper-case letters denote random variables and lower-case letters denote values taken by them. For example, $X_i \in \mathcal{R}$ and $\mathbf{X} = (X_1, \dots, X_n) \in \mathcal{R}^n$ are random variables. $x_i \in \mathcal{R}$ is a value of X_i , and $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{R}^n$ is a value of \mathbf{X} . For a sample \mathbf{x} , we let $F(\mathbf{x}) \stackrel{\text{def}}{=} (F(x_1), \dots, F(x_n)) \in [0, 1]^n$. Order statistics play a central role in our work. We denote random variable order statistics $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ and of a specific sample as $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$.

Given a sample $\mathbf{X} = \mathbf{x}$ of size n and a confidence level $1-\alpha$, we would like to calculate an upper confidence bound for the mean. Let F be the CDF of X_i , i.e., the true distribution and $D \subset \mathcal{R}$ be the support of F . We assume that D has a finite upper bound. Given D and any function $T : D^n \rightarrow \mathcal{R}$ we will calculate an upper confidence bound $b_{D,T}^\alpha(\mathbf{x})$ for the mean of F .

We show in Lemma 2.1 that if D^+ is a superset of D with finite upper bound, then $b_{D^+,T}^\alpha(\mathbf{x}) \geq b_{D,T}^\alpha(\mathbf{x})$. Therefore we only need to know a superset of the support with finite upper bound to obtain a guaranteed bound.

Let $s_D \stackrel{\text{def}}{=} \sup\{x : x \in D\}$. We next describe a method for pairing the sample \mathbf{x} with another vector $\ell \in [0, 1]^n$ to produce a staircase CDF function $G_{\mathbf{x},\ell}$. Let $x_{(n+1)} \stackrel{\text{def}}{=} s_D$. Consider the step function $G_{\mathbf{x},\ell} : \mathcal{R} \rightarrow [0, 1]$ defined from ℓ and \mathbf{x} as follows (see Figure 1):

$$G_{\mathbf{x},\ell}(y) = \begin{cases} 0, & \text{if } y < x_{(1)} \\ \ell_{(i)}, & \text{if } x_{(i)} \leq y < x_{(i+1)} \\ 1, & \text{if } y \geq s_D. \end{cases} \quad (4)$$

In particular, when $\ell = (1/n, \dots, n/n)$, $G_{\mathbf{x},\ell}$ becomes the empirical CDF. Also note that when $\ell = F(\mathbf{x})$, $\forall x, G_{\mathbf{x},\ell}(x) \leq F(x)$, as illustrated in Figure 2.

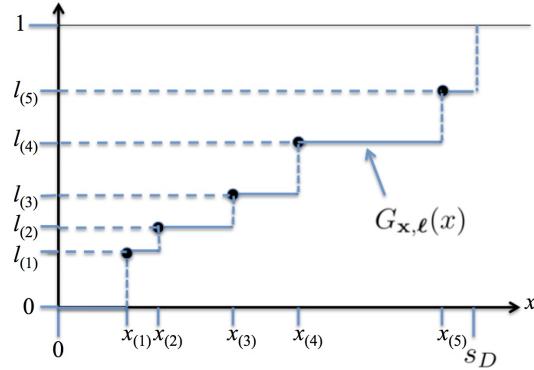


Figure 1. The staircase function $G_{\mathbf{x},\ell}$, which is a function of the sample \mathbf{x} and a vector ℓ of values between 0 and 1. When $\ell = (1/n, \dots, n/n)$, $G_{\mathbf{x},\ell}$ becomes the empirical CDF.

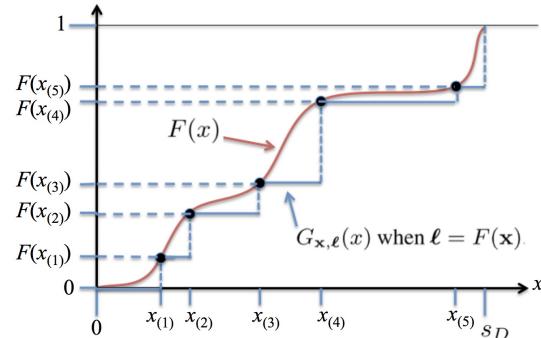


Figure 2. The CDF of a distribution F in red, with a random sample of five order statistics on the x-axis. The blue staircase function shows the function $G_{\mathbf{x},\ell}(x)$ when $\ell = F(\mathbf{x})$. Notice that for all x , $G_{\mathbf{x},\ell}(x) \leq F(x)$.

Following Learned-Miller & Thomas (2019), if we consider $G_{\mathbf{x},\ell}$ to be a CDF, we can compute the mean of the resulting distribution as a function of two vectors \mathbf{x} and ℓ as

$$m_D(\mathbf{x}, \ell) \stackrel{\text{def}}{=} \sum_{i=1}^{n+1} x_{(i)}(\ell_{(i)} - \ell_{(i-1)}) \quad (5)$$

$$= s_D - \sum_{i=1}^n \ell_{(i)}(x_{(i+1)} - x_{(i)}), \quad (6)$$

where $\ell_{(0)} \stackrel{\text{def}}{=} 0$, $\ell_{(n+1)} \stackrel{\text{def}}{=} 1$ and $x_{(n+1)} \stackrel{\text{def}}{=} s_D$. When s_D is finite, this is well-defined. Notice that this function is defined in terms of the *order statistics* of \mathbf{x} and ℓ . Learned-Miller & Thomas (2019) refer to this as the *induced mean* for the sample \mathbf{x} by the vector ℓ .

An ordering on D^n . Next, we introduce a scalar-valued function T which we will use to define a total order on samples in D^n , and define a set of samples less than or equal to another sample. In particular, for any function $T : \mathcal{R}^n \rightarrow \mathcal{R}$, let $\mathbb{S}_{D,T}(\mathbf{x}) = \{\mathbf{y} \in D^n | T(\mathbf{y}) \leq T(\mathbf{x})\}$.

The greatest induced mean for a given \mathbf{U} . Let $\mathbf{U} =$

165 U_1, \dots, U_n be a sample of size n from the continuous uniform distribution on $[0, 1]$, with $\mathbf{u} \stackrel{\text{def}}{=} (u_1, \dots, u_n)$ being a particular sample of \mathbf{U} .
 166
 167
 168

169 Now consider the random quantity

$$170 \quad b_{D,T}(\mathbf{x}, \mathbf{U}) \stackrel{\text{def}}{=} \sup_{\mathbf{z} \in \mathbb{S}_{D,T}(\mathbf{x})} m_D(\mathbf{z}, \mathbf{U}), \quad (7)$$

171
 172

173 which depends upon a fixed sample \mathbf{x} (non-random) and also on the random variable \mathbf{U} .
 174
 175

176 **Our upper confidence bound.** Let $0 < p < 1$. Let $177 Q(p, Y)$ be the *quantile function* of the scalar random variable Y , i.e.,
 178
 179

$$180 \quad Q(p, Y) \stackrel{\text{def}}{=} \inf\{y \in \mathbb{R} : F_Y(y) \geq p\}, \quad (8)$$

181
 182

183 where $F_Y(y)$ is the CDF of Y . We define $b_{D,T}^\alpha(\mathbf{x})$ to be the $(1 - \alpha)$ -quantile of the random quantity $b_{D,T}(\mathbf{x}, \mathbf{U})$.
 184
 185

186 **Definition 2.1** (Upper confidence bound on the mean). Given a sample \mathbf{x} and a confidence level $1 - \alpha$:

$$187 \quad b_{D,T}^\alpha(\mathbf{x}) \stackrel{\text{def}}{=} Q(1 - \alpha, b_{D,T}(\mathbf{x}, \mathbf{U})), \quad (9)$$

188
 189

190 where $b_{D,T}(\mathbf{x}, \mathbf{U})$ is defined in Eq. 7.

191
 192

193 To simplify notation, we drop the superscript and subscripts whenever clear. We show in Section 2.1 that this upper confidence bound has guaranteed coverage for all sample sizes n , for all confidence levels $0 < 1 - \alpha < 1$ and for all distributions F and support D where s_D is finite.
 194
 195

196 We show below that a bound computed from a superset $197 D^+ \supseteq D$ will be looser than a bound computed from the 198 support D . Therefore it is enough to know a superset of the 199 support D to obtain a bound with guaranteed coverage.
 200

201 **Lemma 2.1.** Let $D^+ \supseteq D$ where s_{D^+} is finite. For any 202 sample \mathbf{x} :

$$203 \quad b_D^\alpha(\mathbf{x}) \leq b_{D^+}^\alpha(\mathbf{x}). \quad (10)$$

204
 205

206 *Proof.* Since s_{D^+} is finite, $m_{D^+}(\mathbf{y}, \mathbf{u})$ is well-defined. 207 Since $D \subseteq D^+$, for any \mathbf{y} and \mathbf{u} , $m_D(\mathbf{y}, \mathbf{u}) \leq m_{D^+}(\mathbf{y}, \mathbf{u})$.
 208 Then

$$209 \quad \sup_{\mathbf{y} \in \mathbb{S}_D(\mathbf{x})} m_D(\mathbf{y}, \mathbf{u}) \leq \sup_{\mathbf{y} \in \mathbb{S}_D(\mathbf{x})} m_{D^+}(\mathbf{y}, \mathbf{u}) \quad (11)$$

$$210 \quad \leq \sup_{\mathbf{y} \in \mathbb{S}_{D^+}(\mathbf{x})} m_{D^+}(\mathbf{y}, \mathbf{u}), \quad (12)$$

211
 212

213 where the last inequality is because $\mathbb{S}_D(\mathbf{x}) \subseteq \mathbb{S}_{D^+}(\mathbf{x})$.
 214 Let $b_D(\mathbf{x}, \mathbf{U}) = \sup_{\mathbf{z} \in \mathbb{S}_D(\mathbf{x})} m_D(\mathbf{z}, \mathbf{U})$ and $b_{D^+}(\mathbf{x}, \mathbf{U}) = \sup_{\mathbf{z} \in \mathbb{S}_{D^+}(\mathbf{x})} m_{D^+}(\mathbf{z}, \mathbf{U})$. Then $b_D^\alpha(\mathbf{x})$ and $b_{D^+}^\alpha(\mathbf{x})$ are the $(1 - \alpha)$ -quantiles of $b_D(\mathbf{x}, \mathbf{U})$ and $b_{D^+}(\mathbf{x}, \mathbf{U})$. Since $b_D(\mathbf{x}, \mathbf{U}) \leq b_{D^+}(\mathbf{x}, \mathbf{U})$ for any \mathbf{u} , $b_D^\alpha(\mathbf{x}) \leq b_{D^+}^\alpha(\mathbf{x})$. \square
 215
 216
 217
 218
 219

In Section 2.1 we show that the bound has guaranteed coverage. In Section 3 we discuss how to efficiently compute the bound. In Section 4 we show that when T is a certain linear function, the bound is equal to or tighter than Anderson's for any sample. In addition, we show that when the support is known to be $\{0, 1\}$, our bound recovers the well-known Clopper-Pearson confidence bound for binomial distributions (Clopper & Pearson, 1934). In Section 5, we present simulations that show the consistent superiority of our bounds over previous bounds.

2.1. Guaranteed Coverage

In this section we show that our bound has guaranteed coverage in Theorem 2.7. We omit superscripts and subscripts if they are clear from context.

2.1.1. PREVIEW OF PROOF

We explain the idea behind our bound at a high level using a special case. Note that our proof is more general than our special case, which makes assumptions such as the continuity of F to simplify the intuition.

Suppose that F is continuous, then the *probability integral transform* $F_X(X)$ of X is uniformly distributed on $[0, 1]$ (Angus, 1994). If there exists a sample \mathbf{x}_{\max} with the largest $b^\alpha(\mathbf{x})$ such that $b^\alpha(\mathbf{x}) < \mu$, then the probability a sample \mathbf{Z} outputs $b^\alpha(\mathbf{Z}) < \mu$ is equal to the probability \mathbf{Z} outputs $b^\alpha(\mathbf{Z}) \leq b^\alpha(\mathbf{x}_{\max})$ (the yellow region on the left of Fig. 3). This is the region where the bound fails, and we would like to show that the probability of this region is at most α .

Let $\mathbf{U} \stackrel{\text{def}}{=} F(\mathbf{Z})$ and $\mathbf{u} \stackrel{\text{def}}{=} F(\mathbf{z})$. Then U_i is uniformly distributed on $[0, 1]$. If F is invertible, we can transform the region $\{\mathbf{z} : b^\alpha(\mathbf{z}) \leq b^\alpha(\mathbf{x}_{\max})\}$ to $\{\mathbf{u} : b^\alpha(F^{-1}(\mathbf{u})) \leq b^\alpha(\mathbf{x}_{\max})\}$ where $F^{-1}(\mathbf{u}) \stackrel{\text{def}}{=} (F^{-1}(u_1), \dots, F^{-1}(u_n))$ (the yellow region on the right of Fig. 3).

Through some calculations using the definition of function b , we can show that the yellow region $\{\mathbf{u} : b^\alpha(F^{-1}(\mathbf{u})) \leq b^\alpha(\mathbf{x}_{\max})\}$ is a subset of the striped region $\{\mathbf{u} : b(\mathbf{x}_{\max}, \mathbf{u}) \geq \mu\}$.

Note that since $b^\alpha(\mathbf{x}_{\max}) < \mu$, μ is larger than the $1 - \alpha$ quantile of $b(\mathbf{x}_{\max}, \mathbf{U})$. Therefore, by the definition of quantile, the probability of the striped region is at most α :

$$\mathbb{P}_{\mathbf{U}}(b(\mathbf{x}_{\max}, \mathbf{U}) \geq \mu) \leq \alpha, \quad (13)$$

and thus the probability of the yellow region is at most α .

2.1.2. MAIN RESULT

In this section, we present some supporting lemmas and then the main result in Theorem 2.7. The proofs of the simpler lemmas have been deferred to the supplementary material.

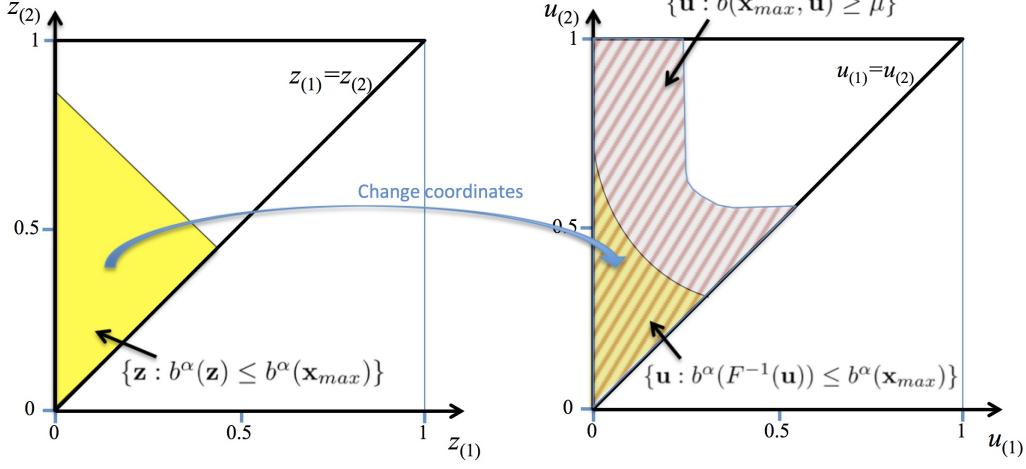


Figure 3. Illustrations of Section 2.1.1. **Left.** The yellow region shows samples of $\mathbf{z} = [z_{(1)}, z_{(2)}]$ such that $b^\alpha(\mathbf{z}) \leq b^\alpha(\mathbf{x}_{max})$. **Right.** The same yellow region, but in the coordinates $\mathbf{u} = F^{-1}(\mathbf{z})$. We will show that the yellow region is a subset of the striped, which contains \mathbf{u} such that $b(\mathbf{x}_{max}, \mathbf{u}) \geq \mu$.

Lemma 2.2. Let X be a random variable with CDF F and $Y \stackrel{\text{def}}{=} F(X)$, known as the probability integral transform of X . Let U be a uniform random variable on $[0, 1]$. Then for any $0 \leq y \leq 1$,

$$\mathbb{P}(Y \leq y) \leq \mathbb{P}(U \leq y). \quad (14)$$

If F is continuous, then Y is uniformly distributed on $[0, 1]$.

The next lemma is illustrated by Fig. 2. It shows that by building a ‘stairstep CDF’ using a random sample and points on the true CDF, the resulting distribution has a mean greater than or equal to the original distribution’s mean.

Lemma 2.3. For any $\mathbf{x} \in D^n$,

$$m_D(\mathbf{x}, F(\mathbf{x})) \geq \mu. \quad (15)$$

For use in the next lemma, we define a partial order for the samples on D^n . Note that it is defined with respect to the *order statistics* of the sample, not the original components.

Definition 2.2 (Partial Order). For any two samples \mathbf{z} and \mathbf{y} , we define $\mathbf{z} \preceq \mathbf{y}$ to indicate that $z_{(i)} \leq y_{(i)}$, $1 \leq i \leq n$.

Lemma 2.4. Let \mathbf{Z} be a random sample of size n from F . Let $\mathbf{U} = U_1, \dots, U_n$ be a sample of size n from the continuous uniform distribution on $[0, 1]$. For any function $T : D^n \rightarrow R$ and any $\mathbf{x} \in D^n$:

$$\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{x})) \leq \mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu). \quad (16)$$

Proof Sketch. Let \cup denote the union of events and $\{\}$ de-

note an event. Then for any $\mathbf{x} \in D^n$:

$$\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{x})) \quad (17)$$

$$= \mathbb{P}_{\mathbf{Z}}(\mathbf{Z} \in \mathbb{S}(\mathbf{x})) \quad (18)$$

$$= \mathbb{P}_{\mathbf{Z}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{\mathbf{Z} = \mathbf{y}\}) \quad (19)$$

$$\leq \mathbb{P}_{\mathbf{Z}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{\mathbf{Z} \preceq \mathbf{y}\}) \quad (20)$$

$$\leq \mathbb{P}_{\mathbf{Z}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{F(\mathbf{Z}) \preceq F(\mathbf{y})\}) \quad (\text{by monotone } F) \quad (21)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{\mathbf{U} \preceq F(\mathbf{y})\}). \quad (22)$$

The last step is by an extension of Lemma 2.2. Recall that $m_D(\mathbf{y}, \mathbf{u}) = s_D - \sum_{i=1}^n u_{(i)}(y_{(i+1)} - y_{(i)})$ where $\forall i, y_{(i+1)} - y_{(i)} \geq 0$. Therefore if $\mathbf{u} \preceq F(\mathbf{y})$ then $m_D(\mathbf{y}, \mathbf{u}) \geq m_D(\mathbf{y}, F(\mathbf{y}))$:

$$\mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{\mathbf{U} \preceq F(\mathbf{y})\}) \quad (23)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{m_D(\mathbf{y}, \mathbf{U}) \geq m_D(\mathbf{y}, F(\mathbf{y}))\}) \quad (24)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{m_D(\mathbf{y}, \mathbf{U}) \geq \mu\}), \text{ by Lem. 2.3} \quad (25)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\sup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} m_D(\mathbf{y}, \mathbf{U}) \geq \mu) \quad (26)$$

$$= \mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu). \quad (27)$$

□

We include a more detailed version of the proof for the above lemma in the supplementary material.

Lemma 2.5. Let $\mathbf{U} = U_1, \dots, U_n$ be a sample of size n from the continuous uniform distribution on $[0, 1]$. Let \mathbf{X} and \mathbf{Z} denote i.i.d. samples of size n from F . For any function $T : D^n \rightarrow \mathcal{R}$ and any $\alpha \in (0, 1)$,

$$\begin{aligned} & \mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{U}}(b_{D,T}(\mathbf{X}, \mathbf{U}) \geq \mu) \leq \alpha) \\ & \leq \mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{X})) \leq \alpha). \end{aligned} \quad (28)$$

275 *Proof.* From Lemma 2.4 for any sample \mathbf{x} ,

$$\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{x})) \leq \mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu). \quad (29)$$

276 Therefore,

$$\mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{X})) \leq \alpha) \quad (30)$$

$$\geq \mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{U}}(b(\mathbf{X}, \mathbf{U}) \geq \mu) \leq \alpha). \quad (31)$$

□

285 **Lemma 2.6.** Let $\mathbf{U} = U_1, \dots, U_n$ be a sample of size n
 286 from the continuous uniform distribution on $[0, 1]$. Let \mathbf{X}
 287 be a random sample of size n from F . For any function
 288 $T : D^n \rightarrow \mathcal{R}$ and any $\alpha \in (0, 1)$,

$$\mathbb{P}_{\mathbf{X}}(b_{D,T}^{\alpha}(\mathbf{X}) < \mu) \quad (32)$$

$$\leq \mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{U}}(b_{D,T}(\mathbf{X}, \mathbf{U}) \geq \mu) \leq \alpha). \quad (33)$$

293 *Proof.* Because $b^{\alpha}(\mathbf{x})$ is the $1 - \alpha$ quantile of $b(\mathbf{x}, \mathbf{U})$,
 294 by the definition of quantile: $\mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \leq b^{\alpha}(\mathbf{x})) \geq$
 295 $1 - \alpha$. Therefore $\mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq b^{\alpha}(\mathbf{x})) \leq \alpha$. If $b^{\alpha}(\mathbf{x}) < \mu$
 296 then $\mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu) \leq \alpha$. Since $b^{\alpha}(\mathbf{x}) < \mu$ implies
 297 $\mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu) \leq \alpha$, we have

$$\mathbb{P}_{\mathbf{X}}(b^{\alpha}(\mathbf{X}) < \mu) \quad (34)$$

$$\leq \mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{U}}(b(\mathbf{X}, \mathbf{U}) \geq \mu) \leq \alpha). \quad (35)$$

□

302 We now show that the bound has guaranteed coverage.

303 **Theorem 2.7.** Let \mathbf{X} be a random sample of size n from F .
 304 For any function $T : D^n \rightarrow \mathcal{R}$ and for any $\alpha \in (0, 1)$:

$$\mathbb{P}_{\mathbf{X}}(b_{D,T}^{\alpha}(\mathbf{X}) < \mu) \leq \alpha. \quad (36)$$

308 *Proof.* Let \mathbf{Z} be a random sample of size n from F .

$$\mathbb{P}_{\mathbf{X}}(b^{\alpha}(\mathbf{X}) < \mu) \quad (37)$$

$$\leq \mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{U}}(b(\mathbf{X}, \mathbf{U}) \geq \mu) \leq \alpha) \text{ by Lemma 2.6} \quad (38)$$

$$\leq \mathbb{P}_{\mathbf{X}}(\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{X})) \leq \alpha) \text{ by Lemma 2.5} \quad (39)$$

$$= \mathbb{P}(W \leq \alpha) \text{ where } W \stackrel{\text{def}}{=} \mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{X})) \quad (40)$$

$$\leq \alpha \text{ by Lemma 2.2.} \quad (41)$$

□

3. Computation

320 Let the superset of the support D^+ be a closed interval with
 321 a finite upper bound. If m is a continuous function,

$$\sup_{\mathbf{y} \in \mathbf{S}_{D^+}(\mathbf{x})} m(\mathbf{y}, \mathbf{u}) = \max_{\mathbf{y} \in \mathbf{S}_{D^+}(\mathbf{x})} m(\mathbf{y}, \mathbf{u}). \quad (42)$$

325 Therefore $b_{D^+}(\mathbf{x}, \mathbf{u})$ is the solution to

$$\max_{y_{(1)}, \dots, y_{(n)}} m(\mathbf{y}, \mathbf{u}) \quad (43)$$

327 subject to:

Algorithm 1 Monte Carlo estimation of $m_{D^+, T}^{\alpha}(\mathbf{x})$ where $D^+ = [0, 1]$. This pseudocode uses 1-based array indexing.

Input: A sample $\mathbf{x} \in D^n$, confidence parameter $1 - \alpha < 1$, a function $T : [0, 1]^n \rightarrow \mathcal{R}$ and Monte Carlo sampling parameter l .

Output: An estimation of $m_{D^+, T}^{\alpha}(\mathbf{x})$

$n \leftarrow \text{length}(\mathbf{x})$.

Create array \mathbf{ms} to hold l floating point numbers, and initialize it to zero.

Create array \mathbf{u} to hold n floating point numbers.

for $i \leftarrow 1$ to l **do**

for $j \leftarrow 1$ to n **do**

$\mathbf{u}[j] \sim \text{Uniform}(0,1)$.

end for

Sort(\mathbf{u} , ascending).

Solve: $M = \max_{y_{(1)}, \dots, y_{(n)}} m(\mathbf{y}, \mathbf{u})$ subject to:

1) $T(\mathbf{y}) \leq T(\mathbf{x})$.

2) $\forall i : 1 \leq i \leq n, 0 \leq y_{(i)} \leq 1$.

3) $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n)}$.

$\mathbf{ms}[i] = M$.

end for

Sort(\mathbf{ms} , ascending).

Return $\mathbf{ms}[\lceil (1 - \alpha)l \rceil]$.

1. $T(\mathbf{y}) \leq T(\mathbf{x})$,
2. $\forall i \in \{1, \dots, n\}, y_{(i)} \in D^+$,
3. $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n)}$.

When D^+ is an interval and T is linear, this is a linear programming problem and can be solved efficiently.

We can compute the $1 - \alpha$ quantile of a random variable M using Monte Carlo simulation, sampling $M l$ times. Letting $m_{(1)} \leq \dots \leq m_{(l)}$ be the sorted values, we output $m_{(\lceil (1 - \alpha)l \rceil)}$ as an approximation of the $1 - \alpha$ quantile.

Running time. When T is linear, the algorithm needs to solve a linear programming problem with n variables and $2n$ constraints l times. For sample size $n = 50$, computing the bound for each sample $\mathbf{x} \in D^n$ takes just a few seconds using $l = 10,000$ Monte Carlo samples.

4. Relationships with Existing Bounds

In this section, we compare our bound to previous bounds including those of Clopper and Pearson, Hoeffding, and Anderson. Proofs omitted in this section can be found in the supplementary material.

4.1. Special Case: Bernoulli Distribution

When we know that $D = \{0, 1\}$, the distribution is Bernoulli. If we choose T to be the sample mean, our bound becomes the same as the Clopper-Pearson confidence

330 bound for binomial distributions (Clopper & Pearson, 1934).
 331 See supplementary material for details.

333 4.2. Comparisons with Anderson and Hoeffding

335 In this section we show that for any sample size n , any
 336 confidence level α and for any sample \mathbf{x} , our method pro-
 337 duces a bound at most Anderson's bound (Theorem 4.3) and
 338 Hoeffding's bound (Theorem 4.4).

339 Note that if we only know an upper bound b of the support,
 340 we can set $D^+ = (-\infty, b]$ and our method still dominates
 341 Anderson's and Hoeffding's. Our bound becomes tighter
 342 or remains constant as the lower support bound increases,
 343 whereas Anderson's remains constant because it only uses
 344 the upper support bound as an input.

345 Anderson's bound constructs an upper bound for the mean
 346 by constructing a lower bound for the CDF. We defined a
 347 lower bound for the CDF as follows.

349 **Definition 4.1** (Lower confidence bound for the CDF). Let
 350 $\mathbf{X} = (X_1, \dots, X_n)$ be a sample of size n from the
 351 distribution with unknown CDF F . Let $\alpha \in (0, 1)$. Let
 352 $H_{\mathbf{X}} : \mathcal{R} \rightarrow [0, 1]$ be a function computed from the sam-
 353 ple \mathbf{X} such that for any CDF F ,

$$354 \mathbb{P}_{\mathbf{X}}(\forall x \in R, F(x) \geq H_{\mathbf{X}}(x)) \geq 1 - \alpha. \quad (44)$$

355 Then $H_{\mathbf{X}}$ is called a $(1 - \alpha)$ lower confidence bound for
 356 the CDF.
 357

358 In Figs. 1 and 2, it is easy to see that if the staircase function
 359 $G_{\mathbf{X}, \ell}$ is a lower confidence bound for the CDF then its
 360 induced mean $m(\mathbf{X}, \ell)$ is an upper confidence bound for μ .

361 **Lemma 4.1.** Let $\mathbf{X} = (X_1, \dots, X_n)$ be a sample of size n
 362 from a distribution with mean μ . Let $\ell \in [0, 1]^n$. If $G_{\mathbf{X}, \ell}$ is
 363 a $(1 - \alpha)$ lower confidence bound for the CDF then

$$364 \mathbb{P}_{\mathbf{X}}(m(\mathbf{X}, \ell) \geq \mu) \geq 1 - \alpha. \quad (45)$$

365 Anderson defined a $1 - \alpha$ lower CDF confidence bound
 366 $G_{\mathbf{X}, \mathbf{u}^{And}}$ where $\mathbf{u}^{And} \in [0, 1]^n$ is defined as

$$367 u_i^{And} \stackrel{\text{def}}{=} \max \{0, i/n - \beta(n)\}. \quad (46)$$

368 Anderson identifies $\beta(n)$ as the one-sided Kolmogorov-
 369 Smirnov statistic that results in equality for Eq. 44. Using
 370 the fact that $G_{\mathbf{X}, \mathbf{u}^{And}}$ is a $1 - \alpha$ lower confidence bound
 371 for the CDF, and letting $D^+ = (-\infty, b]$, Anderson's $1 - \alpha$
 372 upper bound for the mean is

$$373 b^{\alpha, \text{Anderson}}(\mathbf{x}) \stackrel{\text{def}}{=} m_{D^+}(\mathbf{x}, \mathbf{u}^{And}). \quad (47)$$

374 Learned-Miller & Thomas (2019) show that for any sample
 375 \mathbf{x} , Anderson's bound is better than Hoeffding's:

376 **Lemma 4.2** (Theorem 2 from (Learned-Miller & Thomas,
 377 2019)). For any sample size n , for any sample value $\mathbf{x} \in$
 378 D^n , for all $\alpha \in (0, 1)$:

$$379 b^{\alpha, \text{Anderson}}(\mathbf{x}) \leq b^{\alpha, \text{Hoeffding}}(\mathbf{x}). \quad (48)$$

We now compare our method with Anderson's bound:

380 **Theorem 4.3.** Let $D^+ = (-\infty, b]$. For any sample size
 381 n , for any sample value $\mathbf{x} \in D^n$, for all $\alpha \in (0, 1)$, using
 382 $T(\mathbf{x}) = b^{\alpha, \text{Anderson}}(\mathbf{x})$ yields:

$$383 b_{D^+, T}^{\alpha}(\mathbf{x}) \leq b^{\alpha, \text{Anderson}}(\mathbf{x}). \quad (49)$$

We explain briefly why this is true. First, from Figure 2, we
 384 can see that if $G_{\mathbf{X}, \mathbf{u}^{And}}$ is a lower confidence bound then
 385 $\forall i, F(X_{(i)}) \geq \mathbf{u}_{(i)}^{And}$. Note that $G_{\mathbf{X}, \mathbf{u}^{And}}$ must be a lower
 386 bound for all unknown CDF F , so we can pick a continuous
 387 F where, according to Lemma 2.2, $U \stackrel{\text{def}}{=} F(X)$ is uniformly
 388 distributed on $[0, 1]$. Therefore \mathbf{u}^{And} satisfies:

$$389 \mathbb{P}_{\mathbf{U}}(\forall i, U_{(i)} \geq \mathbf{u}_{(i)}^{And}) \geq 1 - \alpha, \quad (50)$$

390 where $U_{(i)}$'s are the order statistics of the uniform distribu-
 391 tion. Since $b(\mathbf{x}, \mathbf{U})$ is defined from linear functions of \mathbf{U}
 392 with negative coefficients (Eq. 6), if $\forall i, U_{(i)} \geq \mathbf{u}_{(i)}^{And}$ then
 393 $b(\mathbf{x}, \mathbf{U}) \leq b(\mathbf{x}, \mathbf{u}^{And})$. Therefore with probability at least
 394 $1 - \alpha$, $b(\mathbf{x}, \mathbf{U}) \leq b(\mathbf{x}, \mathbf{u}^{And})$. So $b(\mathbf{x}, \mathbf{u}^{And})$ is at least the
 395 $1 - \alpha$ quantile of $b(\mathbf{x}, \mathbf{U})$, which is the value of our bound.
 396 Therefore $b(\mathbf{x}, \mathbf{u}^{And})$ is at least the value of our bound.

397 Finally, if T is Anderson's bound, through some calcula-
 398 tions we can show that $b_{D^+, T}^{\alpha}(\mathbf{x}, \mathbf{u}^{And}) = m_{D^+}(\mathbf{x}, \mathbf{u}^{And})$,
 399 which is Anderson's bound. The result follows.

400 The comparison with Hoeffding's bound follows directly
 401 from Lemma 4.2 and Theorem 4.3:

402 **Theorem 4.4.** Let $D^+ = (-\infty, b]$. For any sample size
 403 n , for any sample value $\mathbf{x} \in D^n$, for all $\alpha \in (0, 1)$, using
 404 $T(\mathbf{x}) = b^{\alpha, \text{Anderson}}(\mathbf{x})$ yields:

$$405 b_{D^+, T}^{\alpha}(\mathbf{x}) \leq b^{\alpha, \text{Hoeffding}}(\mathbf{x}). \quad (51)$$

406 Note that the result from Theorem 4.3 can be generalized
 407 for bounds $m(\mathbf{X}, \ell)$ constructed from a $(1 - \alpha)$ confidence
 408 lower bound $G_{\mathbf{X}, \ell}$ using Lemma 4.1. Diouf & Dufour
 409 (2005) present several such bounds with different ℓ com-
 410 puted from the Anderson-Darling or the Eicker statistics
 411 (Theorem 4, 5 and Theorem 6 with constant ϵ). Theorem 4.3
 412 is also true for other value of $\beta(n)$ used to compute \mathbf{u}^{And} as
 413 long as $G_{\mathbf{X}, \mathbf{u}^{And}}$ is a $1 - \alpha$ lower CDF confidence bound.
 414 We show the general case in the supplementary material.

5. Simulations

We perform simulations to compare our bounds to Hoeffding's inequality, Anderson's bound, Maurer and Pontil's,

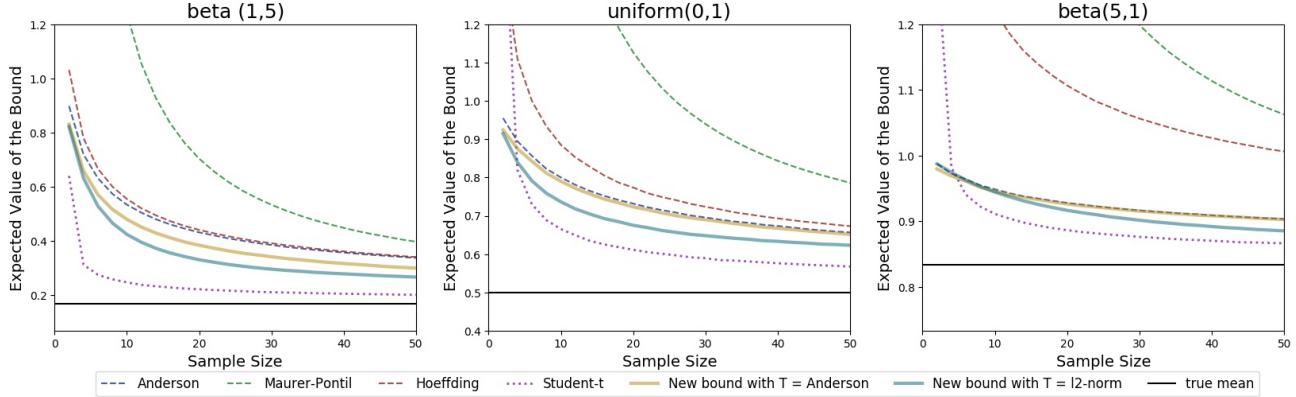


Figure 4. The expected value of the bounds. For each sample size, we sample \mathbf{X} 10,000 times, compute the bound for each sample, and take the average. Our new bound with T being Anderson’s bound consistently has lower expected value than Anderson’s (Theorem 4.3), Hoeffding’s (Theorem 4.4) and Maurer and Pontil’s. With T being the l_2 -norm, the bound is substantially tighter in these examples, and also has guaranteed coverage.

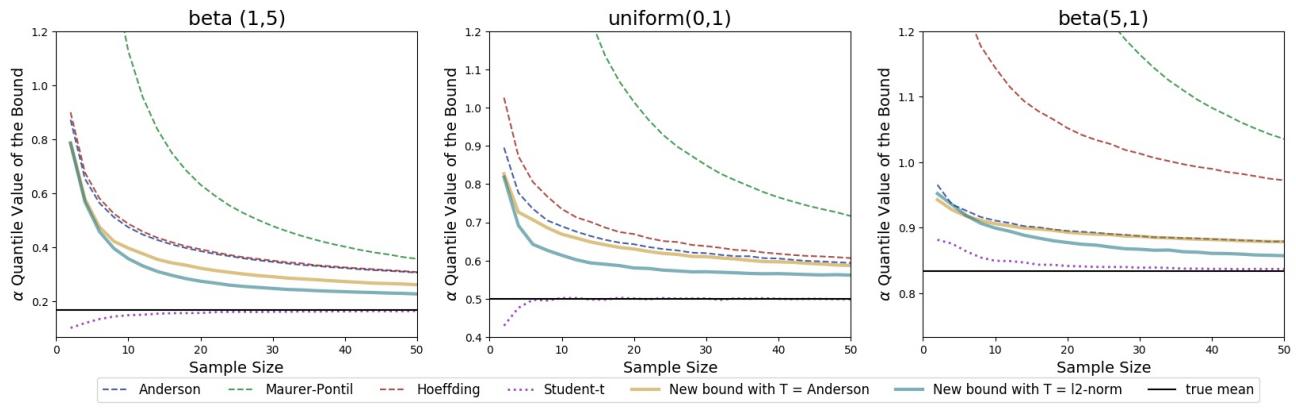


Figure 5. The α -quantile of the bound distribution. For each sample size, we sample \mathbf{X} 10,000 times, compute the bound for each sample, and take the α quantile. If the α -quantile is below the true mean, the bound does not have guaranteed coverage. For the uniform($0, 1$) and beta($1, 5$) distribution, when the sample size is small, Student-t does not have guarantee.

and Student-t’s bound (Student, 1908), the latter being

$$b^{\alpha, \text{Student}}(\mathbf{X}) \stackrel{\text{def}}{=} \bar{X}_n + \sqrt{\frac{\hat{\sigma}^2}{n}} t_{1-\alpha, n-1}. \quad (52)$$

For Anderson’s bound, we use the Dvoretzky-Kiefer-Wolfowitz inequality (Dvoretzky et al., 1956) to define the $1 - \alpha$ CDF lower bound via $\beta(n) = \sqrt{\ln(1/\alpha)/(2n)}$. We use $D = [0, 1]$ and $l = 10,000$ Monte Carlo samples. We consider two functions T :

1. Anderson: $T(\mathbf{x}) = b^{\alpha, \text{Anderson}}(\mathbf{x})$, again with $\beta(n) = \sqrt{\ln(1/\alpha)/(2n)}$. Because this T is linear in \mathbf{x} , it can be computed with the linear program in Eq. 42.
2. l_2 norm: $T(\mathbf{x}) = (\sum_{i=1}^n x_i^2)/n$. In this case, T requires the optimization of a linear functional over a convex region, which results in a simple convex optimization problem.

We perform experiments on three distributions: beta($1, 5$) (skewed right), uniform($0, 1$) and beta($5, 1$) (skewed left). Their PDFs are included in the supplementary material for reference. Additional experiments are in the supplementary material.

In Figure 4 and Figure 5 we plot the expected value and the α -quantile value of the bounds as the sample size increases. Consistent with Theorem 4.3, our bound with T being Anderson’s bound outperforms Anderson’s bound. Our new bound performs better than Anderson’s in distributions that are skewed left, and becomes similar to Anderson’s in right-skewed distributions. Our bound outperforms Hoeffding and Maurer and Pontil’s for all three distributions. Student-t fails (the error rate exceeds α) for beta($1, 5$) and uniform($0, 1$) when the sample size is small (Figure 5).

References

- 440 Anderson, T. W. Confidence limits for the value of an
441 arbitrary bounded random variable with a continuous
442 distribution function. *Bulletin of The International and*
443 *Statistical Institute*, 43:249–251, 1969a.
- 444 Anderson, T. W. Confidence limits for the value of an arbitrary bounded random variable with a continuous distribution function. *Technical Report Number 1, Department of Statistics, Stanford University*, 1969b.
- 445 Angus, J. E. The probability integral transform and related results. *SIAM Rev.*, 36(4):652–654, December 1994. ISSN 0036-1445. doi: 10.1137/1036146. URL <https://doi.org/10.1137/1036146>.
- 446 Bennett, G. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57(297):33–45, 1962.
- 447 Clopper, C. and Pearson, E. S. The use of confidence or fiducial limits illustrated in the case of the binomial. *Biometrika*, 26(4):404–413, 1934.
- 448 Diouf, M. A. and Dufour, J. M. Improved nonparametric inference for the mean of a bounded random variable with application to poverty measures. 2005. URL <http://web.hec.ca/scse/articles/Diouf.pdf>.
- 449 Dvoretzky, A., Kiefer, J., and Wolfowitz, J. Asymptotic minimax character of a sample distribution function and of the classical multinomial estimator. *Annals of Mathematical Statistics*, 27:642–669, 1956.
- 450 Efron, B. and Tibshirani, R. J. *An Introduction to the Bootstrap*. Chapman and Hall, London, 1993.
- 451 Fienberg, S. E., Neter, J., and Leitch, R. A. Estimating the total overstatement error in accounting populations. *Journal of the American Statistical Association*, 72(358): 295–302, 1977.
- 452 Frost, J. Statistics by Jim: Central limit theorem explained, January 2021. URL <https://statisticsbyjim.com/basics/central-limit-theorem/>.
- 453 Gaffke, N. Three test statistics for a nonparametric one-sided hypothesis on the mean of a nonnegative variable. *Mathematical Methods of Statistics*, 14(4):451–467, 2005.
- 454 Hoeffding, W. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- 455 Learned-Miller, E. and DeStefano, J. A probabilistic upper bound on differential entropy. *IEEE Transactions on Information Theory*, 54(11):5223–5230, 2008.
- 456 Learned-Miller, E. and Thomas, P. S. A new confidence interval for the mean of a bounded random variable. *arXiv preprint arXiv:1905.06208*, 2019.
- 457 Maurer, A. and Pontil, M. Empirical Bernstein bounds and sample variance penalization. In *Proceedings of the Twenty-Second Annual Conference on Learning Theory*, pp. 115–124, 2009.
- 458 Pap, G. and van Zuijlen, M. C. A. The Stringer bound in case of uniform taintings. *Computers and Mathematics with Applications*, 29(10):51–59, 1995.
- 459 Romano, J. P. and Wolf, M. Finite sample nonparametric inference and large sample efficiency. *Annals of Mathematical Statistics*, 28(3):756–778, 2000.
- 460 Serfling, R. *Approximation theorems of mathematical statistics*. Wiley series in probability and mathematical statistics : Probability and mathematical statistics. Wiley, New York, NY [u.a.], [nachdr.] edition, 1980. ISBN 0471024031. URL http://gso.gbv.de/DB=2.1/CMD?ACT=SRCHA&SRT=YOP&IKT=1016&TRM=ppn+024353353&sourceid=fbw_bibsonomy.
- 461 Stringer, K. W. Practical aspects of statistical sampling. *Proceedings of Business and Economic Statistics Section, American Statistical Association*, 1963.
- 462 Student. The probable error of a mean. *Biometrika*, pp. 1–25, 1908.

Supplementary Material: Towards Practical Mean Bounds for Small Samples

- In Section A we present more experiments.
- In Section B, as noted in Section 1, we show a log-normal distribution where the sample mean distribution is visibly skewed when $n = 80$.
- In Section C we present the proofs of Section 2.1.2.
- In Section D we discuss the Monte Carlo convergence result of our approximation in Section 3.
- In Section E.1 we show that our bound reduces to the Clopper-Pearson bound for binomial distributions as mentioned in Section 4.1. In Section E.2 we present the proofs of Section 4.2.

A. Other Experiments

In this section we perform experiments to find an upper bound of the mean of distributions given a finite upper bound of the support, or to a lower bound of the mean of distributions given a finite lower bound of the support. We find the lower bound of the mean of a random variable X by finding the upper mean bound of $-X$ and negating it to obtain the lower mean bound of X .

For each experiment we plot the following:

- The expected value of the bounds versus the sample size. For each sample size, we draw 10,000 samples of \mathbf{x} , compute the bound for each \mathbf{x} and compute the average.
- For the upper bound of the mean, we plot the α -quantile of the bound distribution versus the sample size. For each sample size, we draw 10,000 samples of \mathbf{x} , compute the bound for each \mathbf{x} and take the α quantile. If the α -quantile is below the true mean, the bound does not have guaranteed coverage.

For the lower bound of the mean, we plot the $1 - \alpha$ -quantile of the bound distribution versus the sample size. For each sample size, we draw 10,000 samples of \mathbf{x} , compute the bound for each \mathbf{x} and take the $1 - \alpha$ quantile. If the $1 - \alpha$ -quantile is above the true mean, the bound does not have guaranteed coverage.

- Coverage of the bounds. For each value of α from 0.02 to 1 with a step size of 0.02, we draw 10,000 samples of \mathbf{x} , compute the bound for each \mathbf{x} and plot the percentage of the bounds that are greater than or equal to the true mean (denoted *coverage*). If this percentage is larger than $1 - \alpha$, the bound has guaranteed coverage.

We perform the following experiments:

1. For the case in which we know a superset D^+ of the distribution's support with a finite lower bound and a finite upper bound, we compare the following bounds:

- Anderson's bound.
- New bound with T being Anderson's bound.
- Student's t .
- Hoeffding's bound.
- Maurer and Pontil's bound.

We find an upper bound of the mean for the following distributions:

- $\beta(1, 5)$, `uniform(0, 1)` and $\beta(5, 1)$. The known superset of the support is $[0, 1]$. The result is in Figure 6.
- $\beta(0.5, 0.5)$, $\beta(1, 1)$ and $\beta(2, 2)$. The known superset of the support is $[0, 1]$. The result is in Figure 7.
- `binomial(10, 0.1)`, `binomial(10, 0.5)` and `binomial(10, 0.9)`. The known superset of the support is the interval $[0, 10]$. The result is in Figure 8.

2. We also consider the case in which we want an upper bound of the mean without knowing the lower bound of the support (or to find a lower bound without knowing an upper bound of the support). Since Hoeffding's and Maurer and Pontil's bounds require knowing both a finite lower bound and upper bound, they are not applicable in this setting. We compare the following bounds:

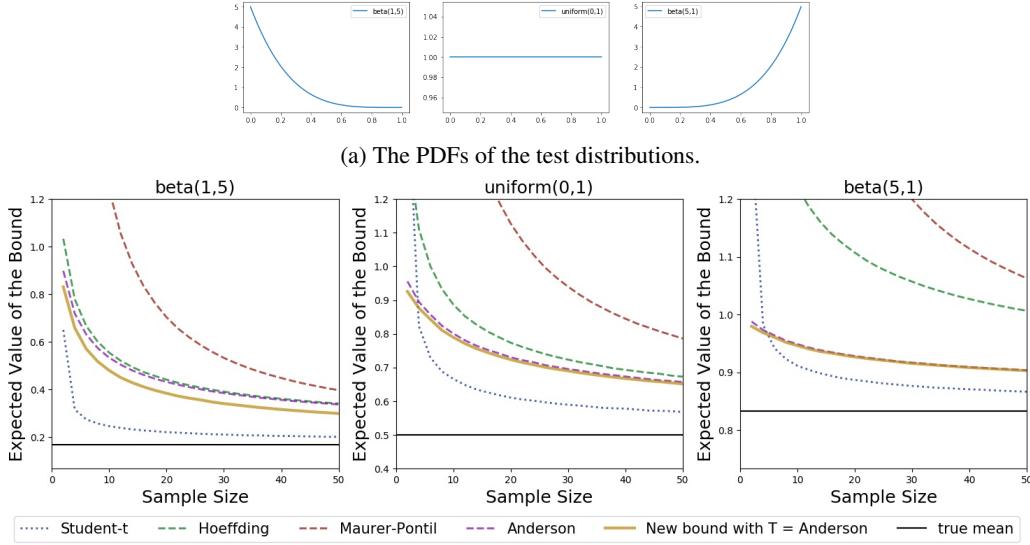
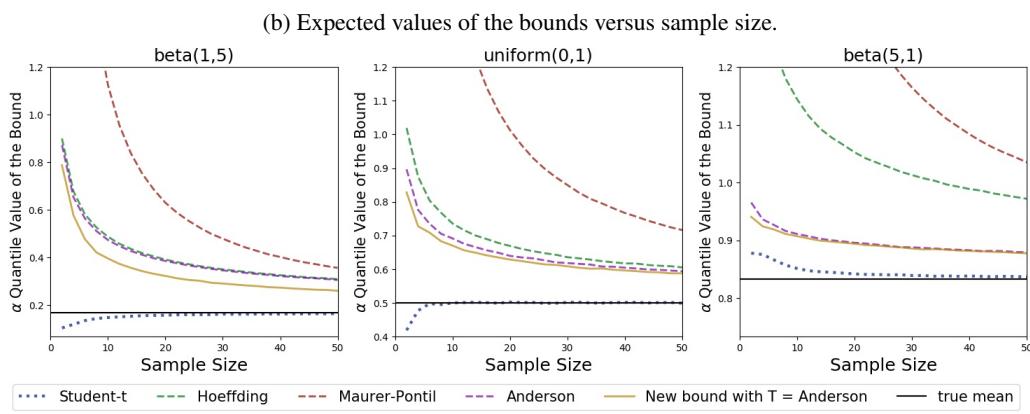
- Anderson's bound.
- New bound with T being Anderson's bound.
- Student's t

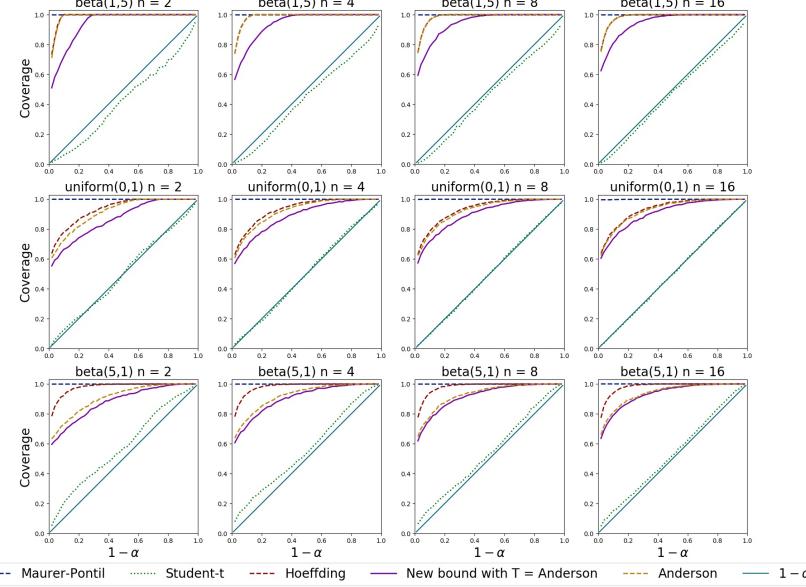
We address the following distributions:

- $\beta(1, 5)$, `uniform(0, 1)` and $\beta(5, 1)$. The known superset of the support is $(-\infty, 1]$. We find the upper bound of the mean. The result is in Figure 9.
- `binomial(10, 0.1)`, `binomial(10, 0.5)` and `binomial(10, 0.9)`. The known superset of the support is $(-\infty, 10]$. We find the upper bound of the mean. The result is in Figure 10.
- `poisson(2)`, `poisson(10)` and `poisson(50)`. The known superset of the support is $[0, \infty)$. We find the lower bound of the mean. The result is in Figure 11.

All the experiments confirm that our bound has guaranteed coverage and is tighter than Anderson's and Hoeffding's.

From the experiments, our upper bound performs the best in distributions that are skewed right (respectively, our lower bound will perform the best in distributions that are skewed left), when we know a tight lower bound and upper bound of the support.

550
 551
 552
 553
 554
 555

 556
 557
 558
 559
 560
 561
 562
 563
 564
 565

 566
 567
 568

 (c) The α -quantiles of bound distributions. If the α -quantile is below the true mean, the bound does not have guaranteed coverage.

 599
 600
 601
 602
 603
 604

 (d) The coverage of the bound. If the coverage is below the line $1 - \alpha$, the bound does not have guaranteed coverage.

 Figure 6. Finding the upper bound of the mean with $D^+ = [0, 1]$

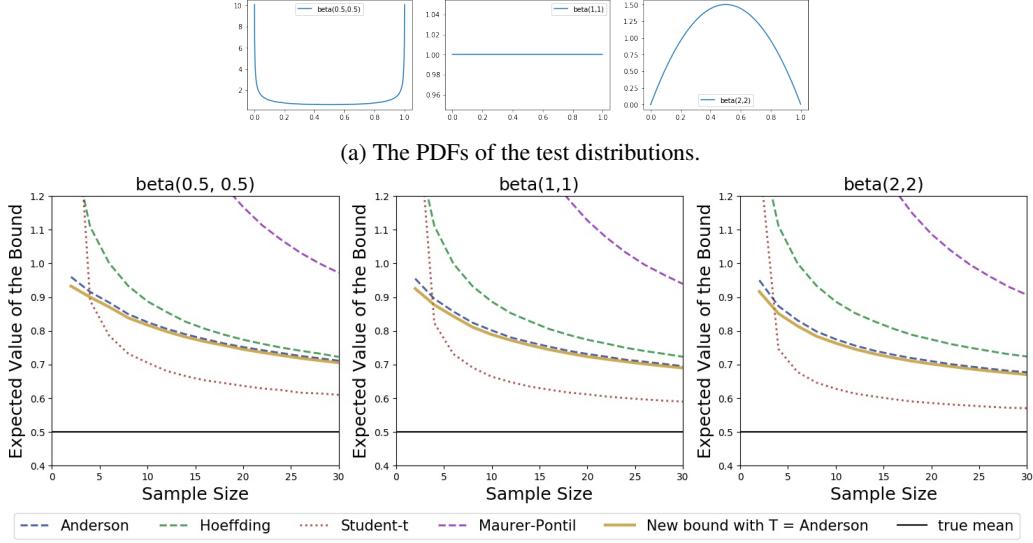


Figure 7. Finding the upper bound of the mean with $D^+ = [0, 1]$

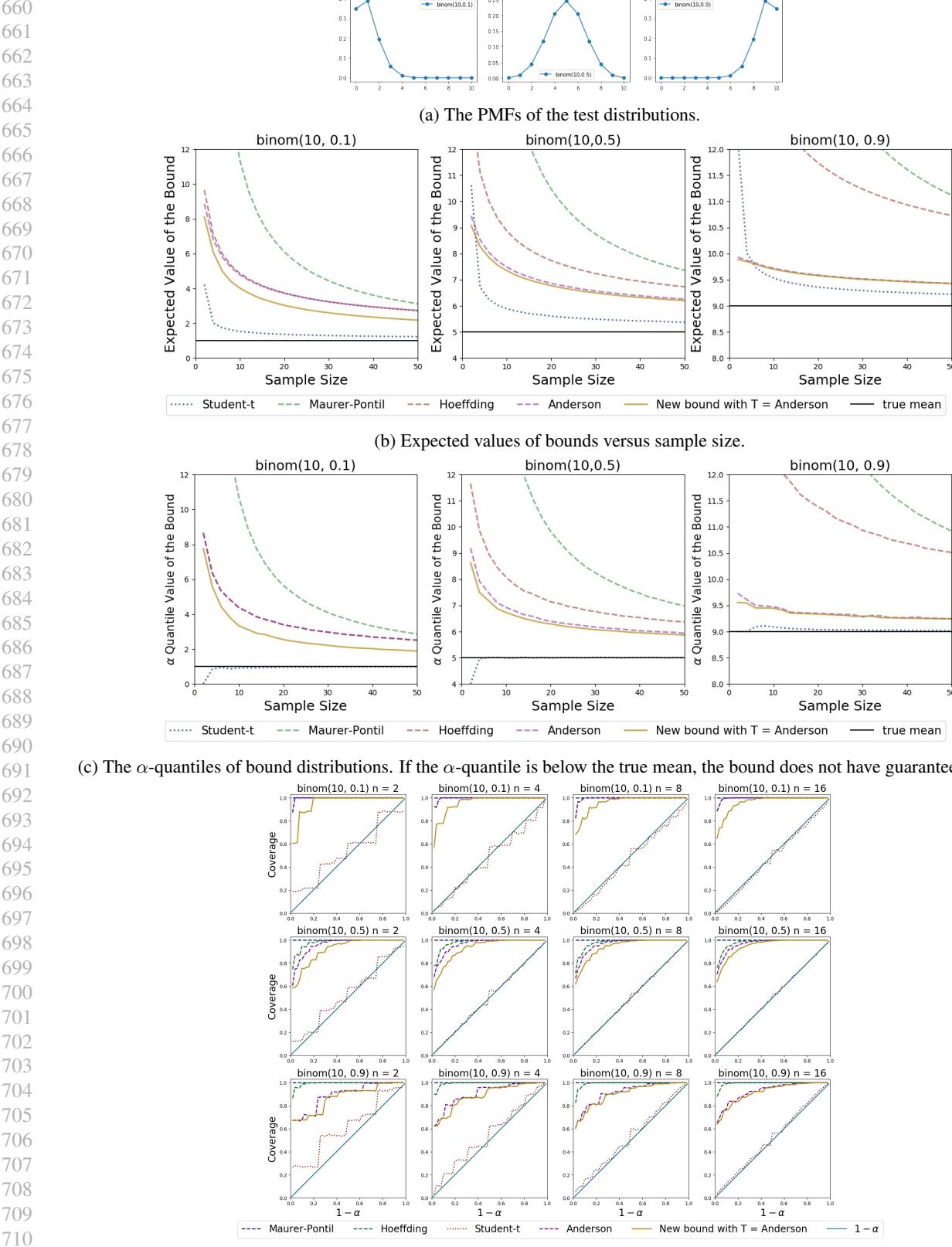


Figure 8. Finding the upper bound of the mean with $D^+ = [0, 10]$

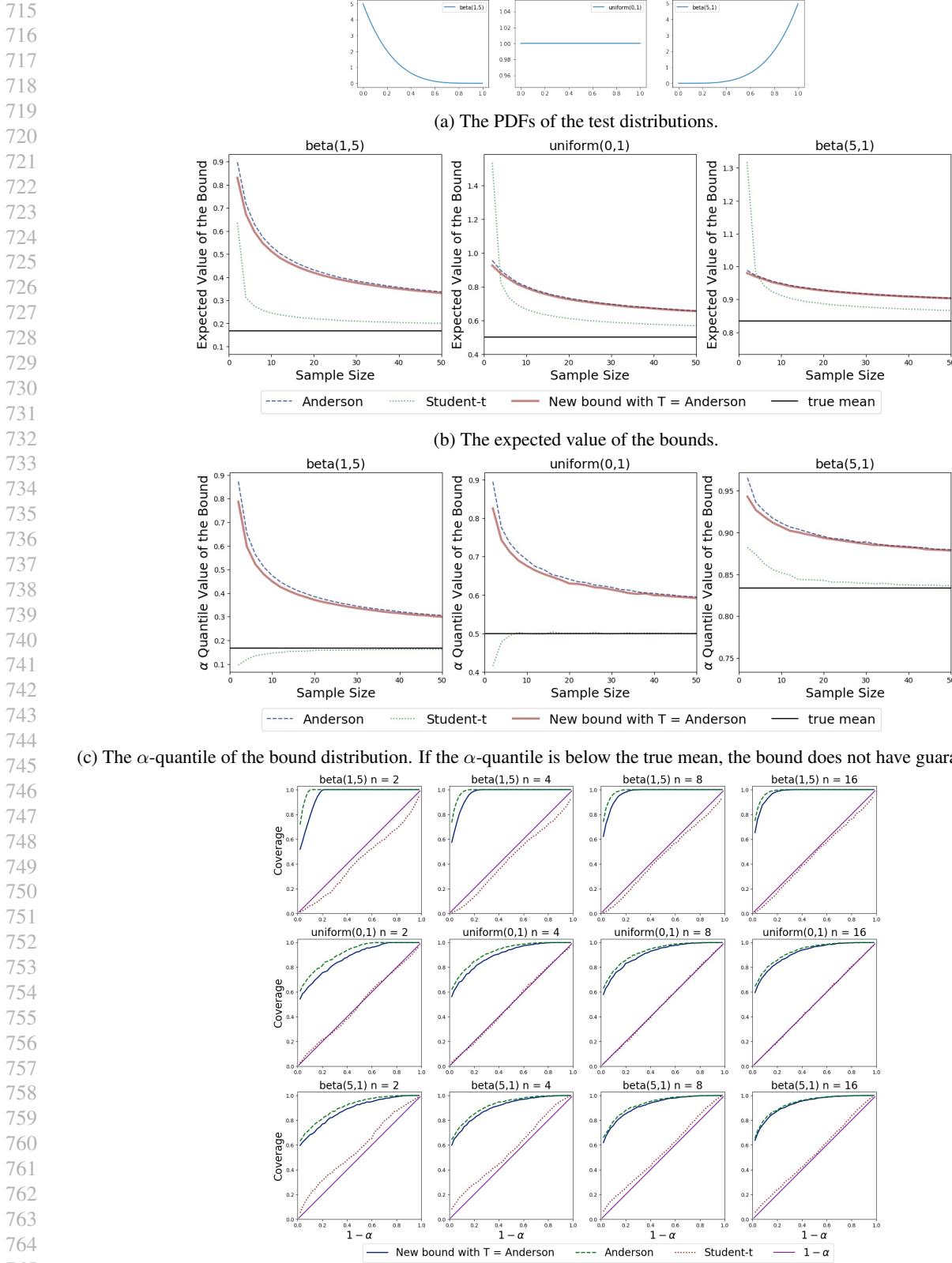


Figure 9. Finding the upper bound of the mean with $D^+ = (-\infty, 1]$

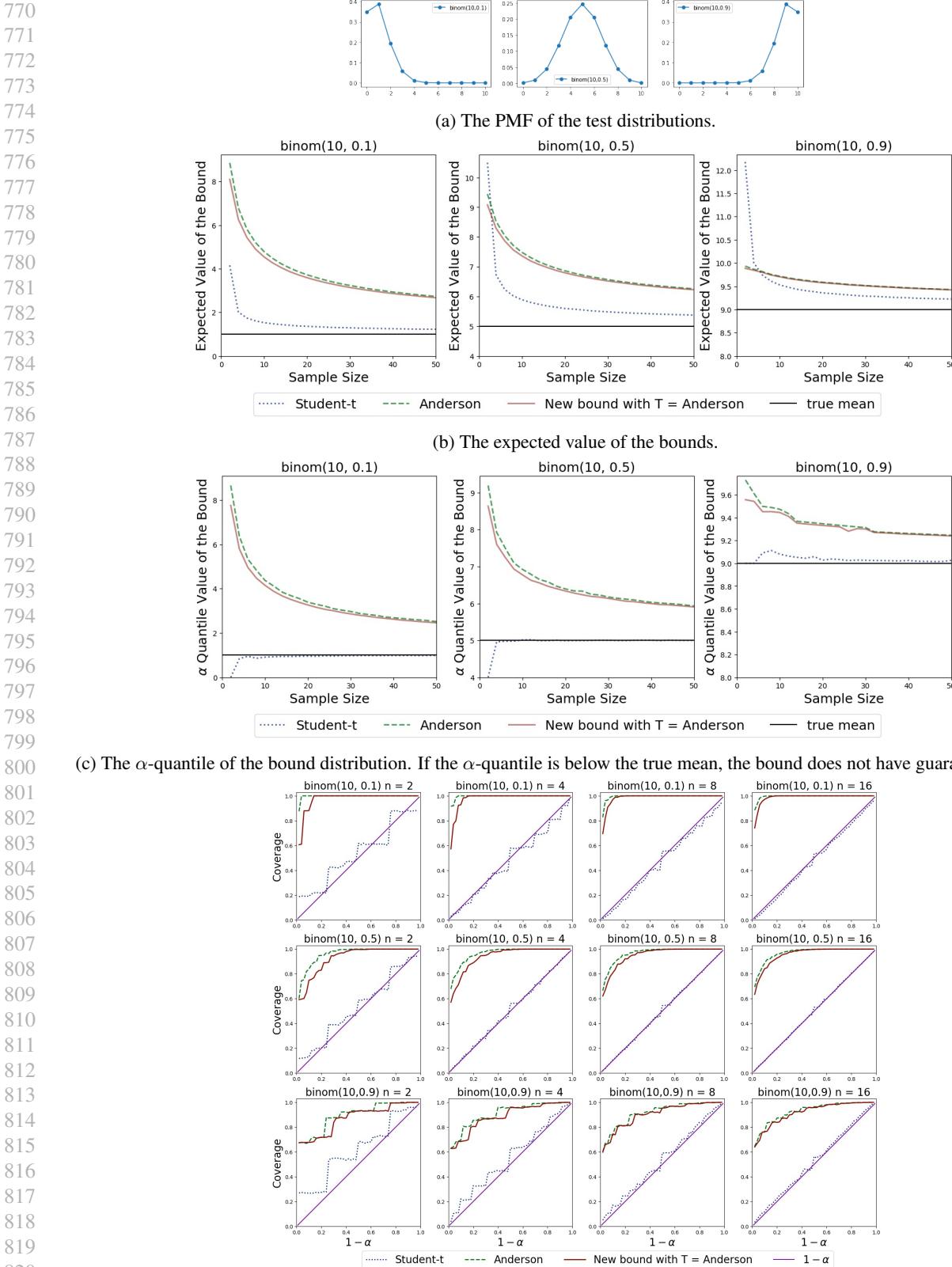


Figure 10. Finding the upper bound of the mean with $D^+ = (-\infty, 10]$

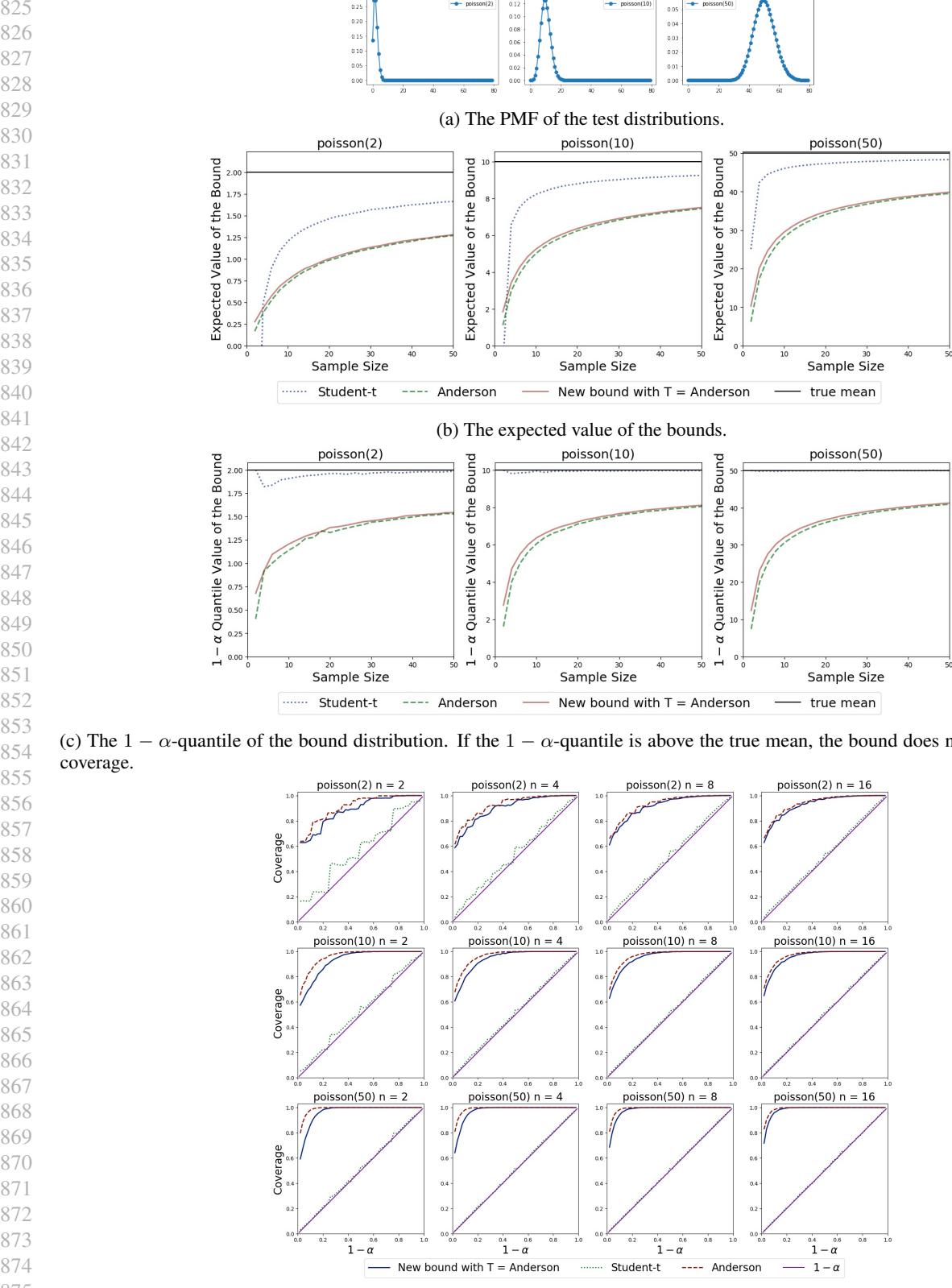


Figure 11. Finding the lower bound of the mean with $D^+ = [0, \infty)$

B. Discussion on Section 1: Skewed Sample Mean Distribution with $n = 80$

In this section, as noted in Section 1, we show a log-normal distribution where the sample mean distribution is visibly skewed when $n = 80$ (Figure 12). Student's t is not a good candidate in this case because the sample mean distribution is not approximately normal. This example is a variation on the one provided by Frost (2021).

While the log-normal distribution is an extreme example of skew, this example illustrates the danger of assuming the validity of arbitrary thresholds on the sample size, such as the traditional threshold of $n = 30$, for using the Student's t method. Clearly there are cases where such a threshold, and even much larger thresholds, are not adequate.

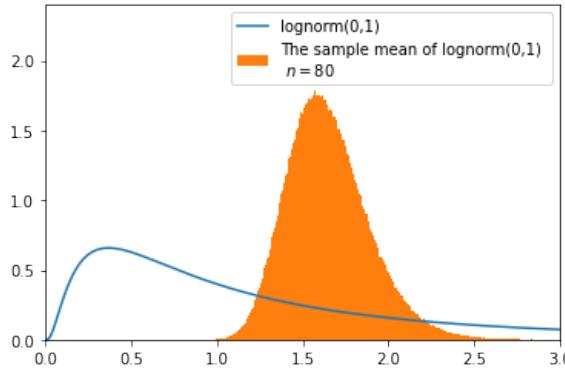


Figure 12. The PDFs of $\text{lognorm}(0, 1)$ and the sample mean distribution of $\text{lognorm}(0, 1)$. The sample mean distribution of $\text{lognorm}(0, 1)$ is visibly skewed when the sample size $n = 80$.

C. Proof of Section 2.1.2

We restate the lemma and theorem statements for convenience.

Lemma C.1 (Lemma 2.2). *Let X be a random variable with CDF F and $Y \stackrel{\text{def}}{=} F(X)$, known as the probability integral transform of X . Let U be a uniform random variable on $[0, 1]$. Then for any $0 \leq y \leq 1$,*

$$\mathbb{P}(Y \leq y) \leq \mathbb{P}(U \leq y). \quad (53)$$

If F is continuous, then Y is uniformly distributed on $(0, 1)$.

Proof. Since $\mathbb{P}(U \leq y) = y$, we will show that $\mathbb{P}(Y \leq y) \leq y$.

We will first show that if $F(x) \leq y$, then $x \leq \sup\{z : F(z) \leq y\}$. Suppose that $x > \sup\{z : F(z) \leq y\}$. Then, $F(x) > y$. Therefore,

$$F(x) \leq y \text{ implies } x \leq \sup\{z : F(z) \leq y\}. \quad (54)$$

Now we have

$$\mathbb{P}(Y \leq y) = \mathbb{P}(F(X) \leq y) \quad (55)$$

$$\leq \mathbb{P}(X \leq \sup\{z : F(z) \leq y\}) \quad (56)$$

$$= F(z^*) \text{ where } z^* = \sup\{z : F(z) \leq y\} \quad (57)$$

$$\leq y. \quad (58)$$

If F is continuous, (Angus, 1994) shows that Y is uniformly distributed on $(0, 1)$. \square

Lemma C.2 (Lemma 2.3). *For any $\mathbf{x} \in D^n$,*

$$m_D(\mathbf{x}, F(\mathbf{x})) \geq \mu. \quad (59)$$

Proof. We have

$$m_D(\mathbf{x}, F(\mathbf{x})) = \sum_{i=1}^{n+1} x_{(i)}(F(x_{(i)}) - F(x_{(i-1)})) \quad (60)$$

where $x_{(0)} \stackrel{\text{def}}{=} -\infty$ and $x_{(n+1)} \stackrel{\text{def}}{=} s_D$. Then:

$$\mu = \int_x x dF(x) \quad (61)$$

$$= \sum_{i=1}^{n+1} \int_{x=x_{(i-1)}}^{x_{(i)}} x dF(x) \quad (62)$$

$$\leq \sum_{i=1}^{n+1} \int_{x=x_{(i-1)}}^{x_{(i)}} x_{(i)} dF(x) \quad (63)$$

$$= \sum_{i=1}^{n+1} x_{(i)}(F(x_{(i)}) - F(x_{(i-1)})) \quad (64)$$

$$= m_D(\mathbf{x}, F(\mathbf{x})). \quad (65)$$

\square

Lemma C.3 (Lemma 2.4). *Let \mathbf{Z} be a random sample of size n from F . Let $\mathbf{U} = U_1, \dots, U_n$ be a sample of size n from the continuous uniform distribution on $[0, 1]$. For any function $T : D^n \rightarrow R$ and any $\mathbf{x} \in D^n$:*

$$\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{x})) \leq \mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu). \quad (66)$$

Proof. Let \cup denote the union of events and $\{\}$ denote an event. Let \mathbf{Z} be a sample from F . Then for any sample \mathbf{x} :

$$\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{x})) \quad (67)$$

$$= \mathbb{P}_{\mathbf{Z}}(\mathbf{Z} \in \mathbb{S}(\mathbf{x})) \quad (68)$$

$$= \mathbb{P}_{\mathbf{Z}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{\mathbf{Z} = \mathbf{y}\}) \quad (69)$$

$$\leq \mathbb{P}_{\mathbf{Z}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \mathbf{Z} \preceq \mathbf{y}) \text{ because } \mathbf{Z} = \mathbf{y} \text{ implies } \mathbf{Z} \preceq \mathbf{y} \quad (70)$$

$$\leq \mathbb{P}_{\mathbf{Z}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} \{F(\mathbf{Z}) \preceq F(\mathbf{y})\}) \quad (71)$$

because F is non-decreasing, so $Z_{(i)} \leq y_{(i)}$ implies $F(Z_{(i)}) \leq F(y_{(i)})$. Let U_1, \dots, U_n be n samples from the uniform distribution on $(0, 1)$. From Lemma 2.2, for any $u \in (0, 1)$, $\mathbb{P}(F(Z_i) \leq u) \leq \mathbb{P}(U_i \leq u)$. Therefore:

$$\mathbb{P}_{\mathbf{Z}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})}\{F(\mathbf{Z}) \leq F(\mathbf{y})\}) \quad (72)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})}\{\mathbf{U} \leq F(\mathbf{y})\}). \quad (73)$$

Recall that $m_D(\mathbf{y}, \mathbf{U}) = s_D - \sum_{i=1}^n U_{(i)}(y_{(i+1)} - y_{(i)})$ where $\forall i, y_{(i+1)} - y_{(i)} \geq 0$. Therefore if $\forall i, U_{(i)} \leq F(y_{(i)})$ then $m_D(\mathbf{y}, \mathbf{U}) \geq m_D(\mathbf{y}, F(\mathbf{y}))$:

$$\mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})}\{\mathbf{U} \leq F(\mathbf{y})\}) \quad (74)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})}\{m_D(\mathbf{y}, \mathbf{U}) \geq m_D(\mathbf{y}, F(\mathbf{y}))\}) \quad (75)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})}\{m_D(\mathbf{y}, \mathbf{U}) \geq \mu\}), \text{ by Lemma 2.3} \quad (76)$$

$$\leq \mathbb{P}_{\mathbf{U}}(\sup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} m_D(\mathbf{y}, \mathbf{U}) \geq \mu) \quad (77)$$

$$= \mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu) \quad (78)$$

The inequality in Eq. 77 is because if there exists $\mathbf{y} \in \mathbb{S}(\mathbf{x})$ such that $m_D(\mathbf{y}, \mathbf{U}) \geq \mu$, then $\sup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} m_D(\mathbf{y}, \mathbf{U}) \geq \mu$. Therefore the event $\cup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})}\{m_D(\mathbf{y}, \mathbf{U}) \geq \mu\}$ is a subset of the event $\sup_{\mathbf{y} \in \mathbb{S}(\mathbf{x})} m_D(\mathbf{y}, \mathbf{U}) \geq \mu$, and Eq. 77 follows.

From Eqs. 71, 73 and Eq. 78:

$$\mathbb{P}_{\mathbf{Z}}(T(\mathbf{Z}) \leq T(\mathbf{x})) \quad (79)$$

$$\leq \mathbb{P}_{\mathbf{U}}(b(\mathbf{x}, \mathbf{U}) \geq \mu). \quad (80)$$

□

D. Discussion on Section 3: Monte Carlo Convergence

In Section 3, we discuss the use of Monte Carlo sampling of the induced mean function $m(\mathbf{z}, U)$, via sampling of the uniform random variable U , to approximate the $1 - \alpha$ quantile of $m(\mathbf{z}, U)$.

The Monte Carlo approximation error is quantified in the following lemma due to Serfling (1980). Let $F(m-) \stackrel{\text{def}}{=} \lim_{x \rightarrow m-} F(x)$.

Lemma D.1 (Theorem 2.3.2 in Serfling (1980)). *Let $0 < p < 1$. If $Q(p, M)$ is the unique solution m of $F(m-) \leq p \leq F(m)$, then for every $\epsilon > 0$,*

$$\mathbb{P}(|M_{[pl]} - Q(p, M)| > \epsilon) \leq 2e^{-2l\delta}, \quad (81)$$

where

$$\delta = \min(p - F(Q(p, M) - \epsilon), F(Q(p, M) + \epsilon) - p).$$

Note that when the condition that $Q(p, M)$ is the unique solution m of $F(m-) \leq p \leq F(m)$ is satisfied, $\delta > 0$. Let $M \stackrel{\text{def}}{=} b_{D, T}(\mathbf{x}, \mathbf{U}) \in [0, 1]$. In Lemma D.2 we will show

that the CDF of M satisfies the condition in Lemma D.1. Therefore the error incurred by computing the bound via Monte Carlo sampling can be decreased to an arbitrarily small value by choosing a large enough number of Monte Carlo samples l . The Monte Carlo estimation of $b_{D+, T}^\alpha(\mathbf{x})$ where $D^+ = [0, 1]$ is presented in Algorithm 1.

Lemma D.2. *Let $M \stackrel{\text{def}}{=} b(\mathbf{x}, \mathbf{U})$. Let F_M be the CDF of M .*

For any \mathbf{x} , for any scalar function T , either:

1. *M is a constant, or*
2. *For any $p \in (0, 1)$, $Q(p, M)$ is the unique solution m of $F_M(m-) \leq p \leq F_M(m)$.*

Proof. We will show that for any \mathbf{x} , for any T , for any $p \in (0, 1)$, $F_M(m-) \leq p \leq F_M(m)$ has a unique solution by showing that for any \mathbf{x} and T , F_M is strictly increasing on its support. To do so, for any c_1, c_2 in the support such that $c_1 < c_2$ we will show that

$$F_M(c_2) - F_M(c_1) > 0. \quad (82)$$

Recall the definition of the induced mean as

$$b(\mathbf{x}, \ell) = \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} \sum_{i=1}^{n+1} z_{(i)}(\ell_{(i)} - \ell_{(i-1)}), \quad (83)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - \sum_{i=1}^n \ell_{(i)}(z_{(i+1)} - z_{(i)}), \quad (84)$$

where $\ell_{(0)} \stackrel{\text{def}}{=} 0$, $\ell_{(n+1)} \stackrel{\text{def}}{=} 1$ and $z_{(n+1)} \stackrel{\text{def}}{=} s_D$.

We now find the support of M . Let $\phi \stackrel{\text{def}}{=} \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} z_{(1)}$. We will show that for any \mathbf{u} where $0 \leq u_i \leq 1$, we have $\phi \leq b(\mathbf{x}, \mathbf{u}) \leq s_D$, and therefore the support of M is a subset of $[\phi, s_D]$. We have

$$b(\mathbf{x}, \mathbf{u}) = \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - \sum_{i=1}^n u_{(i)}(z_{(i+1)} - z_{(i)}) \quad (85)$$

$$\leq \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - \sum_{i=1}^n 0(z_{(i+1)} - z_{(i)}) \quad (86)$$

$$= s_D. \quad (87)$$

990 Similarly we have

$$b(\mathbf{x}, \mathbf{u}) = \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - \sum_{i=1}^n u_{(i)}(z_{(i+1)} - z_{(i)}) \quad (88)$$

$$\geq \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - \sum_{i=1}^n 1(z_{(i+1)} - z_{(i)}) \quad (89)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - (z_{(n+1)} - z_{(1)}) \quad (90)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} z_{(1)} \quad (91)$$

$$= \phi. \quad (92)$$

Therefore $M = b(\mathbf{x}, \mathbf{U}) \in [\phi, s_D]$. We consider two cases: where $\phi = s_D$ and where $\phi < s_D$.

Case 1: $\phi = s_D$.

Then for all i , $\sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} z_{(i)} \geq \phi = s_D$. Since $z_{(i)} \leq s_D$, we have for all i , $\sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} z_{(i)} = s_D$. Therefore

$$b(\mathbf{x}, \ell) = \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} \sum_{i=1}^{n+1} z_{(i)}(\ell_{(i)} - \ell_{(i-1)}) \quad (93)$$

$$= \sum_{i=1}^{n+1} s_D(\ell_{(i)} - \ell_{(i-1)}) \quad (94)$$

$$= s_D. \quad (95)$$

Therefore $M = b(\mathbf{x}, \mathbf{U})$ is a constant s_D , and the $1 - \alpha$ quantile of M is s_D .

Case 2: $\phi < s_D$.

Let $c_1, c_2 \in \mathcal{R}$ be such that $\phi \leq c_1 < c_2 \leq s_D$. We will now show that

$$F_M(c_2) - F_M(c_1) > 0. \quad (96)$$

Let $v \stackrel{\text{def}}{=} \frac{s_D - c_2}{s_D - \phi}$ and $w \stackrel{\text{def}}{=} \frac{s_D - c_1}{s_D - \phi}$. If $\phi \leq c_1 < c_2 \leq s_D$ then $v < w$ and $v, w \in [0, 1]$.

Let $\mathbf{v} \stackrel{\text{def}}{=} (v_1, \dots, v_n)$ and $\mathbf{w} \stackrel{\text{def}}{=} (w_1, \dots, w_n)$ where $\forall i, v_i = v$ and $w_i = w$. Then

$$b(\mathbf{x}, \mathbf{v}) \quad (97)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} \sum_{i=1}^{n+1} z_{(i)}(v_{(i)} - v_{(i-1)}) \quad (98)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} z_{(n+1)}(v_{(n+1)} - v_{(n)}) + z_{(1)}(v_{(1)} - v_{(0)}) \quad (99)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D(1 - v) + z_{(1)}(v - 0) \quad (100)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - (s_D - z_{(1)}) \frac{s_D - c_2}{s_D - \phi} \quad (101)$$

$$= s_D - (s_D - \phi) \frac{s_D - c_2}{s_D - \phi} \text{ because } \frac{s_D - c_2}{s_D - \phi} \geq 0 \quad (102)$$

$$= c_2. \quad (103)$$

Similarly,

$$b(\mathbf{x}, \mathbf{w}) \quad (104)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} \sum_{i=1}^{n+1} z_{(i)}(w_{(i)} - w_{(i-1)}) \quad (105)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} z_{(n+1)}(w_{(n+1)} - w_{(n)}) + z_{(1)}(w_{(1)} - w_{(0)}) \quad (106)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D(1 - w) + z_{(1)}(w - 0) \quad (107)$$

$$= \sup_{\mathbf{z} \in \mathbb{S}(\mathbf{x})} s_D - (s_D - z_{(1)}) \frac{s_D - c_1}{s_D - \phi} \quad (108)$$

$$= s_D - (s_D - \phi) \frac{s_D - c_1}{s_D - \phi} \text{ because } \frac{s_D - c_1}{s_D - \phi} \geq 0 \quad (109)$$

$$= c_1. \quad (110)$$

Since $b(\mathbf{x}, \mathbf{u})$ is constructed from linear function of \mathbf{u} with non-positive coefficients, for any \mathbf{u} such that $v \leq u_{(1)} \leq \dots \leq u_{(n)} < w$ we have:

$$b(\mathbf{x}, \mathbf{w}) < b(\mathbf{x}, \mathbf{u}) \leq b(\mathbf{x}, \mathbf{v}), \quad (111)$$

which is equivalent to:

$$c_1 < b(\mathbf{x}, \mathbf{u}) \leq c_2. \quad (112)$$

So we have $v \leq u_{(1)} \leq \dots \leq u_{(n)} < w$ implies $c_1 < b(\mathbf{x}, \mathbf{u}) \leq c_2$. Therefore for any c_1, c_2 such that $\phi \leq c_1 < c_2 \leq s_D$:

$$F_M(c_2) - F_M(c_1) \quad (113)$$

$$= \mathbb{P}(c_1 < M \leq c_2) \quad (114)$$

$$= \mathbb{P}_{\mathbf{U}}(c_1 < b(\mathbf{x}, \mathbf{U}) \leq c_2) \quad (115)$$

$$\geq \mathbb{P}_{\mathbf{U}}(v \leq U_{(1)} \leq \dots \leq U_{(n)} < w) \quad (116)$$

$$> 0 \text{ because } 0 \leq v < w \leq 1. \quad (117)$$

Since the support of M is in $[\phi, s_D]$ we have that F_M is strictly increasing on the support. \square

In summary, the Monte Carlo estimate of our bound will converge to the correct value as the number of samples grows, allowing the user to make the error arbitrarily small.

E. Discussion on Section 4

E.1. Special Case: Bernoulli Distribution

When we know that $D = \{0, 1\}$, the distribution is Bernoulli. If we choose T to be the sample mean, we will show that our bound becomes the same as the Clopper-Pearson confidence bound for binomial distributions (Clopper & Pearson, 1934).

1045 If $\mathbf{x}, \mathbf{z} \in \{0, 1\}^n$ and $T(\mathbf{z}) \leq T(\mathbf{x})$ then $m(\mathbf{z}, \mathbf{u}) \leq$
 1046 $m(\mathbf{x}, \mathbf{u})$. Therefore for any $\mathbf{u} \in [0, 1]^n$,

$$1048 b_{D,T}(\mathbf{x}, \mathbf{u}) = \sup_{\mathbf{z} \in \{0,1\}^n : T(\mathbf{z}) \leq T(\mathbf{x})} m_D(\mathbf{z}, \mathbf{u}) = m_D(\mathbf{x}, \mathbf{u}). \quad (118)$$

1050 Let $p_{\mathbf{x}}$ be the number of 0's in \mathbf{x} . Therefore the bound
 1051 becomes the $1 - \alpha$ quantile of $m_D(\mathbf{x}, \mathbf{U})$ where

$$1054 m_D(\mathbf{x}, \mathbf{U}) = 1 - \sum_{i=1}^n U_{(i)}(x_{(i+1)} - x_{(i)}) = 1 - U_{(p_{\mathbf{x}})}. \quad (119)$$

1055 Therefore the bound is the $1 - \alpha$ quantile of $1 - U_{(p_{\mathbf{x}})}$. Then

$$1060 \mathbb{P}(U_{(p_{\mathbf{x}})} \leq 1 - b^\alpha(\mathbf{x})) = \mathbb{P}(1 - U_{(p_{\mathbf{x}})} \geq b^\alpha(\mathbf{x})) = \alpha. \quad (120)$$

1061 Let $\beta(i, j)$ denote a beta distribution with parameters i and
 1062 j . We use the fact that the order statistics of a uniform
 1063 distribution are beta-distributed. Since $U_{(p_{\mathbf{x}})} \sim \beta(p_{\mathbf{x}}, n +$
 1064 $1 - p_{\mathbf{x}})$, we have $1 - U_{(p_{\mathbf{x}})} \sim \beta(n - p_{\mathbf{x}} + 1, p_{\mathbf{x}})$

$$1065 b^\alpha(\mathbf{x}) = Q(1 - \alpha, \beta(n - p_{\mathbf{x}} + 1, p_{\mathbf{x}})). \quad (121)$$

1066 This is the same as the Clopper-Pearson upper confidence
 1067 interval for binomial distributions.

E.2. Proof of Section 4.2

1073 **Lemma E.1** (Lemma 4.1). *Let $\mathbf{X} = (X_1, \dots, X_n)$ be a
 1074 sample of size n from a distribution with mean μ . Let $\ell \in$
 1075 $[0, 1]^n$. If $G_{\mathbf{X}, \ell}$ is a $(1 - \alpha)$ lower confidence bound for the
 1076 CDF then*

$$1078 \mathbb{P}_{\mathbf{X}}(m(\mathbf{X}, \ell) \geq \mu) \geq 1 - \alpha. \quad (122)$$

1080 *Proof.* If $\forall y \in \mathcal{R}, F(y) \geq G_{\mathbf{X}, \ell}(y)$ then

$$1082 \forall i : 1 \leq i \leq n, F(X_{(i)}) \geq \ell_{(i)}. \quad (123)$$

1083 Recall that $m_D(\mathbf{X}, \ell) = s_D - \sum_{i=1}^n \ell_{(i)}(z_{(i+1)} - z_{(i)})$.
 1084 Therefore if $\forall i : 1 \leq i \leq n, F(X_{(i)}) \geq \ell_{(i)}$ then
 1085 $m(\mathbf{X}, \ell) \geq m(\mathbf{X}, F(\mathbf{X}))$.

1086 From Lemma 2.3, $m(\mathbf{X}, F(\mathbf{X})) \geq \mu$. Therefore
 1087 $m(\mathbf{X}, \ell) \geq \mu$. And hence, finally,

$$1090 \mathbb{P}(m(\mathbf{X}, \ell) \geq \mu) \geq \mathbb{P}_{\mathbf{X}}(\forall y \in \mathcal{R}, F(y) \geq G_{\mathbf{X}, \ell}(y)) \quad (124)$$

$$1092 = 1 - \alpha. \quad (125)$$

1094 We now show that if $G_{\mathbf{X}, \ell}$ (Figure 1) is a lower confidence
 1095 bound, then the order statistics of ℓ are element-wise smaller
 1096 than the order statistics of a sample of size n from the
 1097 uniform distribution with high probability:
 1098

Lemma E.2. *Let $\mathbf{U} = U_1, \dots, U_n$ be a sample of size n from
 the continuous uniform distribution on $[0, 1]$. Let $\ell \in [0, 1]^n$
 and $\alpha \in (0, 1)$. If $G_{\mathbf{X}, \ell}$ is a $(1 - \alpha)$ lower confidence bound
 for the CDF then:*

$$\mathbb{P}_{\mathbf{U}}(\forall i : 1 \leq i \leq n, U_{(i)} \geq \ell_{(i)}) \geq 1 - \alpha. \quad (126)$$

Proof. Let K be the CDF of a distribution such that K is continuous and strictly increasing (such as the standard Gaussian distribution). Let $\mathbf{X} = (X_1, \dots, X_n)$ be a sample of size n from the distribution with CDF K . By Lemma 2.2, $K(X)$ is uniformly distributed on $[0, 1]$.

By the definition of $G_{\mathbf{X}, \ell}$, if $\forall x \in C, K(x) \geq G_{\mathbf{X}, \ell}(x)$ then:

$$K(y) \geq 0, \quad \text{if } y < X_{(1)} \quad (127)$$

$$K(y) \geq \ell_{(i)}, \quad \text{if } X_{(i)} \leq y < X_{(i+1)} \quad (128)$$

$$K(y) \geq 1, \quad \text{if } y \geq s_C. \quad (129)$$

which is equivalent to:

$$\forall i : 1 \leq i \leq n, K(y) \geq \ell_{(i)}, \text{ if } X_{(i)} \leq y < X_{(i+1)} \quad (130)$$

Since $K(y)$ is non-decreasing, this is equivalent to:

$$\forall i : 1 \leq i \leq n, K(X_{(i)}) \geq \ell_{(i)} \quad (131)$$

Since $G_{\mathbf{X}, \ell}$ is a lower confidence bound,

$$1 - \alpha \leq \mathbb{P}_{\mathbf{X}}(\forall y \in \mathcal{R}, K(y) \geq G_{\mathbf{X}, \ell}(y)) \quad (132)$$

$$= \mathbb{P}_{\mathbf{X}}(\forall i : 1 \leq i \leq n, K(X_{(i)}) \geq \ell_{(i)}) \quad (133)$$

$$= \mathbb{P}_{\mathbf{U}}(\forall i : 1 \leq i \leq n, U_{(i)} \geq \ell_{(i)}) \text{ by Lemma 2.2.} \quad (134)$$

□

To prove Theorem 4.3, we prove the more general version.
 In the following theorem, we compare $b_{D+, T}^\alpha$ to $m_{D+}(\mathbf{x}, \ell)$
 if $G_{\mathbf{X}, \ell}$ is a lower confidence bound for the CDF.

Theorem E.3. *Let $\ell \in [0, 1]^n$. Let $D^+ = [-\infty, b]$. If $G_{\mathbf{X}, \ell}$ is a $1 - \alpha$ lower confidence bound for the CDF, then for any sample size n , for all sample values $\mathbf{x} \in D^n$ and all $\alpha \in (0, 1)$, using $T(\mathbf{x}) = m_{D+}(\mathbf{x}, \ell)$ to compute $b_{D+, T}^\alpha(\mathbf{x})$ yields:*

$$b_{D+, T}^\alpha(\mathbf{x}) \leq m_{D+}(\mathbf{x}, \ell). \quad (135)$$

Proof. Since $G_{\mathbf{X}, \ell}$ is a lower confidence bound for the cdf F , from Lemma E.2,

$$\mathbb{P}(\forall i, U_{(i)} \geq \ell_{(i)}) \geq 1 - \alpha. \quad (136)$$

1100 First we note that

$$1102 \quad b_{D^+,T}(\mathbf{x}, \boldsymbol{\ell}) = \sup_{\mathbf{y}: \mathbf{y} \in \mathbb{S}_{D^+,T}(\mathbf{x})} m_{D^+}(\mathbf{y}, \boldsymbol{\ell}) \quad (137)$$

$$1104 \quad = \sup_{m_{D^+}(\mathbf{y}, \boldsymbol{\ell}) \leq m_{D^+}(\mathbf{x}, \boldsymbol{\ell})} m_{D^+}(\mathbf{y}, \boldsymbol{\ell}) \quad (138)$$

$$1106 \quad = m_{D^+}(\mathbf{x}, \boldsymbol{\ell}). \quad (139)$$

1107

1108 Recall that $b_{D^+,T}^\alpha(\mathbf{x})$ is the $1 - \alpha$ quantile of $b_{D^+,T}(\mathbf{x}, \mathbf{U})$.

1109 In order to show that $b_{D^+,T}^\alpha(\mathbf{x}) \leq b_{D^+,T}(\mathbf{x}, \boldsymbol{\ell})$, we will
1110 show that

$$1112 \quad \mathbb{P}(b_{D^+,T}(\mathbf{x}, \mathbf{U}) \leq b_{D^+,T}(\mathbf{x}, \boldsymbol{\ell})) \geq 1 - \alpha. \quad (140)$$

1113 Recall that $b_{D^+,T}(\mathbf{x}, \mathbf{U}) = \sup_{\mathbf{y} \in \mathbb{S}_T(\mathbf{x})} t_{D^+} -$
1114 $\sum_{i=1}^n U_{(i)}(x_{(i+1)} - x_{(i)})$. Then if $\forall i, U_{(i)} \geq \ell_{(i)}$
1115 then $b_{D^+,T}(\mathbf{x}, \mathbf{U}) \leq b_{D^+,T}(\mathbf{x}, \boldsymbol{\ell})$. Therefore,

1116

$$1118 \quad \mathbb{P}(b_{D^+,T}(\mathbf{x}, \mathbf{U}) \leq b_{D^+,T}(\mathbf{x}, \boldsymbol{\ell})) \quad (141)$$

$$1120 \quad \geq \mathbb{P}(\forall i, U_{(i)} \geq \ell_{(i)}) \quad (142)$$

$$1121 \quad \geq 1 - \alpha, \text{ by Lemma E.2.} \quad (143)$$

1122

□

1123

1124

1125 We can now show the comparison with Anderson's bound
1126 and Hoeffding's bound.

1127 **Theorem E.4** (Theorem 4.3). *Let $D^+ = (-\infty, b]$. For
1128 any sample size n , for any sample value $\mathbf{x} \in D^n$, for all
1129 $\alpha \in (0, 1)$, using $T(\mathbf{x}) = b^{\alpha, \text{Anderson}}(\mathbf{x})$ yields:*

$$1131 \quad b_{D^+,T}^\alpha(\mathbf{x}) \leq b^{\alpha, \text{Anderson}}(\mathbf{x}). \quad (144)$$

1132

1133 *Proof.* We have $b^{\alpha, \text{Anderson}}(\mathbf{x}) = m_{D^+}(\mathbf{x}, \mathbf{u}^{\text{And}})$ where
1134 \mathbf{u}^{And} satisfies $G_{\mathbf{X}, \mathbf{u}^{\text{And}}}$ is a $1 - \alpha$ lower confidence bound
1135 for the CDF. Therefore applying Theorem E.3 with $\boldsymbol{\ell} =$
1136 \mathbf{u}^{And} yields the result. □

1137 **Theorem E.5** (Theorem 4.4). *Let $D^+ = (-\infty, b]$. For
1138 any sample size n , for any sample value $\mathbf{x} \in D^n$, for all
1139 $\alpha \in (0, 1)$, using $T(\mathbf{x}) = b^{\alpha, \text{Hoeffding}}(\mathbf{x})$ yields:*

$$1142 \quad b_{D^+,T}^\alpha(\mathbf{x}) \leq b^{\alpha, \text{Hoeffding}}(\mathbf{x}). \quad (145)$$

1143

1144 *Proof.* The proof follows directly from Lemma 4.2 and
1145 Theorem 4.3. □

1146

1147

1148

1149

1150

1151

1152

1153

1154