

Statistics 583, Problem Set 4 Solutions

Wellner; 5/14/2015

1. (a) Given n distinct data items, show that the probability that a given data item does not appear in a bootstrap sample is $e_n = (1 - 1/n)^n$
- (b) Show that $e_n \rightarrow e^{-1} \approx .368$ as $n \rightarrow \infty$.
- (c) Hence show that the probability that each of B bootstrap samples contains an item i is $(1 - e_n)^B$. Evaluate this quantity for $n = 10, 20, 50, 100$ and $B = 10, 20, 50, 100$.
- (d) Let $N_n \equiv \sum_{j=1}^n 1_{[M_j=0]}$ where $\underline{M} \equiv (M_1, \dots, M_n) \sim \text{Mult}_n(n, \underline{1}/n)$. Show that $E(n^{-1}N_n) = e_n$ as computed in (a).

Solution: (a) The probability that X_i does not appear in a bootstrap sample X_1^*, \dots, X_n^* from \mathbb{F}_n is just $e_n = P(M_i = 0)$ where $M_i \sim \text{Binomial}(n, 1/n)$. Thus we have $e_n = P(M_i = 0) = \binom{n}{0}(1/n)^0(1 - 1/n)^n = (1 - 1/n)^n$.

(b) Since $(1 + x/n)^n \rightarrow e^x$ for any x , it follows immediately that $e_n \rightarrow e^{-1} \approx .368$.

(c) The probability that each of B bootstrap samples contains X_i is clearly $(1 - e_n)^B$. The following table gives values of this for $n = 10, 20, 50, 100$ and $B = 10, 20, 50, 100$.

B/n	10	20	50	100
10	.0137	.0118	.0108	.0105
20	.000189	.000139	.000117	.000110
50	4.89×10^{-10}	2.29×10^{-10}	1.47×10^{-10}	1.27×10^{-10}
100	2.39×10^{-19}	5.26×10^{-20}	2.16×10^{-20}	1.61×10^{-20}

(d) N_n/n is the proportion of the original sample not appearing in the bootstrap sample. Since $N_n \equiv \sum_{j=1}^n 1_{[M_j=0]}$ where each M_j is marginally $\text{Binomial}(n, 1/n)$, it follows immediately that

$$E(N_n/n) = P(M_1 = 0) = (1 - 1/n)^n \rightarrow e^{-1}.$$

Furthermore, from occupancy theory for urn models,

$$\sqrt{n}(n^{-1}N_n - (1 - 1/n)^n) \rightarrow_d N(0, e^{-1}(1 - 2e^{-1}));$$

see e.g. Johnson and Kotz (1977), page 317.

2. Suppose that $T(F) = \text{Var}_F(X)$ so that $T_n \equiv T(\mathbb{F}_n) = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$. Show that the jackknife estimate of the variance $\sigma_n^2(F) \equiv \text{Var}_F(T_n)$ is

$$\widehat{\text{Var}} = \frac{n^2}{(n-1)^3} (\widehat{\mu}_4 - \widehat{\mu}_2^2)$$

where $\widehat{\mu}_k \equiv n^{-1} \sum_{i=1}^n (X_i - \bar{X})^k$ for $k = 1, 2, \dots$. Hence, assuming that $EX^4 < \infty$, the jackknife estimate of variance is consistent for this T :

$$n\widehat{\text{Var}} \rightarrow_p \mu_4 - \mu_2^2 = \mu_2^2 \left\{ 2 + \frac{\mu_4}{\mu_2^2} - 3 \right\} = T_2(F)(2 + \gamma_2).$$

Solution: If $T_n = n^{-1} \sum_{j=1}^n (X_j - \bar{X})^2$, then some algebra yields

$$T_{n,i}^* = nT_n - (n-1)T_{n,i} = \frac{n}{n-1} (X_i - \bar{X})^2$$

and hence

$$\bar{T}_n^* = \frac{n}{n-1} \widehat{\mu}_2.$$

Furthermore,

$$\begin{aligned} \widehat{\text{Var}}_n &= \frac{1}{n(n-1)} \sum_{i=1}^n (T_{n,i}^* - \bar{T}_n^*)^2 \\ &= \frac{1}{n(n-1)} \sum_{i=1}^n T_{n,i}^{*2} - \frac{1}{n-1} (\bar{T}_n^*)^2 \\ &= \frac{1}{n(n-1)} \sum_{i=1}^n \left(\frac{n}{n-1} (X_i - \bar{X})^2 \right)^2 - \frac{1}{n-1} \left(\frac{n}{n-1} \widehat{\mu}_2 \right)^2 \\ &= \frac{1}{n-1} \frac{n^2}{(n-1)^2} \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^4 - \frac{1}{n-1} \frac{n^2}{(n-1)^2} \widehat{\mu}_2^2 \\ &= \frac{n^2}{(n-1)^3} (\widehat{\mu}_4 - \widehat{\mu}_2^2). \end{aligned}$$

Thus we have

$$n\widehat{\text{Var}}_n = \frac{n^3}{(n-1)^3} (\widehat{\mu}_4 - \widehat{\mu}_2^2) \rightarrow_p \mu_4 - \mu_2^2 = \mu_2^2 \left\{ 2 + \frac{\mu_4}{\mu_2^2} - 3 \right\} = T_2(F)(2 + \gamma_2).$$

where the (excess of) kurtosis is

$$\gamma_2 \equiv \frac{\mu_4}{\mu_2^2} - 3.$$

3. (a) Wasserman, problem 3.8.9, page 40: Let X_1, \dots, X_n be n distinct observations (no ties). Let X_1^*, \dots, X_n^* denote a bootstrap sample (from the empirical d.f. \mathbb{F}_n of the X_i 's), and let $\bar{X}_n^* = n^{-1} \sum_{i=1}^n X_i^*$. Find: $E\{\bar{X}_n^* | X_1, \dots, X_n\}$, $Var(\bar{X}_n^* | X_1, \dots, X_n)$, and $Var(\bar{X}_n^*)$.
- (b) Wasserman, problem 3.8.13, page 41.

Solution: (a) First, as we computed in class,

$$E(\bar{X}_n^* | X_1, \dots, X_n) = n^{-1} \sum_{i=1}^n E(X_i^* | X_1, \dots, X_n) = n^{-1} \sum_{i=1}^n \bar{X}_n = \bar{X}_n.$$

Similarly,

$$\begin{aligned} Var(\bar{X}_n^* | X_1, \dots, X_n) &= n^{-2} \sum_{i=1}^n Var(X_i^* | X_1, \dots, X_n) = n^{-2} \sum_{i=1}^n n^{-1} \sum_{j=1}^n (X_i - \bar{X}_n)^2 \\ &= n^{-1} S_n^2 \end{aligned}$$

where $S_n^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. Then it follows that

$$E(\bar{X}_n^*) = E\{E(\bar{X}_n^* | X_1, \dots, X_n)\} = E\{\bar{X}_n\} = \mu(F),$$

and

$$\begin{aligned} Var(\bar{X}_n^*) &= E\{Var(\bar{X}_n^* | X_1, \dots, X_n)\} + Var(E(\bar{X}_n^* | X_1, \dots, X_n)) \\ &= E\{n^{-1} S_n^2\} + Var(\bar{X}_n) \\ &= n^{-1} \frac{n-1}{n} \sigma^2(F) + n^{-1} \sigma^2(F) \\ &= n^{-1} \sigma^2(F) \left\{ \frac{n-1}{n} + 1 \right\} \\ &= n^{-1} \sigma^2(F) \frac{2n-1}{n} = \frac{2n-1}{n} Var(\bar{X}_n) \approx 2Var(\bar{X}_n). \end{aligned}$$

Another way to organize this calculation is as follows:

$$\begin{aligned} Var(\bar{X}_n^*) &= Var(\bar{X}_n^* - \bar{X}_n + \bar{X}_n - \mu + \mu) \\ &= Var(\bar{X}_n^* - \bar{X}_n) + Var(\bar{X}_n - \mu) + 2Cov(\bar{X}_n^* - \bar{X}_n, \bar{X}_n - \mu) \\ &= E\{Var(\bar{X}_n^* - \bar{X}_n | X_1, \dots, X_n)\} + Var\{E(\bar{X}_n^* - \bar{X}_n | X_1, \dots, X_n)\} \\ &\quad + n^{-1} \sigma^2(F) \\ &\quad + 2E\{(\bar{X}_n^* - \bar{X}_n)(\bar{X}_n - \mu)\} \\ &= E\{n^{-1} S_n^2\} + 0 + n^{-1} \sigma^2(F) + 0 \\ &= \frac{n-1}{n} n^{-1} \sigma^2(F) + n^{-1} \sigma^2(F), \end{aligned}$$

so that the contributions of $Var(\overline{X}_n^* - \overline{X}_n)$ and $Var(\overline{X}_n - \mu)$ are approximately equal. This is important, especially since we want the first of these to estimate the second! [The marginal behavior of bootstrap estimators is largely irrelevant, since this accounts for both the deviations $T(\mathbb{F}_n^*) - T(\mathbb{F}_n)$ and $T(\mathbb{F}_n) - T(F)$ via $T(\mathbb{F}_n) - T(F) = T(\mathbb{F}_n^*) - T(\mathbb{F}_n) + T(\mathbb{F}_n) - T(F)$. What is important is that $T(\mathbb{F}_n^*) - T(\mathbb{F}_n)$ mimics (or estimates) $T(\mathbb{F}_n) - T(F)$!]

(b) The marginal distribution of the X_i^* 's (separately, or marginally) agrees with the marginal distribution of the (separate) X_i 's, but the *joint distribution* of the X_i^* 's is dependent. For example,

$$\begin{aligned}
G_2(x_1, x_2) &\equiv P(X_1^* \leq x_1, X_2^* \leq x_2) = E\{P(X_1^* \leq x_1, X_2^* \leq x_2 | X_1, \dots, X_n)\} \\
&= E\{\mathbb{F}_n(x_1)\mathbb{F}_n(x_2)\} = E\{n^{-2} \sum_{i=1}^n \sum_{j=1}^n 1\{X_i \leq x_1\}1\{X_j \leq x_2\}\} \\
&= E\{n^{-2}[\sum_{i=1}^n 1\{X_i \leq x_1\}1\{X_i \leq x_2\} + \sum_{i \neq j} 1\{X_i \leq x_1\}1\{X_i \leq x_2\}]\} \\
&= n^{-2}\{nF(x_1 \wedge x_2) + n(n-1)F(x_1)F(x_2)\} \\
&= (1 - n^{-1})F(x_1)F(x_2) + n^{-1}F(x_1 \wedge x_2).
\end{aligned}$$

Thus $G_2(x_1, x_2)$ is a mixture of the independence distribution $F(x_1)F(x_2)$ and the (Fréchet bound) distribution concentrated on the diagonal, $F(x_1 \wedge x_2)$, with mixing probabilities $(1 - n^{-1})$ and n^{-1} . More generally,

$$\begin{aligned}
G(x_1, \dots, x_n) &\equiv P(X_1^* \leq x_1, \dots, X_n^* \leq x_n) \\
&= E\{P(X_1^* \leq x_1, \dots, X_n^* \leq x_n | X_1, \dots, X_n)\} \\
&= E\{\mathbb{F}_n(x_1) \cdots \mathbb{F}_n(x_n)\} = E\left\{\prod_{k=1}^n \mathbb{F}_n(x_k)\right\} \\
&= E\left\{n^{-n} \sum_{j_1=1}^n \cdots \sum_{j_n=1}^n \prod_{k=1}^n 1\{X_{j_k} \leq x_k\}\right\} \\
&= n^{-n} \sum_{j_1=1}^n \cdots \sum_{j_n=1}^n E\left\{\prod_{k=1}^n 1\{X_{j_k} \leq x_k\}\right\}
\end{aligned}$$

where now the computation becomes somewhat complicated. In particular, though X_1^* and X_2^* are correlated: from our formula (1.4.14) on page 19 of Chapter 1

(with $G(x) = H(x) = x$, $F = G_2$, and $F_X = F_Y = F$), it follows that

$$\begin{aligned}
 \text{Cov}(X_1^*, X_2^*) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{G_2(x, y) - F(x)F(y)\} dx dy \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{(1 - n^{-1})F(x)F(y) + n^{-1}F(x \wedge y) - F(x)F(y)\} dx dy \\
 &= n^{-1} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{F(x \wedge y) - F(x)F(y)\} dx dy \\
 &= n^{-1} \sigma^2(F)
 \end{aligned}$$

where the last equality follows from (1.14.16)-(1.14.17) on page 19, chapter 1. Thus

$$\text{Var}(\bar{X}_n^*) = n^{-2} \{n\sigma^2(F) + n(n-1)(n^{-1}\sigma^2(F))\} = n^{-1}\sigma^2(F) \left\{1 + \frac{n-1}{n}\right\}$$

in agreement with our calculation in (a).

4. Wasserman, example 3.10, page 29.

(a) In the 3rd line of this example, Wasserman says that the jackknife estimate of the standard error of the skewness estimator for the nerve data is .17. Check this. (Show your code or method of calculation.)

(b) I claimed in the solution to problem set #3, problem 3 (see page 9 of the solution set), that the estimated standard error using the delta method is .163 rather than .18 as Wasserman claimed. Check this. (Show your code or method of calculation.)

(c) On page 31, Wasserman claims that the bootstrap based on $B = 10^3$ replications gives .16 as a bootstrap estimate of the standard error. Try it yourself to see what you get. (Show your code or method of calculation.)

The data is posted in two forms at:

<http://www.stat.washington.edu/jaw/COURSES/580s/583/sp15.probsets.html/nrvdat>
and

<http://www.stat.washington.edu/jaw/COURSES/580s/583/sp15.probsets.html/nerve.dat>.

(d) Plot the empirical distribution function of this data together with the MLE of the distribution function assuming: (i) an exponential distribution; (ii) a Weibull distribution. What alternative methods for (a) and (b) does this suggest?

Solution: (a) When I calculate the jackknife estimator of the standard error of the skewness estimator for the nerve data, I get .172; see the posted Mathematica notebook *NerveDataAnal-2.nb*.

(b) As mentioned before, I get .163 for the delta-method estimator of the standard error; see the posted Mathematica notebook *NerveDataAnal-1.nb* which also

computes the delta-method estimator as in Wasserman (without the additional term in the influence function) to get .179.

(c) With a bootstrap sample size of 10^4 I get .1626 as the bootstrap estimator of standard error of the skewness estimator, almost exactly the same as the (corrected) delta method estimator; see the posted Mathematica notebook *NerveDataAnal-3.nb*.

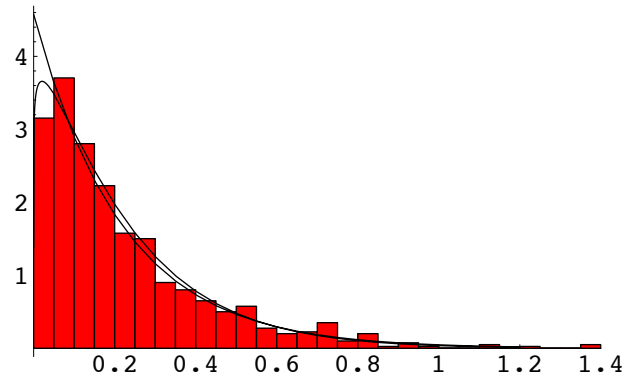


Figure 1: Histogram of the nerve data with exponential and Weibull fitted densities (via MLE).

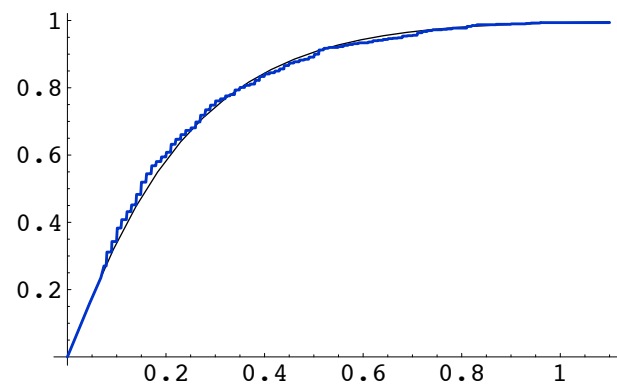


Figure 2: Empirical distribution function (blue) and fitted Weibull c.d.f. (via MLE).

(d) From the plot of the empirical distribution function of the nerve data on page 14 of Wasserman, one might guess that this data could be modeled by an

exponential distribution. The skewness of the exponential distribution is 2, and it should be noted that all of the 95% confidence intervals for skewness given by Wasserman on page 34 include 2. On the other hand, a histogram of the nerve data (Figure 1) shows a slight dip in the frequency relative to exponential in the cell just to the right of zero, and it turns out that a Weibull distribution with (α, β) estimated by $(\hat{\alpha}, \hat{\beta}) = (.22557, 1.08181)$ via maximum likelihood gives an excellent fit. Under the Weibull model the estimated skewness is 1.7778... (to be compared to the sample skewness of 1.761.... See the posted Mathematica notebook *NerveDataAnal4ML.nb*. Thus parametric inference and parametric bootstrap methods based on the Weibull model should perform quite well here.

5. **Optional bonus problem:** Show that for linear statistics the jackknife and bootstrap estimates of bias are zero. (A linear statistic corresponds to a functional $T(F)$ of the form $T(F) = \int \psi dF$ for some function ψ for a distribution function F on R , or $T(P) = \int \psi dP$ for a probability distribution P on a general sample space \mathcal{X} .)

Solution: Suppose that $T(F) = \int \psi dF$. Then

$$T(\mathbb{F}_n) = \int \psi d\mathbb{F}_n = n^{-1} \sum_{i=1}^n \psi(X_i) \equiv n^{-1} \sum_{i=1}^n Y_i,$$

so it suffices to prove the claim for the mean functional (of the Y 's).

(i) For the jackknife,

$$T_{n,i} = T(\mathbb{F}_{n-1,i}) = \frac{1}{n-1} \sum_{j \neq i} \psi(X_j) = \frac{1}{n-1} \sum_{j \neq i} Y_j = \frac{1}{n-1} (n\bar{Y}_n - Y_i).$$

Thus

$$\begin{aligned} T_{n,\cdot} &= \frac{1}{n} \sum_{i=1}^n T_{n,i} = \frac{n^2}{n(n-1)} \bar{Y}_n - \frac{1}{n-1} \bar{Y}_n \\ &= \frac{1}{n-1} (n-1) \bar{Y}_n = \bar{Y}_n. \end{aligned}$$

Furthermore, $T_{n,i}^* = nT_n - (n-1)T_{n,i} = Y_i$, and hence $\bar{T}_n^* = \bar{Y}_n$. Thus the jackknife estimator of bias is

$$\hat{\text{bias}}_n = T_n - \bar{T}_n^* = \bar{Y}_n - \bar{Y}_n = 0.$$

(ii) For the bootstrap, the (normalized) bias is

$$b_n(F) \equiv n\{E_F(T_n) - T(F)\}.$$

Hence the bootstrap estimator of (normalized) bias is given by

$$\begin{aligned} b_n(\mathbb{F}_n) &= n\{E_{\mathbb{F}_n}(T(\mathbb{F}_n^*)) - T(\mathbb{F}_n)\} = n\{E_{\mathbb{F}_n}(\bar{Y}_{n^*}) - \bar{Y}_n\} \\ &= n\{\bar{Y}_n - \bar{Y}_n\} = 0. \end{aligned}$$

6. **Optional bonus problem:** Consider a V -functional of order $r = 2$ given by $T(P) = \iint h(x, y)dP(x)dP(y)$ where h is permutation symmetric. Find the jackknife estimate of bias for the (V -statistic) estimator $T(\mathbb{P}_n)$ of $T(P)$.

Solution: Now we have $T(\mathbb{P}_n) = n^{-2} \sum_{i=1}^n \sum_{j=1}^n h(X_i, X_j) \equiv T_n$. Then

$$\begin{aligned} T_{n,i} &= T(\mathbb{P}_{n-1,i}) = \frac{1}{(n-1)^2} \sum_{j \neq i} \sum_{j' \neq i} h(X_j, X_{j'}) \\ &= \frac{1}{(n-1)^2} \left\{ n^2 T_n - \sum_{j'=1}^n h(X_i, X_{j'}) - \sum_{j=1}^n h(X_j, X_i) + h(X_i, X_i) \right\}. \end{aligned}$$

Thus

$$\begin{aligned} T_{n,\cdot} &= n^{-1} \sum_{i=1}^n T_{n,i} \\ &= \frac{1}{(n-1)^2} \left\{ n^2 T_n - \frac{1}{n} \sum_{i=1}^n \sum_{j'=1}^n h(X_i, X_{j'}) - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n h(X_j, X_i) + \frac{1}{n} \sum_{i=1}^n h(X_i, X_i) \right\}. \end{aligned}$$

Thus

$$\begin{aligned} T_{n,i}^* &= nT_n - (n-1)T_{n,i} \\ &= nT_n - \frac{1}{n-1} \left\{ n^2 T_n - \sum_{j'=1}^n h(X_i, X_{j'}) - \sum_{j=1}^n h(X_j, X_i) + h(X_i, X_i) \right\}, \end{aligned}$$

so that

$$\begin{aligned} \bar{T}_n^* &= nT_n - (n-1)T_{n,\cdot} \\ &= nT_n - \frac{1}{n-1} \left\{ n^2 T_n - 2nT_n + n^{-1} \sum_{i=1}^n h(X_i, X_i) \right\} \\ &= \left(n - \frac{n^2}{n-1} + \frac{2n}{n-1} \right) T_n - \frac{1}{n-1} \sum_{i=1}^n h(X_i, X_i) \\ &= \frac{n}{n-1} T_n - \frac{1}{n(n-1)} \sum_{i=1}^n h(X_i, X_i). \end{aligned}$$

It follows that the jackknife estimate of bias is given by

$$\begin{aligned}\widehat{\text{bias}}_n &= T_n - T_n^* = \left(1 - \frac{n}{n-1}\right) T_n + \frac{1}{n(n-1)} \sum_1^n h(X_i, X_i) \\ &= \frac{1}{n-1} \left\{ \frac{1}{n} \sum_1^n h(X_i, X_i) - T_n \right\}.\end{aligned}$$

Note that $T_n = V_n$ and

$$\begin{aligned}T_n^* &= \frac{n}{n-1} T_n - \frac{1}{n(n-1)} \sum_1^n h(X_i, X_i) \\ &= \frac{1}{n(n-1)} \left\{ \sum_{i=1}^n \sum_{j=1}^n h(X_i, X_j) - \sum_{i=1}^n h(X_i, X_i) \right\} \\ &= \frac{1}{n(n-1)} \sum_{i \neq j} h(X_i, X_j) = U_n,\end{aligned}$$

so that the estimator we just found is exactly $V_n - U_n$.