

For this assignment you can do it in a group of maximum three people. Objectives of this assignment:

- To gain experience in applying the bootstrap methodology.
- Last chance for Thomas to torture you.

Altogether there are 4+1 questions.

1. The following data were used to illustrate the bootstrap by Bradley Efron, the inventor of the bootstrap. The data are LSAT scores (for entrance to law school) and GPA.

LSAT:	576	635	558	578	666	580	555	661	651	605	653	575	545	572	594
GPA:	3.39	3.30	2.81	3.03	3.44	3.07	3.00	3.43	3.36	3.13	3.12	2.74	2.76	2.88	3.96

Each data point is of the form  $X_i = (Y_i, Z_i)$ , where  $Y_i = \text{LSAT}_i$  and  $Z_i = \text{GPA}_i$ .

- (a) Find the estimate of the correlation coefficient  $\rho_{YZ}$ .
  - (b) Estimate the standard error of your estimate in (a) using both jackknife and bootstrap.
  - (c) Compute 95% bootstrap confidence intervals using “normal theory” and bootstrap  $t$ -interval approaches.
2. Let  $X_1, \dots, X_n \sim \text{Uniform}(0, \theta)$ . The MLE is

$$\hat{\theta} = X_{\max} = \max(X_1, \dots, X_n).$$

- (a) Find the distribution of  $\hat{\theta}$  (in terms of  $\theta$  and  $n$ ).
- (b) Derive the analytic expression for the variance of  $\hat{\theta}$ . Call it  $\text{Var}_{F_\theta}(\hat{\theta})$ .
- (c) Generate a data set of size  $n = 50$  and  $\theta = 3$ . Then generate  $B = 5000$  bootstrap samples using parametric bootstrap. Use the bootstrap samples to approximate  $\text{Var}_{F_\theta}(\hat{\theta})$ . Compare your answer to (b).
- (d) With the same data set: repeat (c) with nonparametric bootstrap.
- (e) With the same data set: plot the histograms of  $\hat{\theta}^*$  obtained from the parametric and nonparametric bootstraps.
- (f) Compare the true distribution of  $\hat{\theta}$  to those histograms obtained in (e).

Note: this is an example where the nonparametric bootstrap fails. Can you guess why?

3. Generate two regression data sets  $(x_i, y_i)$ 's with  $n = 512$ , one from the test function in Assignment 2 and the other from

$$f(x) = (4x - 2) + 2 \exp\{-16(4x - 2)^2\},$$

where  $x_i = (i - 1)/n$ . Set the noise variance as  $\sigma^2 = (\|f\|/5)^2$ .

- (a) Obtain the regression curve estimates for both test functions using the genetic algorithm you implemented for Assignment 2.
- (b) Construct 95% pointwise confidence bands for both curve estimates using the bootstrap. Use both “bootstrapping residuals” and “bootstrapping pairs” approaches. For the test function from Assignment 2, comment on the shape of the confidence bands near jump points.

- (c) Describe how you would obtain a confidence interval for the location of a jump point using the bootstrap. You do not need to do any programming for this part.
4. Suppose  $X_1, \dots, X_n$  is a random sample from an exponential distribution with parameter  $\lambda$ , so its density is  $f(x) = \lambda e^{-\lambda x}$  for  $x \geq 0$ .

- (a) Show that the MLE is  $\hat{\lambda}_n = n/(X_1 + \dots + X_n)$ .
- (b) First show that

$$\sqrt{n}(\hat{\lambda}_n - \lambda) \xrightarrow{D} N(0, \lambda^2).$$

Then, use the delta method or otherwise, show that

$$\sqrt{n}(\log \hat{\lambda}_n - \log \lambda) \xrightarrow{D} N(0, 1).$$

- (c) Show that an asymptotic confidence interval for  $\lambda$  is

$$\left( \hat{\lambda}_n e^{-z_{\alpha/2}/\sqrt{n}}, \quad \hat{\lambda}_n e^{z_{\alpha/2}/\sqrt{n}} \right),$$

with  $z_{\alpha}$  denoting the  $1 - \alpha$  quantile of the standard normal distribution.

- (d) Note that  $\lambda(X_1 + \dots + X_n)$  has a Gamma distribution with parameters  $n$  and 1. Deduce from this that an exact confidence interval for  $\lambda$  is given by

$$\left( \hat{\lambda}_n G^{-1}(\alpha/2)/n, \quad \hat{\lambda}_n G^{-1}(1 - \alpha/2)/n \right),$$

where  $G$  denotes the CDF of a Gamma distribution with parameters  $n$  and 1.

- (e) Do a simulation study to compare the coverages and lengths of the exact and asymptotic confidence intervals as derived above, and the  $BC_a$  confidence interval.
5. (Note: This question is optional, and will not be graded.) In this exercise we will follow the necessary steps for deriving the bias corrected and accelerated ( $BC_a$ ) percentile confidence interval. Let  $X_1, \dots, X_n \sim F$ , where  $F$  is a distribution parametrized by  $\theta \in \mathbb{R}$ . Let  $\hat{\theta}_n = \theta(\hat{F}_n)$  be an estimator for  $\theta$ , where  $\hat{F}_n$  is the empirical distribution function. Now suppose there exists a monotonic transformation  $g(\cdot)$  and constants  $a$  and  $b$  such that  $\sigma_{g(\theta)} = 1 + ag(\theta)$  and

$$P_F \left( \frac{g(\hat{\theta}_n) - g(\theta)}{\sigma_{g(\theta)}} + b \leq x \right) = G(x),$$

where  $G(x)$  is a known continuous distribution function.

- (a) Assume that  $a$  and  $b$  are known. Show that a  $100\%(1 - 2\alpha)$  confidence interval for  $\theta$  is given by

$$[L_n, U_n] \stackrel{\text{def}}{=} \left[ g^{-1} \left( \frac{g(\hat{\theta}_n) - [G^{-1}(1 - \alpha) - b]}{1 + a[G^{-1}(1 - \alpha) - b]} \right), g^{-1} \left( \frac{g(\hat{\theta}_n) - [G^{-1}(\alpha) - b]}{1 + a[G^{-1}(\alpha) - b]} \right) \right].$$

- (b) Let  $\hat{\theta}_n^*$  be the bootstrap version of  $\hat{\theta}_n$  and  $\sigma_{g(\hat{\theta}_n)} = 1 + ag(\hat{\theta}_n)$ . The bootstrap argument asserts that  $\frac{g(\hat{\theta}_n) - g(\theta)}{\sigma_{g(\theta)}}$  and  $\frac{g(\hat{\theta}_n^*) - g(\hat{\theta}_n)}{\sigma_{g(\hat{\theta}_n)}} | \hat{F}_n$  have approximately the same distribution.

Denote now  $H(x) = P_{\hat{F}_n}(\hat{\theta}_n^* \leq x)$ , the distribution function of the bootstrap sample of  $\hat{\theta}_n^*$ . Show that

$$H(L_n) \approx G \left( b - \frac{G^{-1}(1 - \alpha) - b}{1 + a[G^{-1}(1 - \alpha) - b]} \right), \quad H(U_n) \approx G \left( b - \frac{G^{-1}(\alpha) - b}{1 + a[G^{-1}(\alpha) - b]} \right).$$

- (c) Now assume that  $G(x)$  is the standard normal distribution function. Show that  $L_n$  and  $U_n$  can be approximated by

$$\left[ H^{-1}(\alpha_1), H^{-1}(\alpha_2) \right] = \left[ \hat{\theta}_n^{*(\alpha_1)}, \hat{\theta}_n^{*(\alpha_2)} \right],$$

where  $\alpha_1 = \Phi \left( b + \frac{b+z^{(\alpha)}}{1-a(b+z^{(\alpha)})} \right)$  and  $\alpha_2 = \Phi \left( b + \frac{b+z^{(1-\alpha)}}{1-a(b+z^{(1-\alpha)})} \right)$ .

- (d) Give an estimate of  $b$  by showing that  $H(\hat{\theta}_n) \approx G(b)$ .
- (e) (This part is optional. That is, you can still get a perfect score for this assignment without completing this part. However, you will not get any hints or help from your instructor or TA.) If  $a$  and  $b$  are small, it can be shown that (by Taylor's expansion)

$$\alpha_1 \approx \alpha + [2b + a(z^{(\alpha)})^2] \phi(z^{(\alpha)}), \quad \alpha_2 \approx 1 - \alpha + [2b + a(z^{(1-\alpha)})^2] \phi(z^{(1-\alpha)}),$$

where  $\phi$  is the standard normal density. On the other hand, if  $g$  is asymptotically linear, the best  $\alpha_1$  and  $\alpha_2$  (due to Edgeworth and Cornell-Fisher expansions) are

$$\alpha_1 = \alpha + \left[ \frac{1}{3}\gamma + \frac{1}{6}\gamma(z^{(\alpha)})^2 \right] \phi(z^{(\alpha)}), \quad \alpha_2 = 1 - \alpha + \left[ \frac{1}{3}\gamma + \frac{1}{6}\gamma(z^{(1-\alpha)})^2 \right] \phi(z^{(1-\alpha)}),$$

where  $\gamma$  is the skewness of  $\hat{\theta}$ . By comparing the expressions of  $\alpha_1$  or  $\alpha_2$ , we obtain  $a = \gamma$ . Therefore,  $a$  can be estimated by estimating the skewness of  $\hat{\theta}$ .

Note: The actual argument of approximating  $a$  is technical and unintuitive. You can treat it as a way to better approximate  $\hat{\theta}_n$  when the distribution of  $\hat{\theta}_n$  is not symmetric. This argument also reveals that the optimal  $a$  and  $b$  are the same.

—— End of Assignment 6 ——