

1. See class notes.
2. See Chapter 7 and 8 notes.
3. A. Let $\mathbb{H}_{n,i}$ denote the empirical distribution of the data with the i th pair (X_i, Y_i) omitted. Let $T_{n,i} \equiv T(\mathbb{H}_{n,i})$, $T_n \equiv n^{-1} \sum_{i=1}^n T_{n,i}$ and let $T_{n,i}^* \equiv nT_n - (n-1)T_{n,i}$. Then the Jackknife estimator of $nVar_H(T_n)$ is

$$\frac{1}{n-1} \sum_{i=1}^n \{(T_{n,i}^* - \bar{T}_n^*)^2\}.$$

B. The ideal bootstrap estimator of $nVar_H(T_n)$ is $nVar_{\mathbb{H}_n}(T_n)$. To implement this, we would draw B bootstrap samples

$$(X_{j1}^*, Y_{j1}^*), \dots, (X_{jn}^*, Y_{jn}^*), \quad j = 1, \dots, B,$$

let $\mathbb{H}_{j,n}^*(x, y) \equiv n^{-1} \sum_{i=1}^n 1_{[X_{ji}^* \leq x, Y_{ji}^* \leq y]}$ be the empirical distribution function of the j th bootstrap sample, and compute $T_{j,n}^* \equiv T(\mathbb{H}_{j,n}^*)$, $j = 1, \dots, B$. Then the bootstrap estimator of $nVar_H(T_n)$ is just

$$n \frac{1}{B} \sum_{j=1}^B \{T_{j,n}^* - \bar{T}_n^*\}^2.$$

C. Because $T(H)$ is a smooth functional of the marginal first and second moments, and the natural substitution estimators of these parameters are jointly asymptotically normal under the hypotheses $EX^4 < \infty$, $EY^4 < \infty$, it is clear that $\sqrt{n}(T(\mathbb{H}_n) - T(H)) \rightarrow_d N(0, V^2)$ for some V^2 by the delta method. Because $T(H)$ is a ratio of moments, and because the central limit theorem does a better job of approximating sums rather than moments, it seems likely that use of the logarithmic transformation $g(x) = \log x$ might be helpful: $\sqrt{n}(\log(T(\mathbb{H}_n)) - \log(T(H)))$ will probably converge to normality faster than $\sqrt{n}(T(\mathbb{H}_n) - T(H))$.

D. Because $T(H)$ is a smooth functional of moments, and because the bootstrap works for moments, it follows by preservation of bootstrap convergence under differentiable mappings that the bootstrap will "work" in this situation. Similarly, the jackknife will also estimate $nVar_H(T_n)$ consistently in this problem.

4. (This problem is from the midterm exam!)

A. Show that this testing problem is invariant with respect to the group of scale changes, G given by $g_c(\underline{x}, \underline{y}) = (c\underline{x}, c\underline{y})$ where $c > 0$.

Solution: If $X \sim \text{exponential}(\lambda)$, then

$$\begin{aligned} P_\lambda(cX > t) &= P_\lambda(X > t/c) = \exp(-\lambda t/c) \\ &= \exp(-(\lambda/c)t) = P_{\lambda/c}(X > t), \end{aligned}$$

and similarly for $Y \sim \text{exponential}(\mu)$. Hence the induced group on the parameter space is $\bar{g}(\lambda, \mu) = (\lambda/c, \mu/c)$. Note that for any $\bar{g} \in \bar{G}$ we have $\bar{g}\Theta_0 = \{(\lambda/c, \mu/c) : \lambda \leq \mu\} = \{(\lambda, \mu) : \lambda \leq \mu\} = \Theta_0$ and $\bar{g}\Theta = \{(\lambda/c, \mu/c) : (\lambda, \mu) \in R^+ \times R^+\} = \Theta$. Hence the testing problem is invariant under the group G .

B. Find the UMP G -invariant test of H versus K . [Hint: You may use the fact that the family of distributions $\{\delta^{-1}F_{r,s} : \delta > 0\}$ has monotone likelihood ratio.] [I should have said *monotone decreasing likelihood ratio*.]

Solution: By sufficiency we may reduce to consideration of $(S, T) \equiv (\sum_1^m X_i, \sum_1^n Y_j)$. The induced group G^* on the space of the sufficient statistic is given by $G^* = \{g^*(s, t) = (cs, ct) : c > 0\}$, and the maximal invariant for the group G^* is $V \equiv S/T$; the corresponding \bar{G} maximal invariant is $\delta = \lambda/\mu$. Now $2\lambda X_i \sim \chi_2^2$, and similarly $2\mu Y_j \sim \chi_2^2$. hence $2\lambda S \sim \chi_{2m}^2$ and $2\mu T \sim \chi_{2n}^2$. Hence

$$\frac{n}{m}V = \frac{\mu}{\lambda} \cdot \frac{(2\lambda S/2m)}{(2\mu T/2n)} = \delta^{-1}F_{2m,2n}$$

where $F_{2m,2n}$ has an F -distribution with degrees of freedom $2m, 2n$. Since the family $\delta^{-1}F_{r,s}$ has monotone decreasing monotone likelihood ratio, we conclude that the UMP G -invariant test of H versus K is given by "reject H if $nV/m < F_{2m,2n,\alpha}$ " where $P(F_{2m,2n} \leq F_{2m,2n,\alpha}) = \alpha$. (Alternatively, "reject H if $m/(nV) = (n^{-1}T/m^{-1}S) > F_{2m,2n,1-\alpha}$ " where $P(F_{2m,2n} \geq F_{2m,2n,1-\alpha}) = \alpha$.)

C. Specify as exactly as possible how you would carry out the test derived in B.

Solution: See part B above.

5. In the context of testing for a disease, let $X \sim F$ denote the outcome of the test for a diseased individual and let $Y \sim G$ denote the outcome of the test for a non-diseased individual. Assuming that $X > x$ (or $Y > x$) leads to classifying the individual as "diseased", the *Receiver Operating Characteristic* or ROC curve R is a plot of *sensitivity* $\equiv P(X > x) = \bar{F}(x)$ versus $1 - \textit{specificity}$ $\equiv 1 - P(Y \leq x) = P(Y > x) = \bar{G}(x)$. Thus the ROC curve $R = R_{F,G}$ can be written as

$$R(t) = \bar{F}(\bar{G}^{-1}(t)) = 1 - F(G^{-1}(1-t)), \quad 0 < t < 1.$$

A. A good test for a disease has ROC curve with values close to 1 for small t and is everywhere above the line $I(t) = t$. Show that $R(t) \geq t$ for $0 \leq t \leq 1$

with strict inequality for some t if and only if $G <_s F$.

Solution: $R(t) = 1 - F(G^{-1}(1-t)) \geq t$ for all $0 \leq t \leq 1$ if and only if $1-t \geq F(G^{-1}(1-t))$ for all $0 \leq t \leq 1$, if and only if $u \geq F(G^{-1}(u))$ for all $0 \leq u \leq 1$, if and only if $F^{-1}(u) \geq G^{-1}(u)$ for all $0 \leq u \leq 1$, if and only if $F(x) \leq G(x)$ for all $-\infty < x < \infty$, and the latter means $G \leq_s F$. If strict inequality holds for some t in the first inequality ($R(t) > t$), then strict inequality holds in $F(x) < G(x)$ for some x , and hence $G <_s F$.

B. Consider $A \equiv \int_0^1 R(t)dt$ as a measure of the quality of the disease test (values close to 1 indicating an excellent test). Show that A can be expressed in terms of the Mann-Whitney-Wilcoxon functional $\int F dG$.

Solution:

$$\begin{aligned} \int_0^1 R(t) dt &= 1 - \int_0^1 F \circ G^{-1}(1-t) dt \\ &= 1 - \int_0^1 F \circ G^{-1}(u) du \\ &= 1 - EF \circ G^{-1}(U) \quad \text{where } U \sim \text{Uniform}(0,1). \\ &= 1 - EF(Y) \quad \text{since } Y \equiv G^{-1}(U) \text{ has distribution } G \\ &= 1 - \int F dG = P(Y < X). \end{aligned}$$

C. Suppose that X_1, \dots, X_m are i.i.d. F and Y_1, \dots, Y_n are i.i.d. G and consider estimation of the ROC curve $R \equiv R_{F,G}$ on the basis of the data.

(i) Propose a nonparametric estimator $\mathbb{R}_{m,n}(t)$ of $R(t) = R_{F,G}(t)$.

Solution: A natural estimator is $\mathbb{R}_{m,n}(t) \equiv R_{IF_m, IG_n}(t) = 1 - IF_m(IG_n^{-1}(1-t))$.

(ii) Give conditions on F and G which imply that your estimator in (i) is consistent for a fixed $t \in (0,1)$.

Solution: Now $IF_m(x) \rightarrow_{a.s.} F(x)$ uniformly in x by the Glivenko-Cantelli theorem, and similarly, $IG_n(y) \rightarrow_{a.s.} G(y)$ uniformly in y . Furthermore, $IG_n^{-1}(t) \rightarrow_{a.s.} G^{-1}(t)$ for any t for which G^{-1} is continuous at t . It follows that $\mathbb{R}_{m,n}(t) \rightarrow_{a.s.} R(t)$ for any t which is a point of continuity of G^{-1} .

(iii) Compute the Gateaux derivatives of the functionals $T_1(F) \equiv R_{F,G}(t)$ for fixed G and t and $T_2(G) \equiv R_{F,G}(t)$ for fixed F and t .

Solution: Let F_1 be a fixed distribution function $\neq F$, and set $F_\epsilon \equiv (1-\epsilon)F + \epsilon F_1$. Then

$$\begin{aligned} \frac{d}{d\epsilon} T_1(F_\epsilon)|_{\epsilon=0} &= \frac{d}{d\epsilon} \{1 - F_\epsilon \circ G^{-1}(1-t)\}|_{\epsilon=0} \\ &= -(F_1 - F)(G^{-1}(1-t)) \equiv \dot{T}_1(F; F_1 - F). \end{aligned}$$

Similarly, Let G_1 be a fixed distribution function $\neq G$, and set $G_\epsilon \equiv (1-\epsilon)G + \epsilon G_1$. Then

$$\begin{aligned} \frac{d}{d\epsilon} T_2(G_\epsilon)|_{\epsilon=0} &= \frac{d}{d\epsilon} \{1 - F \circ G_\epsilon^{-1}(1-t)\}|_{\epsilon=0} \\ &= -f(G^{-1}(1-t)) \frac{d}{d\epsilon} G_\epsilon^{-1}(1-t)|_{\epsilon=0} \\ &= \frac{f(G^{-1}(1-t))}{g(G^{-1}(1-t))} \{G_1(G^{-1}(1-t)) - (1-t)\} \end{aligned}$$

because

$$\begin{aligned} 0 &= \frac{d}{d\epsilon} (1-t)|_{\epsilon=0} = \frac{d}{d\epsilon} G_\epsilon(G_\epsilon^{-1}(1-t))|_{\epsilon=0} \\ &= \frac{d}{d\epsilon} \{G(G_\epsilon^{-1}(1-t)) + \epsilon(G_1 - G) \circ G_\epsilon^{-1}(1-t)\}|_{\epsilon=0} \\ &= g(G^{-1}(1-t)) \frac{d}{d\epsilon} G_\epsilon^{-1}(1-t)|_{\epsilon=0} \\ &\quad + (G_1 - G) \circ G^{-1}(1-t). \end{aligned}$$

(iv) Give conditions on F and G which imply that your estimator in (ii) is asymptotically normal (for a fixed $t \in (0,1)$). Find the influence function of your estimator (with help from (iii)).

Solution: Suppose that F and G are differentiable at $G^{-1}(1-t)$ with derivatives f and g respectively and suppose that $g(G^{-1}(1-t)) > 0$. Then, for the case $m = n$

$$\begin{aligned} \sqrt{n}(\mathbb{R}_{n,n}(t) - R(t)) &= -\sqrt{n}(IF_n(IG_n^{-1}(1-t)) - F(G^{-1}(1-t))) \\ &= -\sqrt{n}(IF_n(IG_n^{-1}(1-t)) - F(IG_n^{-1}(1-t))) \\ &\quad - \sqrt{n}(F(IG_n^{-1}(1-t)) - F(G^{-1}(1-t))) \\ &= -\sqrt{n}(IF_n(G^{-1}(1-t)) - F(G^{-1}(1-t))) + o_p(1) \end{aligned}$$

$$\begin{aligned}
 & - \frac{(F(IG_n^{-1}(1-t)) - F(G^{-1}(1-t)))}{IG_n^{-1}(1-t) - G^{-1}(1-t)} \sqrt{n}(IG_n^{-1}(1-t) - G^{-1}(1-t)) \\
 &= - \sqrt{n}(IF_n(G^{-1}(1-t)) - F(G^{-1}(1-t))) + o_p(1) \\
 & \quad - \frac{f}{g}(G^{-1}(1-t))\sqrt{n}(IG_n^{-1}(1-t) - G^{-1}(1-t)) + o_p(1) \\
 &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \{ - (1_{(-\infty, G^{-1}(1-t))}(X_i) - F(G^{-1}(1-t))) \\
 & \quad + \frac{f}{g}(G^{-1}(1-t))(1_{(-\infty, G^{-1}(1-t))}(Y_i) - (1-t)) \} + o_p(1),
 \end{aligned}$$

so the influence function in this case is

$$\begin{aligned}
 \psi_{F,G}(x, y) &= - (1_{(-\infty, G^{-1}(1-t))}(x) - F(G^{-1}(1-t))) \\
 & \quad + \frac{f}{g}(G^{-1}(1-t))(1_{(-\infty, G^{-1}(1-t))}(y) - (1-t))
 \end{aligned}$$

When $m \neq n$ it is natural to normalize by $\sqrt{mn/N}$ (as we did in the case of the Mann-Whitney-Wilcoxon statistic), and then, assuming that $\lambda_N \equiv m/N \rightarrow \lambda$, the influence function for $\mathbb{R}_{m,n}(t)$ is given by the pair of functions

$$\left(\begin{array}{l} -\sqrt{1-\lambda}(1_{(-\infty, G^{-1}(1-t))}(x) - F(G^{-1}(1-t))) \\ \sqrt{\lambda} \frac{f}{g}(G^{-1}(1-t))(1_{(-\infty, G^{-1}(1-t))}(y) - (1-t)) \end{array} \right).$$

(v) What can you say about your estimator \mathbb{R}_n as an estimator of the function R uniformly in $0 \leq t \leq 1$?

Solution: First, note that

$$\begin{aligned}
 \sup_{0 \leq t \leq 1} |\mathbb{R}_{n,n}(t) - R(t)| &= \sup_{0 \leq t \leq 1} |IF_n(IG_n^{-1}(1-t)) - F(G^{-1}(1-t))| \\
 &\leq \sup_{0 \leq t \leq 1} |IF_m(IG_n^{-1}(1-t)) - F(IG_n^{-1}(1-t))| \\
 & \quad + \sup_{0 \leq t \leq 1} |F(IG_n^{-1}(1-t)) - F(G^{-1}(1-t))| \\
 &\leq \sup_{-\infty < x < \infty} |IF_m(x) - F(x)|
 \end{aligned}$$

$$+ \sup_{0 \leq t \leq 1} |F(G^{-1}(\Gamma_n^{-1}(1-t))) - F(G^{-1}(1-t))|$$

where, in the last line, we have replaced IG_n^{-1} by something equal in distribution (jointly in n), namely $G^{-1}(\Gamma_n^{-1})$ with $\Gamma_n(t) = n^{-1} \sum_{i=1}^n 1_{[0,t]}(\xi_i)$ the empirical d.f. of n i.i.d. Uniform(0,1) random variables. By the Glivenko-Cantelli theorem we have $\|IF_m - F\|_\infty \rightarrow_{a.s.} 0$, so the first term converges a.s. to 0. Also by the Glivenko-Cantelli theorem and symmetry about the identity, $\|\Gamma_n^{-1} - I\|_\infty \rightarrow_{a.s.} 0$. Thus we see that the second term will converge a.s. to 0 if the function $F(G^{-1}(1-t))$ is uniformly continuous on $[0,1]$; i.e. if $F \circ G^{-1}$ is continuous on the closed interval $[0,1]$.

(vi) Explain how and why (or why not) you could use the jackknife or bootstrap to estimate the variance of the estimator in (i). Be as detailed and explicit as possible.

Solution: Because the jackknife fails for the median and other quantiles, and because $R_{F,G}$ involves $G^{-1}(1-t)$, it seems likely that the jackknife will fail for $R_{F,G}$. On the other hand, the bootstrap "works" for quantiles and for linear statistics (and differentiable functionals more generally), and hence it will work for $R_{F,G}(t)$.