

Statistics 583, Problem Set 4 Solutions

Wellner; 4/29/2009

1. Let $U_{m,n} \equiv T(\mathbb{F}_m, \mathbb{G}_n)$ where $T(F, G) = \int F dG = P(X \leq Y)$ is the Mann-Whitney functional and \mathbb{F}_m and \mathbb{G}_n are the empirical df's of X_1, \dots, X_m i.i.d. with df F , Y_1, \dots, Y_n i.i.d. with df G where F and G are continuous.
- (a) Show that

$$mnU_{m,n} + n(n+1)/2 = W_{m,n} \equiv \sum_{j=1}^n Q_j = \sum_{j=1}^n R_{m+j}.$$

- (b) Show that $EU_{m,n} = P(X \leq Y) = \int F dG$ and that

$$\begin{aligned} \text{Var}(\sqrt{mn}U_{m,n}) &= (n-1) \int (1-G)^2 dF + (m-1) \int F^2 dG - (N-1) \left(\int F dG \right)^2 + \int F dG \\ &= (n-1) \text{Var}[1-G(X)] + (m-1) \text{Var}[F(Y)] + \int F dG \left(1 - \int F dG \right). \end{aligned}$$

- (c) When $F = G$ use the results of A and B to compute $E_{(F,F)}W_{m,n}$ and $\text{Var}_{(F,F)}(W_{m,n})$. (This should agree with calculations for the Wilcoxon rank sum form of the statistic under the null hypothesis via finite sampling calculations.)

Solution: (a) Using empirical distribution function notation, $N\mathbb{H}_N = m\mathbb{F}_m + n\mathbb{G}_n$, so

$$\begin{aligned} mnU_{m,n} &= \int m\mathbb{F}_m d(n\mathbb{G}_n) = \int N\mathbb{H}_N d(n\mathbb{G}_n) - \int n\mathbb{G}_n d(n\mathbb{G}_n) \\ &= \sum_{j=1}^n N\mathbb{H}_N(Y_j) - \sum_{j=1}^n n\mathbb{G}_n(Y_j) \\ &= \sum_{j=1}^n R_{m+j} - \sum_{j=1}^n j \\ &= \sum_{j=1}^n R_{m+j} - n(n+1)/2. \end{aligned}$$

- (b) The expectation is easy:

$$E(U_{m,n}) = \frac{1}{mn} \sum_{j=1}^m \sum_{i=1}^n P(X_i \leq Y_j) = P(X_1 \leq Y_1) = \int F dG.$$

For the variance, we first calculate

$$\begin{aligned}
E[mnU_{m,n}]^2 &= \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^m \sum_{l=1}^n E1_{[X_i \leq Y_j, X_k \leq Y_l]} \\
&= \sum_{i=1}^m \sum_{j=1}^n E1_{[X_i \leq Y_j]} + \sum_{i \neq k} \sum_{j=1}^n P(X_i \leq Y_j, X_k \leq Y_j) \\
&\quad + \sum_{i=1}^m \sum_{j \neq l} P(X_i \leq Y_j, X_i \leq Y_l) + \sum_{i \neq k} \sum_{j \neq l} P(X_i \leq Y_j, X_k \leq Y_l) \\
&= mnP(X_1 \leq Y_1) + m(m-1)nP(X_1 \leq Y_1, X_2 \leq Y_1) \\
&\quad + mn(n-1)P(X_1 \leq Y_1, X_1 \leq Y_2) \\
&\quad + m(m-1)n(n-1)P(X_1 \leq Y_1, X_2 \leq Y_2),
\end{aligned}$$

$$P(X_1 \leq Y_1) = \int FdG,$$

$$P(X_1 \leq Y_1, X_2 \leq Y_1) = EP(X_1 \leq Y_1, X_2 \leq Y_1 | Y_1) = \int F^2(x)dG(x),$$

$$P(X_1 \leq Y_1, X_1 \leq Y_2) = EP(X_1 \leq Y_1, X_1 \leq Y_2 | X_1) = \int (1 - G(x-))^2 dF(x),$$

and

$$P(X_1 \leq Y_1, X_2 \leq Y_2) = P(X_1 \leq Y_1)^2 = \left(\int FdG \right)^2.$$

It follows by algebra that

$$\begin{aligned}
\text{Var}(mnU_{m,n}) &= E(mnU_{m,n})^2 - \{E(mnU_{m,n})\}^2 \\
&= mn \int FdG + m(m-1)n \int F^2dG \\
&\quad + mn(n-1) \int (1 - G(x-))^2 dF(x) \\
&\quad + m(m-1)n(n-1) \{ \int FdG \}^2 - (mn \int FdG)^2 \\
&= m(m-1)n \{ \int F^2dG - (\int FdG)^2 \} \\
&\quad + mn(n-1) \{ \int (1 - G(x-))^2 dF(x) - (\int FdG)^2 \} \\
&\quad - mn \int FdG (1 - \int FdG).
\end{aligned}$$

By noting that

$$\begin{aligned}\int FdG = P(X \leq Y) &= 1 - P(X > Y) \\ &= 1 - \int G(x-)dF(x) = \int (1 - G(x-))dF(x),\end{aligned}$$

this yields the claimed variance formula (to within a left limit):

$$\begin{aligned}\text{Var}(\sqrt{mn}U_{m,n}) &= (m-1)\text{Var}(F(Y)) + (n-1)\text{Var}(1 - G(X-)) \\ &\quad + \int FdG \left(1 - \int FdG\right).\end{aligned}$$

(c) When $F = G$ continuous we find that

$$E(U_{mn}) = \int FdF = 1/2,$$

and, since now $\text{Var}[F(Y)] = \text{Var}[G(X)] = 1/12$,

$$\begin{aligned}\text{Var}(\sqrt{mn}U_{m,n}) &= (m-1)\frac{1}{12} + (n-1)\frac{1}{12} + \frac{1}{4} \\ &= (N-2)\frac{1}{12} + \frac{1}{4} = (N+1)\frac{1}{12}.\end{aligned}$$

Hence from part (a) it follows that

$$E\left(\sum_{j=1}^n Q_j\right) = n(n+1)/2 + mnE(U_{m,n}) = n(N+1)/2$$

and

$$\text{Var}\left(\sum_{j=1}^n Q_j\right) = mn\text{Var}(\sqrt{mn}U_{m,n}) = mn(N+1)\frac{1}{12}$$

both of which agree with finite sampling calculations (drawing a sample of size n balls from an urn without replacement where the numbers on the N balls in the urn are $\{1, 2, \dots, N\}$).

2. Suppose that \mathcal{F}_+ is the class of distribution functions F on \mathbb{R}^+ with mean $\mu_F = E_F X < \infty$, and consider the functional $T(F)$ defined for a fixed $x_0 \in R^+$ by

$$T(F) \equiv e_F(x_0) \equiv E_F(X - x_0 | X > x_0) = \frac{\int_{x_0}^{\infty} (1 - F(t))dt}{1 - F(x_0)}.$$

This functional is the *mean residual life* functional.

(a) For what collection of df's F_0 is T weakly continuous at F_0 ? For what collection of df's F_0 is T continuous at F_0 with respect to the Kolmogorov metric?

(b) Find the influence function of $T(F)$.

Solution: (a) Suppose that $F_n \rightarrow_d F$, and write

$$T(F) = \frac{\int_0^\infty (x - x_0) 1_{(x_0, \infty)}(x) dF(x)}{1 - F(x_0)}.$$

Thus

$$T(F_n) = \frac{\int_0^\infty (x - x_0) 1_{(x_0, \infty)}(x) dF_n(x)}{1 - F_n(x_0)} \quad (1)$$

$$\rightarrow \frac{\int_0^\infty (x - x_0) 1_{(x_0, \infty)}(x) dF(x)}{1 - F(x_0)} = T(F) \quad (2)$$

by the Helly-Bray lemma if x_0 is a continuity point of F and if x is uniformly integrable with respect to the sequence F_n :

$$\limsup_n E_{F_n} \{X 1_{[X \geq \lambda]}\} \rightarrow 0$$

as $\lambda \rightarrow \infty$. If $F_n \rightarrow F$ with respect to the Kolmogorov metric d_K , then $F_n \rightarrow F$, so the previous argument goes through under the same assumptions, but (2) may continue to hold even when F is discontinuous at x_0 . To see this, note that $d_K(F_n, F) \rightarrow 0$ implies that $F_n(x_0) \rightarrow F(x_0)$ (even if F is discontinuous at x_0), while the numerator of $T(F_n)$ can be written as

$$\begin{aligned} \int_{x_0}^\infty (1 - F_n(t)) dt &= \int_{x_0}^M (1 - F_n(t)) dt + \int_M^\infty (1 - F_n(t)) dt \\ &\rightarrow \int_{x_0}^M (1 - F(t)) dt + \epsilon = \int_{x_0}^\infty (1 - F(t)) dt + 2\epsilon \end{aligned}$$

by using $d_K(F_n, F) \rightarrow 0$ and uniform integrability.

(b) First note that with $F_t \equiv (1 - t)F + tG$ we have both

$$\frac{d}{dt}(1 - F_t(x_0)) \Big|_{t=0} = -(G - F)(x_0)$$

and

$$\frac{d}{dt} \int_{x_0}^\infty (1 - F_t(y)) dy \Big|_{t=0} = - \int_{x_0}^\infty (G - F)(y) dy.$$

Thus by the product rule we calculate

$$\begin{aligned}\frac{d}{dt}T(F_t)|_{t=0} &= -\frac{\int_{x_0}^{\infty}(G-F)(y)dy}{1-F(x_0)} + \frac{\int_{x_0}^{\infty}(1-F(y))dy}{(1-F(x_0))^2}(G-F)(x_0) \\ &= e_F(x_0)\frac{(G-F)(x_0)}{1-F(x_0)} - \frac{\int_{x_0}^{\infty}(G-F)(y)dy}{1-F(x_0)}.\end{aligned}$$

Taking $G = \delta_x = 1_{[x, \infty)}$ yields the influence function for T at F :

$$\begin{aligned}IC(x; T, F) &= e_F(x_0)\frac{(1_{[x, \infty)}(x_0) - F(x_0))}{1-F(x_0)} - \frac{\int_{x_0}^{\infty}(1_{[x, \infty)}(y) - F(y))dy}{1-F(x_0)} \\ &= e_F(x_0)\frac{(1_{[0, x_0]}(x) - F(x_0))}{1-F(x_0)} - \frac{\int_{x_0}^{\infty}(1_{[0, y]}(x) - F(y))dy}{1-F(x_0)} \\ &= \begin{cases} e_F(x_0) - \frac{\int_{x_0}^{\infty}1-F(y)dy}{1-F(x_0)} & x \leq x_0 \\ \frac{-F(x_0)}{1-F(x_0)}e_F(x_0) - \frac{\int_{x_0}^{\infty}1_{[0, y]}(x_0)-F(y)dy}{1-F(x_0)} & x > x_0 \end{cases} \\ &= \begin{cases} 0 & x \leq x_0 \\ \frac{-F(x_0)}{1-F(x_0)}e_F(x_0) - \frac{\int_{x_0}^x -F(y)dy}{1-F(x_0)} - \frac{\int_x^{\infty}1_{[0, y]}(x_0)-F(y)dy}{1-F(x_0)} & x > x_0 \end{cases} \\ &= \begin{cases} 0 & x \leq x_0 \\ \frac{(x-x_0)-e_F(x_0)}{1-F(x_0)} & x > x_0 \end{cases} \\ &= \frac{[(x-x_0) - e_F(x_0)]1_{(x_0, \infty)}(x)}{1-F(x_0)}.\end{aligned}$$

Note that

$$E_F[IC^2(X; T, F)] = \frac{Var(X - x_0 | X > x_0)}{1-F(x_0)}.$$

3. Let F be a distribution function on \mathbb{R}^2 with finite second moments, and let $\rho(F)$ be the correlation coefficient

$$\rho(F) = \frac{Cov_F(X, Y)}{\sqrt{Var_F(X)Var_F(Y)}}.$$

Assume that $|\rho(F)| < 1$.

- (a) Give an example of a sequence of bivariate distributions $\{F_n\}$ satisfying $F_n \rightarrow_d F$, but $\rho(F_n) \rightarrow 1 \neq \rho(F)$.
(b) Find a collection \mathcal{F} of distribution functions on \mathbb{R}^2 so that ρ is weakly continuous on \mathcal{F} .

Solution: (a) Without loss of generality we may suppose that F is a bivariate distribution function with zero means, $E_F(X) = E_F(Y) = 0$. Let $F_n = (1 - n^{-1})F + n^{-1}\delta_{(a_n, b_n)}$ with $(a_n, b_n) \in \mathbb{R}^2$. Note that F_n has marginal distribution functions $F_{n,X} = (1 - n^{-1})F_X + n^{-1}\delta_{a_n}$, $F_{n,Y} = (1 - n^{-1})F_Y + n^{-1}\delta_{b_n}$ respectively where F_X and F_Y are the marginal df's of F . Thus we compute

$$\begin{aligned}
Cov_{F_n}(X, Y) &= E_{F_n}(XY) - E_{F_n}(X)E_{F_n}Y \\
&= (1 - n^{-1})E(XY) + n^{-1}a_nb_n - ((1 - n^{-1})E_F X + n^{-1}a_n)((1 - n^{-1})E_F Y + n^{-1}b_n) \\
&= (1 - n^{-1})\{E_F(XY) - E_F X \cdot E_F Y\} \\
&\quad + (1 - n^{-1})E_F X \cdot E_F Y - (1 - n^{-1})^2 E_F X \cdot E_F Y \\
&\quad - (1 - n^{-1})E_F X \cdot n^{-1}b_n - (1 - n^{-1})E_F Y \cdot n^{-1}a_n - n^{-2}a_nb_n \\
&= (1 - n^{-1})Cov_F(X, Y) + n^{-1}(1 - n^{-1})a_nb_n \\
&\quad - n^{-1}(1 - n^{-1})\{E_F X \cdot b_n + E_F Y \cdot a_n\} + n^{-1}(1 - n^{-1})E_F X \cdot E_F Y \\
&= (1 - n^{-1})Cov_F(X, Y) + n^{-1}(1 - n^{-1})a_nb_n
\end{aligned}$$

since $E_F X = E_F Y = 0$. Similarly,

$$\begin{aligned}
Var_{F_n}(X) &= E_{F_n}(X^2) - (E_{F_n}(X))^2 \\
&= (1 - n^{-1})E_F X^2 + n^{-1}a_n^2 - ((1 - n^{-1})E_F X + n^{-1}a_n)^2 \\
&= (1 - n^{-1})\{E_F(X^2) - (E_F X)^2\} + n^{-1}(1 - n^{-1})(E_F X)^2 \\
&\quad + n^{-1}(1 - n^{-1})a_n^2 - 2n^{-1}(1 - n^{-1})E_F X \cdot a_n \\
&= (1 - n^{-1})Var_F(X) + n^{-1}(1 - n^{-1})a_n^2 + n^{-1}(1 - n^{-1})\{(E_F X)^2 - 2E_F X \cdot a_n\} \\
&= (1 - n^{-1})Var_F(X) + n^{-1}(1 - n^{-1})a_n^2, \quad \text{and} \\
Var_{F_n}(Y) &= E_{F_n}(Y^2) - (E_{F_n}(Y))^2 = (1 - n^{-1})Var_F(Y) + n^{-1}(1 - n^{-1})b_n^2.
\end{aligned}$$

Choosing $a_n = b_n = n$ yields

$$\begin{aligned}
Cov_{F_n}(X, Y) &= n + o(n) = n(1 + o(1)), \\
Var_{F_n}(X) &= n + o(n) = n(1 + o(1)), \\
Var_{F_n}(Y) &= n + o(n) = n(1 + o(1)).
\end{aligned}$$

Thus we find that

$$\rho(F_n) = \frac{Cov_{F_n}(X, Y)}{\sqrt{Var_{F_n}(X)Var_{F_n}(Y)}} = \frac{n(1 + o(1))}{n(1 + o(1))} \rightarrow 1$$

as $n \rightarrow \infty$. Thus ρ is weakly discontinuous at every F .

(b) Consider the following collection of distributions on \mathbb{R}^2 : for some $r > 2$ and $M < \infty$

$$\mathcal{F}_{r,M} \equiv \{F : E_F|X|^r \leq M, E_F|Y|^r \leq M\}.$$

Then ρ is weakly-continuous on $\mathcal{F}_{r,M}$ at any F with $Var_F(X) > 0$ and $Var_F(Y) > 0$. Here is a proof: let $\{F_n\} \subset \mathcal{F}_{r,M}$ satisfy $F_n \rightarrow_d F$. Then with $(X_n, Y_n) \sim F_n$ and $(X, Y) \sim F$ we have $(X_n, Y_n) \rightarrow_d (X, Y)$, and by a Skorokhod construction there exist $(X_n^*, Y_n^*) =_d (X_n, Y_n)$ and $(X^*, Y^*) =_d (X, Y)$ defined on a common probability space and satisfying $(X_n^*, Y_n^*) \rightarrow_{a.s.} (X^*, Y^*)$. But because $\{F_n \subset \mathcal{F}_{r,M}, X_n^2, Y_n^2, \text{ and } |X_n Y_n|\}$ are all uniformly integrable: since $r > 2$,

$$EX_n^2 1_{[X_n^2 \geq \lambda]} \leq \frac{1}{\lambda^{r-2}} E|X_n|^r \leq \frac{M}{\lambda^{r-2}}$$

so

$$\lim_{\lambda \rightarrow \infty} \limsup_{n \rightarrow \infty} EX_n^2 1_{[X_n^2 \geq \lambda]} \leq \lim_{\lambda \rightarrow \infty} \frac{M}{\lambda^{r-2}} = 0$$

and similarly for $\{Y_n^2\}$, so the uniform integrability of $|X_n Y_n|$ follows by Cauchy-Schwarz. The same holds true for the (X_n^*, Y_n^*) pairs since the uniform integrability only depends on the (marginal) distributions. Thus by Vitali's theorem it follows that

$$EX_n^s = EX_n^* s \rightarrow EX^* s = EX^s$$

and

$$EY_n^s = EY_n^* s \rightarrow EY^* s = EY^s$$

for $s = 1, 2$, while Vitali also yields

$$EX_n Y_n = EX_n^* Y_n^* \rightarrow EX^* Y^* = EXY.$$

Therefore

$$Var_{F_n}(X_n) \rightarrow Var_F(X), Var_{F_n}(Y_n) \rightarrow Var_F(Y), \quad (3)$$

and

$$Cov_{F_n}(X_n, Y_n) \rightarrow Cov_F(X, Y). \quad (4)$$

Since we have assumed that $Var_F(X) > 0$ and $Var_F(Y) > 0$, (3) and (4) yield

$$\rho(F_n) = \frac{Cov_{F_n}(X_n, Y_n)}{\sqrt{Var_{F_n}(X_n) \cdot Var_{F_n}(Y_n)}} \rightarrow \frac{Cov_F(X, Y)}{\sqrt{Var_F(X) \cdot Var_F(Y)}} = \rho(F);$$

i.e. ρ is continuous on $\mathcal{F}_{r,M}$ at any F with positive variances.

It is interesting to note that the hypothesis $\{F_n\} \subset \mathcal{F}_{r,M}$ *cannot be weakened* to $\{F_n\} \subset \mathcal{F}_{2,M}$ (and hence it can also not be weakened to the still larger class $\mathcal{F}_{2,\infty}$). Here is a counterexample. Let F be a d.f. on R^2 with $EX = 0 = EY$ and $EX^2 = 1 = EY^2$, and $\rho(F) < 1$ where $(X, Y) \sim F$. Let $M > 1$ be a big number, and consider the class

$$\mathcal{F}_{2,M} = \{F \text{ on } R^2 : E_F X^2 \leq M, E_F Y^2 \leq M\}.$$

Let $a_n, b_n > 0$; we will specify them in terms of M shortly. Consider the sequence of d.f.'s $\{F_n\} \subset \mathcal{F}_{2,M}$ defined by

$$F_n(x, y) = \left(1 - \frac{1}{n}\right)F(x, y) + \frac{1}{2n}\delta_{(a_n, b_n)} + \frac{1}{2n}\delta_{(-a_n, -b_n)}.$$

Then for any bounded and continuous function $\psi : R^2 \rightarrow R$,

$$\begin{aligned} \int \psi dF_n &= \left(1 - \frac{1}{n}\right) \int \psi dF + \frac{1}{2n}\psi(a_n, b_n) + \frac{1}{2n}\psi(-a_n, -b_n) \\ &\rightarrow \int \psi dF, \end{aligned}$$

so $F_n \rightarrow_d F$. Furthermore, with $(X_n, Y_n) \sim F_n$,

$$EX_n = (1 - 1/n)EX = 0, EY_n = 0,$$

$$EX_n^2 = (1 - 1/n)EX^2 + \frac{a_n^2}{n} = (1 - 1/n)M + \frac{a_n^2}{n} = M$$

if $a_n^2 = n\{M - (1 - 1/n)\}$. Similarly,

$$EY_n^2 = (1 - 1/n)EY^2 + \frac{b_n^2}{n} = M$$

if $b_n^2 = n\{M - (1 - 1/n)\}$. With these choices of a_n and b_n ,

$$Cov(X_n, Y_n) = (1 - 1/n)Cov(X, Y) + \frac{a_n b_n}{n},$$

$$\begin{aligned} \rho(F_n) &= \frac{Cov(X_n, Y_n)}{\sqrt{Var(X_n)Var(Y_n)}} \\ &= \frac{(1 - 1/n)Cov(X, Y) + M - (1 - 1/n)}{\sqrt{M^2}} \\ &\rightarrow \frac{\rho(F) + M - 1}{M} \neq \rho(F). \end{aligned}$$

Thus $\rho(F)$ is not continuous on $\mathcal{F}_{2,M}$.