

# On the relative efficiency of the intent-to-treat Wilcoxon–Mann–Whitney test in the presence of noncompliance

By LU MAO 

*Department of Biostatistics and Medical Informatics, School of Medicine and Public Health,  
University of Wisconsin–Madison, Madison, Wisconsin 53726, U.S.A.*

[lmao@biostat.wisc.edu](mailto:lmao@biostat.wisc.edu)

## SUMMARY

A general framework is set up to study the asymptotic properties of the intent-to-treat Wilcoxon–Mann–Whitney test in randomized experiments with nonignorable noncompliance. Under location-shift alternatives, the Pitman efficiencies of the intent-to-treat Wilcoxon–Mann–Whitney and  $t$  tests are derived. It is shown that the former is superior if the compliers are more likely to be found in high-density regions of the outcome distribution or, equivalently, if the noncompliers tend to reside in the tails. By logical extension, the relative efficiency of the two tests is sharply bounded by least and most favourable scenarios in which the compliers are segregated into regions of lowest and highest density, respectively. Such bounds can be derived analytically as a function of the compliance rate for common location families such as Gaussian, Laplace, logistic and  $t$  distributions. These results can help empirical researchers choose the more efficient test for existing data, and calculate sample size for future trials in anticipation of noncompliance. Results for nonadditive alternatives and other tests follow along similar lines.

*Some key words:* Causal inference; Instrumental variable; Nonignorable missingness; Pitman efficiency; Randomized controlled trial;  $t$  test.

## 1. MOTIVATION

The asymptotic properties of the two-sample Wilcoxon–Mann–Whitney test (Wilcoxon, 1945; Mann & Whitney, 1947) have been thoroughly investigated in the literature, especially in comparison with the  $t$  test (Sidak et al., 1999). It is well known that the relative performance of the Wilcoxon–Mann–Whitney and  $t$  tests depends on the shape of the underlying outcome distribution, with the former being superior under heavy-tailed distributions. In fact, the asymptotic relative efficiency of the two tests has been explicitly derived for various outcome distributions, such as the Gaussian, logistic and Laplace families, under location-shift alternatives (see, e.g., van der Vaart, 1998, Ch. 14).

The situation may be different in the presence of nonignorable, i.e., informative, noncompliance to group assignment (Angrist et al., 1996). In such settings, an intent-to-treat test that compares the subjects according to their randomization status is commonly employed to circumvent the selection bias in the treatment received (Gupta, 2011). Because of the unknown relationship between the compliance status and potential outcome, however, the power structures of such tests remain opaque. In this paper, we use the instrumental-variable framework (Angrist et al., 1996; Imbens & Rubin, 1997) to study the asymptotic properties of the intent-to-treat Wilcoxon–Mann–Whitney test in the presence of non-compliance, with particular emphasis on comparison with its  $t$ -test counterpart in detecting additive treatment effects.

## 2. ASYMPTOTIC THEORY

## 2.1. General set-up for intent-to-treat tests

Let  $Y(z, a)$  denote the continuous potential outcome under randomization status  $z$  and received treatment  $a$ , where  $z$  and  $a$  are binary variables, each taking value 1 for the active treatment and value 0 for the control (Rubin, 1978). Similarly, let  $A(z)$  denote the potential treatment under randomization status  $z$ . This notation defines four compliance classes: always-taker if  $A(z) = 1$ ; complier if  $A(z) = z$ ; never-taker if  $A(z) = 0$ ; and defier if  $A(z) = 1 - z$  ( $z = 1, 0$ ). Under the stable unit treatment value assumption (Imbens & Rubin, 2015), the observed treatment and outcome are  $A = A(Z)$  and  $Y = Y(Z, A)$ , respectively. In addition, we assume that randomization has a nontrivial effect on the treatment received, i.e.,  $\text{pr}\{A(1) = 1\} \neq \text{pr}\{A(0) = 1\}$ , and that there is no defier in the population, i.e.,  $A(1) \geq A(0)$  almost surely. In the interest of generality, we relax the usual exclusion restriction assumption that  $Y(1, a) = Y(0, a)$  ( $a = 1, 0$ ) with probability 1 (Angrist et al., 1996) to allow for possible randomization direct effects.

Denote the proportions of the three remaining compliance classes by  $p_a = \text{pr}\{A(1) = A(0) = 1\}$ ,  $p_c = \text{pr}\{A(1) = 1, A(0) = 0\}$  and  $p_n = \text{pr}\{A(1) = A(0) = 0\}$ . By Imbens & Rubin (1997), the observed outcomes in each randomized group can be viewed as a mixture of potential outcomes from the three classes. To be specific, let  $\omega_z(\cdot | \theta)$  denote the cumulative distribution function for  $(Y | Z = z)$  ( $z = 1, 0$ ), where  $\theta$  is some real-valued parameter. Then

$$\begin{aligned}\omega_1(\cdot | \theta) &= p_c v_{1c}(\cdot | \theta) + p_a v_{1a}(\cdot | \theta) + p_n v_{1n}(\cdot | \theta), \\ \omega_0(\cdot | \theta) &= p_c v_{0c}(\cdot | \theta) + p_a v_{0a}(\cdot | \theta) + p_n v_{0n}(\cdot | \theta),\end{aligned}\tag{1}$$

where  $\{Y(z, z) | A(1) = 1, A(0) = 0\} \sim v_{zc}(\cdot | \theta)$ ,  $\{Y(z, 1) | A(1) = A(0) = 1\} \sim v_{za}(\cdot | \theta)$  and  $\{Y(z, 0) | A(1) = A(0) = 0\} \sim v_{zn}(\cdot | \theta)$  ( $z = 1, 0$ ). Suppose that  $\theta = 0$  represents a scenario of the sharp null hypothesis

$$H_0 : Y(z, a) = Y(0, 0) \quad (z, a = 1, 0),$$

so that  $v_{zc}(\cdot | 0) = v_c(\cdot)$ ,  $v_{za}(\cdot | 0) = v_a(\cdot)$  and  $v_{zn}(\cdot | 0) = v_n(\cdot)$  ( $z = 1, 0$ ) for some  $v_c$ ,  $v_a$  and  $v_n$ . Then  $\omega_1(\cdot | 0) = \omega_0(\cdot | 0) = v(\cdot)$ , where  $v = p_c v_c + p_a v_a + p_n v_n$ .

With a random sample of size  $n$ , suppose that a level- $\alpha$  test ( $0 < \alpha < 1$ ) rejects  $H_0$  and accepts  $H_A : \theta > 0$  when  $n^{1/2}|T_n - \mu(0)|/\hat{\sigma}_n > z_{1-\alpha/2}$ , where  $T_n$  is a regular estimator for some function  $\mu(\theta)$  of  $\theta$ ,  $\hat{\sigma}_n^2$  is a consistent estimator for its asymptotic variance, and  $z_{1-\alpha/2} = \Phi^{-1}(1 - \alpha/2)$  with  $\Phi(\cdot)$  denoting the standard normal cumulative distribution function. The asymptotic power of the test can be evaluated analytically under a sequence of contiguous alternatives approaching  $H_0$  at rate  $O(n^{-1/2})$ , say  $\theta = n^{-1/2}h$  for some  $h > 0$ . Under this  $\theta$ , the power of the test tends to  $\Phi(\zeta h - z_{1-\alpha/2})$  for some  $\zeta > 0$  as  $n \rightarrow \infty$  (see van der Vaart, 1998, Ch. 14). The quantity  $\zeta^2$ , commonly called the Pitman efficiency, is given by  $\zeta^2 = \dot{\mu}(0)^2/\sigma_0^2$ , where  $\sigma_0^2$  is the asymptotic limit of  $\hat{\sigma}_n^2$  under  $H_0$  and  $\dot{f}(x) = df(x)/dx$  for any function  $f$ . The Pitman efficiency is inversely proportional to the sample size needed to achieve a given power and is used aptly to measure the asymptotic quality of the test.

Nonparametric intent-to-treat tests can often be formulated as contrasts of empirical distributions. Let  $\hat{\omega}_z(\cdot)$  denote the empirical analogue of  $\omega_z(\cdot | \theta)$  ( $z = 1, 0$ ). Then the Wilcoxon–Mann–Whitney and  $t$  test statistics can be written as  $T_{\text{WMW},n} = \int \hat{\omega}_0(y)\hat{\omega}_1(dy) - 2^{-1}$  (Mann & Whitney, 1947) and  $T_{t,n} = \int y\{\hat{\omega}_1(dy) - \hat{\omega}_0(dy)\}$ , respectively. The estimands of  $T_{\text{WMW},n}$  and  $T_{t,n}$  are thus  $\mu_{\text{WMW}}(\theta) = \int \omega_0(y | \theta)\omega_1(dy | \theta) - 2^{-1}$  and  $\mu_t(\theta) = \int y\{\omega_1(dy | \theta) - \omega_0(dy | \theta)\} = -\int \{\omega_1(y | \theta) - \omega_0(y | \theta)\} dy$ , respectively. Write  $\nabla\eta(\cdot) = \partial\eta(\cdot | \theta)/\partial\theta|_{\theta=0}$  for any distribution function  $\eta$  indexed by  $\theta$ . The following lemma, proved in the Appendix, provides general expressions for the Pitman efficiencies of the two tests.

**LEMMA 1.** *Under (1) with  $\theta = O(n^{-1/2})$ , the Pitman efficiencies for the intent-to-treat Wilcoxon–Mann–Whitney and  $t$  tests are*

$$\zeta_{\text{WMW}}^2 = 12q(1-q) \left[ \int \{\nabla\omega_1(y) - \nabla\omega_0(y)\} \dot{v}(y) dy \right]^2,$$

$$\zeta_t^2 = q(1-q)V(v)^{-1} \left[ \int \{\nabla\omega_1(y) - \nabla\omega_0(y)\} dy \right]^2,$$

respectively, where  $q = E(Z)$  and  $V(v) = \int \{y - \int y^* v(dy^*)\}^2 \dot{v}(y) dy$ .

By Lemma 1, the Pitman efficiencies of the two tests are both functionals of

$$\nabla\omega_1 - \nabla\omega_0 = p_c \{\nabla v_{1c} - \nabla v_{0c}\} + p_a \{\nabla v_{1a} - \nabla v_{0a}\} + p_n \{\nabla v_{1n} - \nabla v_{0n}\}, \quad (2)$$

which mixes the treatment and randomization effects across the compliance classes. Under exclusion restriction, for example, the latter two terms on the right-hand side of (2) vanish so that the difference consists only of the treatment effect on the compliers.

## 2.2. The Wilcoxon–Mann–Whitney test versus the $t$ test under additive effects

For ease of comparison with existing theory under perfect compliance (see, e.g., van der Vaart, 1998, § 14.1), we first focus on constant additive effects. Extensions to general patterns of treatment effects are studied in § 3.1. Specifically, consider the location-shift alternatives

$$H_{A,n} : Y(z, a) = Y(0, 0) + (rz + a)\theta, \quad \theta = n^{-1/2}h, \quad (3)$$

where  $h > 0$  and  $r \in \mathbb{R}$ . That is, the treatment and randomization add  $\theta$  and  $r\theta$ , respectively, to the outcome. The case with exclusion restriction corresponds to  $r = 0$ . Combining (2) with (3), we obtain that  $\nabla\omega_1(y) - \nabla\omega_0(y) = p_c \partial \{v_c(y - \theta - r\theta) - v_c(y)\} / \partial \theta|_{\theta=0} + p_a \partial \{v_a(y - \theta - r\theta) - v_a(y - \theta)\} / \partial \theta|_{\theta=0} + p_n \partial \{v_n(y - r\theta) - v_n(y)\} / \partial \theta|_{\theta=0} = -\{(p_c + r)\dot{v}_c(y) + r\dot{v}_a(y) + \dot{v}_n(y)\} = -\{p_c \dot{v}_c(y) + r\dot{v}(y)\}$ . Substituting this into Lemma 1, we arrive at the next theorem.

**THEOREM 1.** *Under the location-shift contiguous alternatives  $H_{A,n}$  in (3), the Pitman efficiencies for the intent-to-treat Wilcoxon–Mann–Whitney and  $t$  tests are*

$$\zeta_{\text{WMW}}^2(p_c, v) = 12q(1-q)(p_c \langle \dot{v}, \dot{v}_c \rangle + r \langle \dot{v}, \dot{v} \rangle)^2,$$

$$\zeta_t^2(p_c, v) = q(1-q)V(v)^{-1}(p_c + r)^2, \quad (4)$$

respectively, where  $\langle f, g \rangle = \int f(y)g(y) dy$  for any Lebesgue square-integrable functions  $f$  and  $g$ . As a result, the relative efficiency is

$$\mathcal{R}(p_c, v) \equiv \frac{\zeta_{\text{WMW}}^2(p_c, v)}{\zeta_t^2(p_c, v)} = 12V(v) \left( \frac{p_c \langle \dot{v}, \dot{v}_c \rangle + r \langle \dot{v}, \dot{v} \rangle}{p_c + r} \right)^2. \quad (5)$$

Simulations described in the [Supplementary Material](#) show that asymptotic power functions based on (4) approximate the actual power fairly accurately under finite samples.

Either under perfect compliance so that  $p_c = 1$ , or under uninformative noncompliance so that  $v_c = v$ , the far right-hand side of (5) reduces to  $12V(v)\langle \dot{v}, \dot{v} \rangle^2$ , recovering the classical results for the two tests (see, e.g., van der Vaart, 1998, Example 14.13). In that sense, the impact of informative noncompliance shows in the difference between  $\mathcal{R}(p_c, v)$  and  $\mathcal{R}(1, v) = 12V(v)\langle \dot{v}, \dot{v} \rangle^2$ , which, in turn because of (5), is determined by the difference between  $\langle \dot{v}, \dot{v}_c \rangle$  and  $\langle \dot{v}, \dot{v} \rangle$ . The following proposition says that this difference depends crucially on how strongly being a complier is correlated with the height of the outcome density.

**PROPOSITION 1.** *Let  $Y = Y(0, 0)$  and  $\delta = I\{A(1) = 1, A(0) = 0\}$ . Then  $\langle \dot{v}, \dot{v}_c \rangle - \langle \dot{v}, \dot{v} \rangle = p_c^{-1} \text{cov}\{\delta, \dot{v}(Y)\}$  so that*

$$\mathcal{R}(p_c, v) = 12V(v) [\langle \dot{v}, \dot{v} \rangle + (p_c + r)^{-1} \text{cov}\{\delta, \dot{v}(Y)\}]^2. \quad (6)$$

*Proof.* With  $Y \sim \nu$ , since  $E(\delta - p_c) = 0$  we have that

$$\text{cov}\{\delta, \dot{\nu}(Y)\} = E\{(\delta - p_c)\dot{\nu}(Y)\} = p_c E\{\dot{\nu}(Y) \mid \delta = 1\} - p_c \langle \dot{\nu}, \dot{\nu} \rangle = p_c (\langle \dot{\nu}, \dot{\nu}_c \rangle - \langle \dot{\nu}, \dot{\nu} \rangle).$$

The identity in (6) then follows from (5).  $\square$

Using Proposition 1, we obtain via standard calculations that

$$\mathcal{R}(p_c, \nu) - \mathcal{R}(1, \nu) = 12V(\nu)(p_c + r)^{-2} \text{cov}\{\delta, \dot{\nu}(Y)\} \{p_c \langle \dot{\nu}, \dot{\nu}_c \rangle + (p_c + 2r) \langle \dot{\nu}, \dot{\nu} \rangle\}.$$

This means that the sign of  $\mathcal{R}(p_c, \nu) - \mathcal{R}(1, \nu)$  coincides with that of  $\text{cov}\{\delta, \dot{\nu}(Y)\}$ . Hence, roughly speaking, noncompliance favours the Wilcoxon–Mann–Whitney test over the  $t$  test if the compliers are more likely to appear in high-density regions or, conversely, if noncompliers are more likely to appear in the tails. Consider a simple situation where  $\dot{\nu}$  is unimodal and symmetric about zero and  $\text{pr}(\delta = 1 \mid Y = y)$  is nonincreasing in  $|y|$ ; then one can easily show that  $\text{cov}\{\delta, \dot{\nu}(Y)\} \geq 0$  so that  $\mathcal{R}(p_c, \nu) \geq \mathcal{R}(1, \nu)$ .

The correspondence between the relative location of compliers/noncompliers and the relative efficiency of the two tests has an intuitive explanation rooted in classical theory, which tells us that the rank-based Wilcoxon–Mann–Whitney test is more robust against noise in the tails. In an intent-to-treat analysis, the compliers are the sole carrier of the treatment signal between the randomized groups, with the noncompliers contributing nothing but random noise. Hence, if the noise contributors tend to reside in the tails, their influence is better contained by the Wilcoxon–Mann–Whitney test than by the  $t$  test, allowing the former to pick out the signal more easily.

This qualitative relationship can be carried over to quantitative limits. By (6), the range of  $\mathcal{R}(p_c, \nu)$  depends on that of  $\text{cov}\{\delta, \dot{\nu}(Y)\}$ . Intuitively,  $\text{cov}\{\delta, \dot{\nu}(Y)\}$  should be maximized or minimized when the compliers are selected according to the largest or smallest density values, respectively. To make this precise, denote the range of  $\nu_c$  under fixed  $\nu$  and  $p_c$  by

$$\mathcal{P}(p_c \mid \nu) = \{\nu_c : \nu = p_c \nu_c + (1 - p_c)\eta \text{ for some distribution } \eta\}.$$

We can obtain bounds for  $\{\langle \dot{\nu}, \dot{\nu}_c \rangle : \nu_c \in \mathcal{P}(p_c \mid \nu)\}$  using the following lemma, whose proof can be found in the [Supplementary Material](#).

LEMMA 2. *Given a  $\nu$ -integrable function  $g$ , let*

$$\begin{aligned} \dot{\nu}_g^{(L)}(y) &= p_c^{-1} I\{g(y) < l\} \dot{\nu}(y) + c^{(L)} I\{g(y) = l\} \dot{\nu}(y), \\ \dot{\nu}_g^{(U)}(y) &= p_c^{-1} I\{g(y) > u\} \dot{\nu}(y) + c^{(U)} I\{g(y) = u\} \dot{\nu}(y), \end{aligned} \quad (7)$$

where  $l, u \in \mathbb{R}$  and  $c^{(L)}, c^{(U)} \in [0, p_c^{-1}]$  are such that  $\int \dot{\nu}_g^{(L)}(y) dy = \int \dot{\nu}_g^{(U)}(y) dy = 1$ . Then

$$\{\langle g, \dot{\nu}_c \rangle : \nu_c \in \mathcal{P}(p_c \mid \nu)\} = [\langle g, \dot{\nu}_g^{(L)} \rangle, \langle g, \dot{\nu}_g^{(U)} \rangle].$$

Replacing  $g$  in Lemma 2 with  $\dot{\nu}$ , we obtain the range of  $\langle \dot{\nu}, \dot{\nu}_c \rangle$  and, by (5), that of  $\mathcal{R}(p_c, \nu)$ .

THEOREM 2. *Under  $H_{A,n}$  in (3) with fixed  $(p_c, \nu)$ , the range of the relative efficiency between the intent-to-treat Wilcoxon–Mann–Whitney and  $t$  tests is*

$$12V(\nu) \left( \frac{p_c \langle \dot{\nu}, \dot{\nu}_c^{(L)} \rangle + r \langle \dot{\nu}, \dot{\nu} \rangle}{p_c + r} \right)^2 \leq \mathcal{R}(p_c, \nu) \leq 12V(\nu) \left( \frac{p_c \langle \dot{\nu}, \dot{\nu}_c^{(U)} \rangle + r \langle \dot{\nu}, \dot{\nu} \rangle}{p_c + r} \right)^2, \quad (8)$$

where  $\{\dot{\nu}_c^{(L)}, \dot{\nu}_c^{(U)}\}$  is  $\{\dot{\nu}_g^{(L)}, \dot{\nu}_g^{(U)}\}$  defined in (7) with  $g$  replaced by  $\dot{\nu}$ .

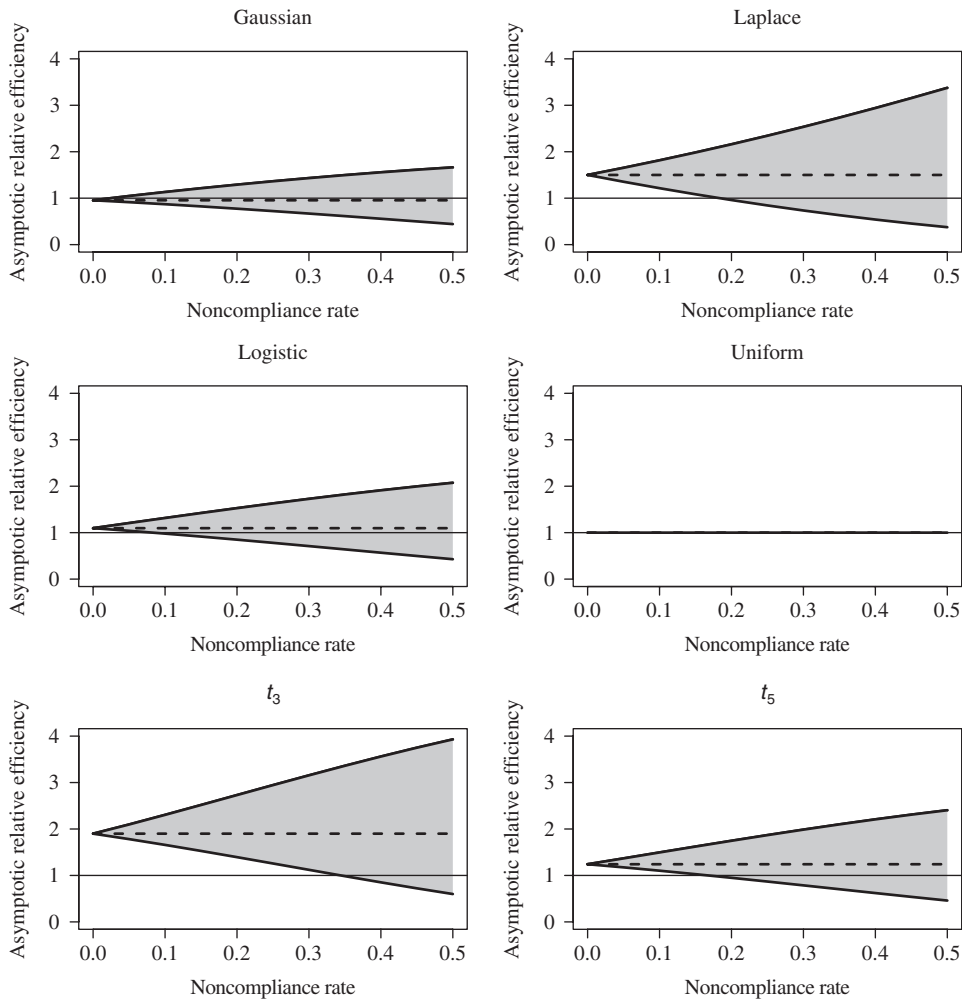


Fig. 1. The range of asymptotic relative efficiency of the intent-to-treat Wilcoxon–Mann–Whitney test versus the  $t$  test against  $H_{A,n}$  under exclusion restriction for different outcome distributions. The dashed line in each panel represents the asymptotic relative efficiency under perfect compliance;  $t_d$  refers to the  $t$  distribution with  $d$  degrees of freedom.

For a unimodal symmetric  $\dot{\nu}$ , it is clear from the above results that the least and most favourable scenarios for the Wilcoxon–Mann–Whitney test are those in which the compliers exclusively occupy the tails or modal regions, respectively. Standard derivation along this line gives us

$$\langle \dot{\nu}, \dot{\nu}_c^{(L)} \rangle = 2p_c^{-1} \int_{-\infty}^{v^{-1}(p_c/2)} \dot{\nu}(y)^2 dy, \quad \langle \dot{\nu}, \dot{\nu}_c^{(U)} \rangle = 2p_c^{-1} \int_0^{v^{-1}\{(1+p_c)/2\}} \dot{\nu}(y)^2 dy. \quad (9)$$

Using (9) in conjunction with (8), we derive the range of  $\mathcal{R}(p_c, \nu)$  as a function of the noncompliance rate  $1 - p_c$  for Gaussian, Laplace, logistic, uniform and  $t$  distributions under exclusion restriction, i.e.,  $r = 0$ . The results are plotted in Fig. 1, with the underlying analytical formulae given in the [Supplementary Material](#). The dashed lines represent the asymptotic relative efficiency values under perfect compliance, which agree with the numbers tabulated in Table 14.12 of [van der Vaart \(1998\)](#). Not surprisingly,  $\mathcal{R}(p_c, \nu) \equiv 1$  for the uniform distribution, whose constant density makes the upper and lower bounds the same. Otherwise, we find by numerical calculation that the lower bounds of  $\mathcal{R}(p_c, \nu)$  under the Laplace, logistic,  $t_3$  and  $t_5$  distributions intersect with the unit horizontal line at  $1 - p_c = 18.4\%$ ,  $8.3\%$ ,  $34.4\%$  and  $16.7\%$ , respectively. This means that the Wilcoxon–Mann–Whitney test is always more efficient

than the  $t$  test for these distributions, as long as the noncompliance rate is below the corresponding threshold.

### 2.3. Practical considerations

The theoretical results in § 2.2 are practically instructive in two ways. First, we can use Theorem 1 to pick the more efficient test for analysing a current trial in the presence of noncompliance. In (5) with  $r = 0$ , for example,  $p_c$  can be estimated by the sample analogue of  $E(A = 1 | Z = 1) - E(A = 1 | Z = 0)$  (Angrist et al., 1996). Moreover,  $\dot{v}(\cdot)$  and  $\dot{v}_c(\cdot)$  can be estimated by a kernel density estimator based on the randomized control group and one using the methods of Imbens & Rubin (1997), respectively. Plugging in these estimates, we obtain an estimate for  $\mathcal{R}(p_c, v)$ , telling us which of the Wilcoxon–Mann–Whitney and  $t$  tests is likely to be more powerful. Intuitively, this data-dependent selection of test should not affect the Type I error rate as it does not use the part of the data that is informative about the treatment effect.

As a demonstration, the above procedure was run on data collected from 11 204 participants in the well-known U.S. National Job Training Partnership Act Study, a randomized trial conducted between 1987 and 1989 to assess the effect of a job-training programme on the trainee's subsequent earnings (Abadie et al., 2002). The estimated compliance rate is  $p_c = 62.7\%$ . As shown in the Supplementary Material, both  $\dot{v}$  and  $\dot{v}_c$  appear heavily bimodal. Using these density estimates, we find that  $\mathcal{R}(p_c, v) \approx 4.9$ , suggesting that the intent-to-treat Wilcoxon–Mann–Whitney test is likely much more powerful than the  $t$  test. This is substantiated by the actual test results, with the former producing a  $p$ -value of 0.001 and the latter 0.023.

Second, one can calculate the sample size for a future trial by inverting the power function with the derived Pitman efficiencies. By Theorem 1 with  $r = 0$ , the Pitman efficiency for the  $t$  test is  $p_c^2$  times that under perfect compliance. Provided that a good estimate for  $p_c$  can be obtained from historical data, the investigator need only inflate the standard sample size by  $p_c^{-2}$ . The same is true for the Wilcoxon–Mann–Whitney test under uninformative noncompliance. With informative noncompliance, and without further knowledge about the compliance structure, the best one can do is use the lower bound on  $\langle \dot{v}, \dot{v}_c \rangle$  constructed by Lemma 2 to obtain a conservative estimate for the sample size. For example, under exclusion restriction with additive treatment effect  $\theta$  in units of standard deviation, it can be shown that the sample sizes needed to achieve a power of  $1 - \beta$  for the Gaussian and Laplace distributions are bounded by

$$\frac{\pi(z_{1-\alpha/2} + z_{1-\beta})^2}{12q(1-q)\Phi\{2^{1/2}\Phi^{-1}(p_c/2)\}^2\theta^2}, \quad \frac{2(z_{1-\alpha/2} + z_{1-\beta})^2}{3q(1-q)p_c^4\theta^2},$$

respectively. As an example, the data from the above job-training study are used as pilot data to calculate the sample size for a new trial. The details are provided in the Supplementary Material.

## 3. EXTENSIONS AND DISCUSSIONS

### 3.1. Inhomogeneous effects

The framework of § 2.1 accommodates all patterns of treatment and randomization effects, to the extent that  $\nabla\omega_1(y) - \nabla\omega_0(y)$  can be derived under the specific contiguous alternatives. Indeed, consider a general alternative hypothesis that transforms the null outcome by a function  $\mathcal{T}$  indexed by treatment and randomization effect parameters  $\theta$  and  $r\theta$ , respectively; that is,

$$H_{\Lambda,n} : Y(z, a) = \mathcal{T}\{Y(0, 0); a\theta, zr\theta\},$$

where  $\mathcal{T}(y; 0, 0) = y$  for all  $y \in \mathbb{R}$ . For example, the location-shift alternative in (3) corresponds to  $\mathcal{T}(y; a\theta, zr\theta) = y + a\theta + zr\theta$ . Multiplicative effects can be formulated as  $\mathcal{T}(y; a\theta, zr\theta) = y \exp(a\theta + zr\theta)$ . Sometimes subjects with lower baseline outcomes stand to gain more from the treatment than do those with already high values. In such cases, it would be appropriate to consider  $\mathcal{T}(y; a\theta, zr\theta) = y + (rz + a)\theta\kappa(y)$ ,

where  $\theta > 0$  and  $\kappa(y)$  is a decreasing function, e.g.,  $\kappa(y) = \{1 + \exp(cy)\}^{-1}$  for some  $c > 0$ . The following proposition, proved in the [Supplementary Material](#), expresses  $\nabla\omega_1(y) - \nabla\omega_0(y)$  in terms of a tangent function of  $\mathcal{T}(y; a\theta, zr\theta)$  as  $\theta \rightarrow 0$ .

LEMMA 3. *Under regularity conditions specified in the [Supplementary Material](#), suppose that  $\partial\mathcal{T}(y; a\theta, zr\theta)/\partial\theta|_{\theta=0} = \dot{\mathcal{T}}_{a,z}(y)$  for some function  $\dot{\mathcal{T}}_{a,z}(y)$ . Write  $\dot{\mathcal{T}}_c(y) = \dot{\mathcal{T}}_{1,1}(y) - \dot{\mathcal{T}}_{0,0}(y)$ ,  $\dot{\mathcal{T}}_a(y) = \dot{\mathcal{T}}_{1,1}(y) - \dot{\mathcal{T}}_{1,0}(y)$  and  $\dot{\mathcal{T}}_n(y) = \dot{\mathcal{T}}_{0,1}(y) - \dot{\mathcal{T}}_{0,0}(y)$ . Then*

$$\nabla\omega_1(y) - \nabla\omega_0(y) = -\{p_c\dot{\mathcal{T}}_c(y)\dot{v}_c(y) + p_a\dot{\mathcal{T}}_a(y)\dot{v}_a(y) + p_n\dot{\mathcal{T}}_n(y)\dot{v}_n(y)\}.$$

Under the inhomogeneous additive model  $\mathcal{T}(y; a\theta, zr\theta) = y + (rz + a)\theta\kappa(y)$ , for example, it is immediate that  $\dot{\mathcal{T}}_{a,z}(y) = (rz + a)\kappa(y)$  so that  $\dot{\mathcal{T}}_c(y) = (1 + r)\kappa(y)$ ,  $\dot{\mathcal{T}}_a(y) = \dot{\mathcal{T}}_n(y) = r\kappa(y)$  and, by Lemma 3,  $\nabla\omega_1(y) - \nabla\omega_0(y) = -\kappa(y)\{p_c\dot{v}_c(y) + r\dot{v}(y)\}$ . Using the last expression in Lemma 1, we obtain that

$$\begin{aligned}\zeta_{\text{WMW}}^2(p_c, v) &= 12q(1 - q) (p_c\langle\kappa, \dot{v}_c\rangle + r\langle\kappa, \dot{v}\rangle)^2, \\ \zeta_t^2(p_c, v) &= q(1 - q)V(v)^{-1} (p_c\langle\kappa, \dot{v}_c\rangle + r\langle\kappa, \dot{v}\rangle)^2.\end{aligned}\tag{10}$$

Taking  $\kappa(y) \equiv 1$  reduces (10) to (4) of Theorem 1. Although the Pitman efficiencies under nonconstant  $\kappa(\cdot)$  appear a bit more complex, one can nonetheless use Lemma 2 to derive bounds for them. The case with multiplicative effects is considered in the [Supplementary Material](#).

### 3.2. General intent-to-treat tests

The intent-to-treat Wilcoxon–Mann–Whitney and  $t$  tests can both be formulated as contrasts between the empirical distributions in the randomized groups. As seen from the main steps in § 2, what determines their respective asymptotic properties is the functional derivative, or more precisely the Hadamard derivative (van der Vaart, 1998, Ch. 20), of the contrast in question. This helps us generalize the previous results to a broad class of nonparametric intent-to-treat tests. The following proposition is proved in the [Supplementary Material](#).

PROPOSITION 2. *Consider an intent-to-treat test that rejects  $H_0$  if  $n^{1/2}|\mathcal{H}(\hat{\omega}_1, \hat{\omega}_0)|/\hat{\sigma}_n > z_{1-\alpha/2}$ , where  $\mathcal{H}(\cdot, \cdot)$  is a functional satisfying  $\mathcal{H}(\omega_1, \omega_0) = -\mathcal{H}(\omega_0, \omega_1)$  for all  $\omega_1$  and  $\omega_0$ , and  $\hat{\sigma}_n^2$  is a consistent variance estimator for  $n^{1/2}|\mathcal{H}(\hat{\omega}_1, \hat{\omega}_0)|$ . Let  $\dot{\mathcal{H}}_v$  denote the Hadamard derivative of  $\mathcal{H}(\omega_1, \omega_0)$  with respect to  $\omega_1$  at  $\omega_1 = \omega_0 = v$ . Then the Pitman efficiency of the test is*

$$\zeta_{\mathcal{H}}^2 = q(1 - q) \left[ \int \dot{\mathcal{H}}_v^2 \{I(y \leq \cdot) - v(\cdot)\} \dot{v}(y) dy \right]^{-1} \dot{\mathcal{H}}_v^2(\nabla\omega_1 - \nabla\omega_0).\tag{11}$$

In (11), the expressions for  $\nabla\omega_1 - \nabla\omega_0$  and  $\dot{\mathcal{H}}_v$  are determined by the contiguous alternatives and the specific test, respectively. To recover the results of Lemma 1 for the Wilcoxon–Mann–Whitney and  $t$  tests under location-shift alternatives, simply plug in  $\dot{\mathcal{H}}_v\{\psi(\cdot)\} = -\int \psi(y)\dot{v}(y) dy$  and  $-\int \psi(y) dy$ , respectively, along with  $\nabla\omega_1 - \nabla\omega_0 = -\{p_c\dot{v}_c(y) + r\dot{v}(y)\}$ . For a new example, consider the  $\tau$ -quantile test based on  $\mathcal{H}(\hat{\omega}_1, \hat{\omega}_0) = \hat{\omega}_1^{-1}(\tau) - \hat{\omega}_0^{-1}(\tau)$  ( $0 < \tau < 1$ ). The Hadamard derivative of this  $\mathcal{H}$  is  $\dot{\mathcal{H}}_v\{\psi(\cdot)\} = -\dot{v}(z_{v,\tau})^{-1}\psi(z_{v,\tau})$ , where  $z_{v,\tau} = v^{-1}(\tau)$  (e.g., van der Vaart, 1998, Lemma 21.3). A straightforward derivation in the [Supplementary Material](#) shows that the Pitman efficiency under the location-shift alternatives is  $\zeta_{\tau}^2 = q(1 - q)\tau^{-1}(1 - \tau)^{-1} \{p_c\dot{v}_c(z_{v,\tau}) + r\dot{v}(z_{v,\tau})\}^2$ . As one might have guessed, the  $\tau$ -quantile test is powerful when the neighbourhood of the corresponding quantile is densely populated by compliers, driving up the value of  $\dot{v}_c(z_{v,\tau})$ . Formal analysis and comparison with the Wilcoxon–Mann–Whitney and  $t$  tests is left to the interested reader.



## ACKNOWLEDGEMENT

This research was supported by the U.S. National Science Foundation and National Institutes of Health. The author very much appreciates the helpful comments from the editor, associate editor and two referees.

## SUPPLEMENTARY MATERIAL

[Supplementary Material](#) includes technical and numerical results.

## APPENDIX

*Proof of Lemma 1.* The null asymptotic variance of  $T_{\text{WMW},n}$  is the distribution-free constant  $12^{-1}q^{-1}(1-q)^{-1}$  by Example 14.13 in [van der Vaart \(1998\)](#). On the other hand,

$$\dot{\mu}_{\text{WMW}}(0) = \frac{\partial}{\partial \theta} \bigg|_{\theta=0} \int \omega_0(y | \theta) \omega_1(dy | \theta) = - \int \{\nabla \omega_1(y) - \nabla \omega_0(y)\} \dot{v}(y) dy,$$

where the second equality follows from integration by parts. The case for  $\zeta_t^2$  can be similarly derived.  $\square$

## REFERENCES

- ABADIE, A., ANGRIST, J. & IMBENS, G. (2002). Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings. *Econometrica* **70**, 91–117.
- ANGRIST, J. D., IMBENS, G. W. & RUBIN, D. B. (1996). Identification of causal effects using instrumental variables. *J. Am. Statist. Assoc.* **91**, 444–55.
- GUPTA, S. K. (2011). Intention-to-treat concept: A review. *Perspect. Clin. Res.* **2**, 109–12.
- IMBENS, G. W. & RUBIN, D. B. (1997). Estimating outcome distributions for compliers in instrumental variables models. *Rev. Econ. Stud.* **64**, 555–74.
- IMBENS, G. W. & RUBIN, D. B. (2015). *Causal Inference in Statistics, Social, and Biomedical Sciences*. Cambridge: Cambridge University Press.
- MANN, H. B. & WHITNEY, D. R. (1947). On a test of whether one of two random variables is stochastically larger than the other. *Ann. Math. Statist.* **18**, 50–60.
- RUBIN, D. B. (1978). Bayesian inference for causal effects: The role of randomization. *Ann. Statist.* **6**, 34–58.
- SIDAK, Z., SEN, P. K. & HAJEK, J. (1999). *Theory of Rank Tests*. London: Academic Press.
- VAN DER VAART, A. W. (1998). *Asymptotic Statistics*. Cambridge: Cambridge University Press.
- WILCOXON, F. (1945). Individual comparisons by ranking methods. *Biomet. Bull.* **1**, 80–3.

[Received on 4 October 2020. Editorial decision on 13 October 2021]