# MULTIPLE IMPROVEMENTS OF MULTIPLE IMPUTATION LIKELIHOOD RATIO TESTS

#### By Kin Wai Chan and Xiao-Li Meng

The Chinese University of Hong Kong and Harvard University

Multiple imputation (MI) inference handles missing data by first properly imputing the missing values m times, and then combining the results from the m complete-data analyses. However, the existing method for combining likelihood ratio tests has multiple defects: (i) the combined test statistic can be negative in practice but its null distribution is approximated by a standard F distribution; (ii) it is not invariant to re-parametrization; (iii) it fails to ensure monotonic power due to its use of an inconsistent estimator of the fraction of missing information (FMI) under the alternative hypothesis; and (iv) it requires non-trivial access to the likelihood ratio test statistic as a function of estimated parameters instead of datasets. This paper shows, via both theoretical derivations and empirical investigations, that essentially all of these problems can be straightforwardly addressed if we are willing to perform an additional likelihood ratio test by stacking the m completed datasets as one big completed dataset. A particularly intriguing finding is that the FMI itself can be estimated consistently by a likelihood ratio statistic for testing whether the m completed datasets produced by MI can be regarded effectively as samples coming from a common model. Practical guidelines are provided based on an extensive comparison of existing MI tests. Intrigued issues regarding nuisance parameters are also discussed.

1. Historical Successes and Failures. Missing-data problems are ubiquitous in practice, to the extent that the absence of any missingness often is a strong indication that the data have been pre-processed or manipulated in some way (e.g., Blocker and Meng, 2013). Multiple imputation (MI) (Rubin, 1978, 2004) has been a preferred method by many practitioners, especially those who are ill-equipped to handle missingness on their own, due to lack of information or skills or resources. MI relies on the data collector (e.g., a census bureau) to build a reliable imputation model to fill in the missing data  $m(\geq 2)$  times, so the users of the data can apply their favorite software or procedures that are designed to handle complete data,

AMS 2000 subject classifications: Primary 62D05; Secondary 62F03, 62E20.

Keywords and phrases: Fraction of missing information; missing data; invariant test; monotonic power; robust estimation.

Meng thanks NSF and JTF for partial financial support. He is also embarrassed by the research immaturity displayed in his thesis, and is very thankful to Keith's (Kin Wai) creativity and diligence which led to the beautiful remedies presented here.

and do so m times. MI inference, e.g., hypothesis testing, is then performed by appropriately combining these m complete-data results.

Although MI was designed initially for public-use datasets, over the past 30 years or so, it has become a method of choice for handling missing data in general, because it separates the handling missingness from conducting analysis (e.g., Tu et al., 1993; Rubin, 1996, 2004; Schafer, 1999; King et al., 2001; Peugh and Enders, 2004; Kenward and Carpenter, 2007; Rose and Fraser, 2008; Holan et al., 2010; Kim and Yang, 2017). Software routines for performing MI are now available in R (van Buuren and Groothuis-Oudshoorn, 2011; Su et al., 2011), Stata (Royston and White, 2011), SAS (Berglund and Heeringa, 2014) and SPSS; also see Harel and Zhou (2007) and Horton and Kleinman (2007) for summaries on software that utilize MI.

This convenient separation, however, creates the thorny issue of uncongeniality, i.e., the incompatibility between the imputation model and the subsequent analysis procedures (Meng, 1994a). This issue is examined in detail by Xie and Meng (2017), which shows that uncongeniality is easiest to deal with when the imputer's model is more saturated than the user's model/procedure, and when the user is conducting efficient analysis, such as likelihood inference. The current paper, therefore, focuses on conducting MI likelihood ratio tests (LRTs), assuming the imputation model is sufficiently saturated to render the validity of the common assumptions, which we shall review, made in the literature about conducting LRTs with MI.

Like many hypothesis testing procedures in common practice, the exact null distributions of various MI test statistics, LRTs or not, are intractable. This intractability is not computational, but rather statistical due to the well-known issue of nuisance parameter, that is, the lack of pivotal quantity, as highlighted, historically, by the Behrens-Fisher problem (Wallace, 1980). Indeed, the nuisance parameter in the MI context is the so-called "the fraction of missing information" (FMI), which is determined by the ratio of the between-imputation variance to within-imputation variance (and its multivariate counterparts), and hence the challenge we face is almost identical to the one faced by the Behrens-Fisher problem, as shown in Meng (1994b).

An added challenge in the MI context is that the user's complete-data procedures can be very restrictive. What is available to the user could vary from the entire likelihood function, to point estimators such as MLE and Fisher information, to a single p-value. Therefore, there have been a variety of procedures proposed in the literature, depending on what quantities we assume the user has access to, as we shall review shortly.

Among them, a promising idea was to directly combine LRT statistics. However, the execution of this idea as presented in Meng and Rubin (1992) relied too heavily on the usual asymptotic equivalence (in terms of the data size, not the number of imputations, m) between the LRT and Wald test under the null. Its asymptotic validity, unfortunately, does not protect it from quick deterioration for small data sizes, such as delivering negative "F test statistic" or FMI. Worst of all, the test can have essentially zero power because the estimator of FMI can be badly inconsistent under some alternative hypotheses. In addition, the combining rule of Meng and Rubin (1992) requires the user to have access to the LRT as a function of parameter values, not just as a function of the data. The former one is often unavailable from standard software packages. This defective MI LRT, however, has been adopted by textbooks (e.g., van Buuren S, 2012; Kim and Shao, 2013) and popular software, e.g., the function pool.compare in the R package mice (van Buuren and Groothuis-Oudshoorn, 2011), the function testModels in the R package mitml (Grund et al, 2017), the function milrtest (Medeiros, 2008) in the Stata module mim (Carlin et al, 2008).

To minimize the negative impact of this defective LRT test, this paper derives MI LRTs that are free of the defects as outlined in the abstract and detailed in  $\S$  1.3 below. We achieve this mainly by switching the order of two main operators in the combining rule of Meng and Rubin (1992): Maximizing the average of the m log-likelihoods instead of averaging the maximizers of them. This switching, guided by the likelihood principle, automatically renders positivity, invariance and monotonic power. Other judicious uses of the likelihood functions permit us to overcome the remaining defects.

The remainder of Section 1 provides background and notation. Section 2 then discusses the defects of the existing MI LRT and our remedies. Section 3 investigates computational requirements for our proposals, including theoretical considerations and comparisons. In particular, Algorithm 1 of Section 3.1 computes our most recommended test. Section 4 provides empirical evidence with simulated and real data. Section 5 calls for future work. Appendices A–C supplement with proofs, and additional investigations.

1.1. Notation and Complete-data Tests. Let  $X_{\text{obs}}$  and  $X_{\text{mis}}$  be, respectively, the observed and missing parts of an intended complete dataset  $X = X_{\text{com}} = \{X_{\text{obs}}, X_{\text{mis}}\}$ , which consists of n observations. Denote the sampling model — probability or density, depending on the data type — of X by  $f(\cdot \mid \psi)$ , where  $\psi \in \Psi \subset \mathbb{R}^h$  is a vector of parameters. The goal is to test  $H_0: \theta = \theta_0$  when only  $X_{\text{obs}}$  is available, where  $\theta = \theta(\psi) \in \Theta \subset \mathbb{R}^k$  is a function of  $\psi$ , and  $\theta_0$  is a specified vector. For simplicity, we will focus on the standard two-sided alternative, but our approach adapts to general complete-data LRTs. Denote the true values of  $\psi$  and  $\theta$  by  $\psi^*$  and  $\theta^*$ . Here

we assume  $X_{\rm obs}$  is rich enough that the missing data mechanism is ignorable (Rubin, 1976), or it has been properly incorporated into the imputation model by the imputer, who may have access to additional confidential data.

Let  $\hat{\theta} = \hat{\theta}(X)$  and  $U = U_{\theta} = U_{\theta}(X)$  be respectively the complete-data MLE of  $\theta$  and an efficient estimator of  $\text{Var}(\hat{\theta})$  (e.g., the inverse of the observed Fisher information). Also, let  $\hat{\psi}_0 = \hat{\psi}_0(X)$  and  $\hat{\psi} = \hat{\psi}(X)$  be respectively the  $H_0$ -constrained and unconstrained complete-data MLEs of  $\psi$ , and  $U_{\psi} = U_{\psi}(X)$  be an efficient estimator of  $\text{Var}(\hat{\psi})$ . For testing  $H_0$  against  $H_1$ , the common choices include the Wald statistic  $D_{\text{W}} = d_{\text{W}}(\hat{\theta}, U)/k$  and the LRT statistic  $D_{\text{L}} = d_{\text{L}}(\hat{\psi}_0, \hat{\psi} \mid X)/k$ , where

$$d_{\mathbf{W}}(\widehat{\theta}, U) = (\widehat{\theta} - \theta_0)^{\mathsf{T}} U^{-1}(\widehat{\theta} - \theta_0), \qquad d_{\mathbf{L}}(\widehat{\psi}_0, \widehat{\psi} \mid X) = 2\log \frac{f(X \mid \widehat{\psi})}{f(X \mid \widehat{\psi}_0)}.$$

Under regularity conditions (RCs), such as those in § 4.2.2 and § 4.4.2 of Serfling (2001) when the rows of X are independent and identically distributed, we have the following classical results.

PROPERTY 1.1. Under  $H_0$ , (i)  $D_W \Rightarrow \chi_k^2/k$  and  $D_L \Rightarrow \chi_k^2/k$ ; and (ii)  $n(D_W - D_L) \stackrel{\text{pr}}{\to} 0$  as  $n \to \infty$  where " $\Rightarrow$ " and " $\stackrel{\text{pr}}{\to}$ " denote convergence in distribution and in probability, respectively.

Testing  $\theta = \theta_0$  based on  $X_{\rm obs}$  is more involved. For MI, let  $X^{\ell} = \{X_{\rm obs}, X_{\rm mis}^{\ell}\}$ ,  $\ell = 1, \ldots, m$ , be the m completed datasets, where  $X_{\rm mis}^{\ell}$  are drawn conditionally independently from a proper imputation model given  $X_{\rm obs}$ ; see Rubin (2004). We then carry out a complete-data estimation or testing procedure on  $X^{\ell}, \ell = 1, \ldots, m$ , resulting in a set of m quantities. The so-called MI inference is to appropriately combine them to obtain a single answer.

1.2. MI Wald Test and Fraction of Missing Information. Let  $d_{\mathbf{W}}^{\ell} = d_{\mathbf{W}}(\widehat{\theta}^{\ell}, U^{\ell}), \ \widehat{\theta}^{\ell} = \widehat{\theta}(X^{\ell})$  and  $U^{\ell} = U(X^{\ell})$  be the imputed counterparts of  $d_{\mathbf{W}}(\widehat{\theta}, U), \ \widehat{\theta}$  and U, respectively, for each  $\ell$ . Also, write their averages as

(1.1) 
$$\overline{d}_{W} = \frac{1}{m} \sum_{\ell=1}^{m} d_{W}^{\ell}, \qquad \overline{\theta} = \frac{1}{m} \sum_{\ell=1}^{m} \widehat{\theta}^{\ell}, \qquad \overline{U} = \frac{1}{m} \sum_{\ell=1}^{m} U^{\ell}.$$

Under congeniality (Meng, 1994a), one can show that asymptotically (Rubin and Schenker, 1986)  $Var(\bar{\theta})$  can be consistently estimated by

(1.2) 
$$T = \overline{U} + (1 + 1/m)B$$
, where  $B = \frac{1}{m-1} \sum_{\ell=1}^{m} (\widehat{\theta}^{\ell} - \overline{\theta})(\widehat{\theta}^{\ell} - \overline{\theta})^{\mathsf{T}}$ 

is known as the between-imputation variance, in contrast to  $\overline{U}$  in (1.1), which measures within-imputation variance. Intriguingly, 2T serves as a universal (estimated) upper bound of  $Var(\overline{\theta})$  under uncongeniality (Xie and Meng, 2017). Under RCs, we have that, as  $m, n \to \infty$ ,

$$n(\overline{U} - \mathcal{U}_{\theta}) \stackrel{\text{pr}}{\to} \mathbf{0}, \qquad n(T - \mathcal{T}_{\theta}) \stackrel{\text{pr}}{\to} \mathbf{0}, \qquad n(B - \mathcal{B}_{\theta}) \stackrel{\text{pr}}{\to} \mathbf{0}$$

for some deterministic matrices  $\mathcal{U}_{\theta}$ ,  $\mathcal{T}_{\theta}$  and  $\mathcal{B}_{\theta} = \mathcal{T}_{\theta} - \mathcal{U}_{\theta}$ , where **0** denotes a matrix of zeros, and the subscript  $\theta$  highlights that these matrices are for estimating  $\theta$ , because there are also corresponding  $\mathcal{T}_{\psi}$ ,  $\mathcal{B}_{\psi}$ ,  $\mathcal{U}_{\psi}$  for the entire parameter  $\psi$ . Similar to  $\overline{U}$ , T and B, we define  $\overline{U}_{\psi}$ ,  $T_{\psi}$  and  $B_{\psi}$  for the component  $\psi$ . Note that if  $\hat{\theta}_{\text{com}}$  and  $\hat{\theta}_{\text{obs}}$  are the MLEs of  $\theta$  based on  $X_{\text{com}}$  and  $X_{\text{obs}}$  (under congeniality), respectively, then  $\mathcal{U}_{\theta} \cong \mathsf{Var}(\hat{\theta}_{\text{com}})$  and  $\mathcal{T}_{\theta} \cong \mathsf{Var}(\hat{\theta}_{\text{obs}})$  as  $n \to \infty$ , where  $A_n \cong B_n$  means that  $A_n - B_n = o_p(A_n + B_n)$ .

The straightforward MI Wald test  $D_{\rm W}(T)=d_{\rm W}(\overline{\theta},T)/k$  is not practical because T is singular when m < k (usually  $3 \le m \le 10$ ). Even when it is not singular, it is usually not a very stable estimator of  $\mathcal{T}_{\theta}$  because m is small. To circumvent this problem, Rubin (1978) adopted the following assumption of equal fraction of missing information (EFMI).

Assumption 1 (EFMI of 
$$\theta$$
). There is  $r \ge 0$  such that  $\mathcal{T}_{\theta} = (1 + r)\mathcal{U}_{\theta}$ .

EFMI clearly is a very strong assumption, implying that the missing data have caused an equal amount of loss of information for estimating every component of  $\theta$ . However, as we shall see shortly, the adoption of this assumption, for the purpose of hypothesis testing, is essentially the same as to summarize the impact of (at least) k nuisance parameters due to FMI by a single nuisance parameter, i.e., the average FMI across different components. How well this reduction strategy works therefore will affect more the power of the test than its validity, as long as we can construct an approximate null distribution that is more robust to the EFMI assumption. The issue of power turns out to be a rather tricky one, because without the reduction strategy we would also lose power when m/k is small or even modest. It is because we simply do not have enough degrees of freedom to estimate all the nuisance parameters well or at all. We will illustrate this point in § 4.2. (To clarify some confusions in literature, r in Assumption 1 is the odds of the missing information, not the FMI, which is f = r/(1 + r).)

Under EFMI, Rubin (2004) replaced T by  $(1 + \widetilde{r}'_{W})\overline{\widetilde{U}}$ , where

$$(1.3) \widetilde{r}'_{W} = \frac{(m+1)}{k(m-1)} (\overline{d}'_{W} - \widetilde{d}'_{W}); \overline{d}'_{W} = \frac{1}{m} \sum_{\ell=1}^{m} d_{W}(\widehat{\theta}^{\ell}, \overline{U}),$$

 $\widetilde{d}'_{\mathrm{W}} = d_{\mathrm{W}}(\overline{\theta}, \overline{U})$ , and the prime " $\ell$ " indicates that  $\overline{U}$  is used instead of individual  $\{U^{\ell}\}_{\ell=1}^{m}$ . Then, Rubin (2004) proposed a simple MI Wald test statistic:

(1.4) 
$$\widetilde{D}'_{W} = \frac{\widetilde{d}'_{W}}{k(1 + \widetilde{r}'_{W})}.$$

The intuition behind (1.3)–(1.4) is important because the forms here are the building blocks for virtually all the subsequent developments. The "obvious" Wald test statistic  $\tilde{d}'_{\rm W}/k$  is too large (compared to the usual  $\chi_k^2/k$ ) because it fails to take into account of missing information. The  $(1+\tilde{r}'_{\rm W})$  factor attempts to correct this, with the amount of correction determined by the (average) amount of between-imputation variance relative to the within-imputation variance. Expression (1.3) shows that this relative amount can be estimated by contrasting the average of individual Wald statistics and the Wald statistic based on an average of individual estimates. Using the difference between "average of functions" and "function of average", namely,

$$(1.5) Ave{G(x)} - G(Ave{x})$$

is a common practice, e.g.,  $G(x) = x^2$  for variance; see Meng (2002).

Since the exact null distribution of  $\widetilde{D}'_{W}$  is intractable, Li *et al.* (1991b) proposed to approximate it by  $F_{k,\widetilde{\mathrm{df}}(\widetilde{r}'_{W},k)}$ , the F distribution with degrees of freedom k and  $\widetilde{\mathrm{df}}(\widetilde{r}'_{W},k)$ , where, denoting  $K_{m}=k(m-1)$ ,

(1.6) 
$$\widetilde{\mathrm{df}}(\boldsymbol{r}_m, k) = \begin{cases} 4 + (K_m - 4)\{1 + (1 - 2/K_m)/\boldsymbol{r}_m\}^2, & \text{if } K_m > 4; \\ (m - 1)(1 + 1/\boldsymbol{r}_m)^2(k + 1)/2, & \text{otherwise.} \end{cases}$$

This approximation assumes n is sufficiently large so that the standard asymptotic  $\chi^2$  distribution in Property 1.1 can be used. If n is small, the small sample degree of freedom in Barnard and Rubin (1999) should be used.

1.3. The Current MI Likelihood Ratio Test and Its Defect. Let  $d_{\rm L}^\ell = d_{\rm L}(\hat{\psi}_0^\ell,\hat{\psi}^\ell\mid X^\ell),~\hat{\psi}_0^\ell=\hat{\psi}_0(X^\ell)$  and  $\hat{\psi}^\ell=\hat{\psi}(X^\ell)$  be the imputed counterparts of  $d_{\rm L}(\hat{\psi}_0,\hat{\psi}\mid X),~\hat{\psi}_0$  and  $\hat{\psi}$ , respectively, for each  $\ell$ . Let their averages be

(1.7) 
$$\overline{d}_{L} = \frac{1}{m} \sum_{\ell=1}^{m} d_{L}^{\ell}, \quad \overline{\psi}_{0} = \frac{1}{m} \sum_{\ell=1}^{m} \widehat{\psi}_{0}^{\ell}, \quad \overline{\psi} = \frac{1}{m} \sum_{\ell=1}^{m} \widehat{\psi}^{\ell}.$$

Similar to  $\widetilde{r}'_{W}$ , Meng and Rubin (1992) proposed to estimate  $r_{m}$  by

$$(1.8) \quad \widetilde{r}_{L} = \frac{m+1}{k(m-1)} (\overline{d}_{L} - \widetilde{d}_{L}), \quad \text{where} \quad \widetilde{d}_{L} = \frac{1}{m} \sum_{\ell=1}^{m} d_{L}(\overline{\psi}_{0}, \overline{\psi} \mid X^{\ell}),$$

and hence it is again in the form of (1.5). Computation of  $\tilde{r}_L$  requires users to have access to (i) a subroutine for  $(X, \psi_0, \psi) \mapsto d_L(\psi_0, \psi \mid X)$ , and (ii) the estimates  $\hat{\psi}_0^{\ell}$  and  $\hat{\psi}^{\ell}$ , rather than the matrices  $\overline{U}$  and B. Therefore computing  $\tilde{r}_L$  is easier than computing  $\tilde{r}_W'$ . The resulting MI LRT is

(1.9) 
$$\widetilde{D}_{L} = \frac{\widetilde{d}_{L}}{k(1 + \widetilde{r}_{L})},$$

whose null distribution can be approximated by  $F_{k,\widetilde{\mathrm{df}}(\widetilde{r}_1,k)}$ .

The main theoretical justification (and motivation) was the asymptotic equivalence between the complete-data Wald test statistic and LRT statistic under the null, as stated in Property 1.1. This equivalence permitted the replacement of  $\overline{d}'_{W}$  and  $\widetilde{d}'_{W}$  in (1.3) respectively by  $\overline{d}_{L}$  and  $\widetilde{d}_{L}$  in (1.8). However, this is also where the problems lie.

First, with finite samples,  $0 \leq \tilde{d}_L \leq \bar{d}_L$  is not guaranteed, consequently nor is  $\widetilde{D}_L \geq 0$  or  $\widetilde{r}_L \geq 0$ . Since  $\widetilde{D}_L$  is referred to an F distribution and  $\widetilde{r}_L$  estimates  $r_m \geq 0$ , clearly negative values of  $\widetilde{D}_L$  or  $\widetilde{r}_L$  will cause trouble.

Second,  $\widetilde{D}_{L}$  is not invariant to re-parameterization of  $\psi$ . For each individual LRT statistic  $d_{L}^{\ell}$  and bijective map g such that  $\varphi = g(\psi)$ , we have

$$(1.10) d_{\mathcal{L}}^{\ell} = d_{\mathcal{L}}(\widehat{\psi}_{0}^{\ell}, \widehat{\psi}^{\ell} \mid X^{\ell}) = d_{\mathcal{L}}(g^{-1}(\widehat{\varphi}_{0}^{\ell}), g^{-1}(\widehat{\varphi}^{\ell}) \mid X^{\ell}),$$

where  $\widehat{\varphi}_0^{\ell}$  and  $\widehat{\varphi}^{\ell}$  are the constrained and unconstrained MLEs of  $\varphi$  based on  $X^{\ell}$ . However, the MI LRT statistic  $\widetilde{d}_L$  no longer has this property because

$$\widetilde{d}_{L} = \frac{1}{m} \sum_{\ell=1}^{m} d_{L}(\overline{\psi}_{0}, \overline{\psi} \mid X^{\ell}) + \frac{1}{m} \sum_{\ell=1}^{m} d_{L}(g^{-1}(\overline{\varphi}_{0}), g^{-1}(\overline{\varphi}) \mid X^{\ell})$$

for most bijective maps g, where  $\overline{\varphi}_0 = m^{-1} \sum_{\ell=1}^m \widehat{\varphi}_0^{\ell}$  and  $\overline{\varphi} = m^{-1} \sum_{\ell=1}^m \widehat{\varphi}^{\ell}$ . See § 4 how  $\widetilde{D}_{\rm L}$  vary dramatically with parametrizations in finite samples.

Third, the estimator  $\tilde{r}_L$  involves the estimators of  $\psi$  under  $H_0$ , i.e.,  $\hat{\psi}_0^{\ell}$  and  $\overline{\psi}_0$ . When  $H_0$  fails, they may not be consistent for  $\psi$ . As a result,  $\tilde{r}_L$  is no longer consistent for  $r_m$ . A serious consequence is that the power of the test statistic  $\tilde{D}_L$  is not guaranteed to monotonically increase as  $H_1$  moves away from  $H_0$ . Indeed our simulations (see § 3.2) show that under certain parametrizations, the power may nearly vanish for obviously false  $H_0$ .

Fourth, in order to compute  $d_{\rm L}$  in (1.8), users need to have access to the LRT function  $\widetilde{\mathcal{D}}_{\rm L}$ , but, in most software, the function is built as  $\mathcal{D}_{\rm L}$ , where (1.11)

$$\widetilde{\mathscr{D}}_{L}: (X, \psi_{0}, \psi) \mapsto d_{L}(\psi_{0}, \psi \mid X), \quad \mathscr{D}_{L}: X \mapsto d_{L}(\widehat{\psi}_{0}(X), \widehat{\psi}(X) \mid X).$$

Hence users would need to write themselves a subroutine for evaluating  $\widetilde{\mathcal{D}}_{L}$ . This may not be feasible for users because of lack of information or skills.

In short, four problems need to be resolved: (i) lack of non-negativity, (ii) lack of invariance, (iii) lack of consistency and power, and (iv) lack of a computationally feasible algorithm. Problems (i)–(iii) are resolved in § 2 below, where § 2.1 presents an invariant combining rule, which fully resolves (ii). Next, we propose two estimators of  $\mathcal{F}_m$  (or equivalently  $\mathcal{F}$ ) in § 2.2 and § 2.4. We start with a quick ad hoc fix that requires slightly less assumption but only addresses (i), and then construct a test that fully resolves (i) and (iii). Finally, in § 3, we derive a very handy algorithm, which resolves (iv).

## 2. Improved MI Likelihood Ratio Tests.

2.1. An Invariant Combining Rule. To derive a MI LRT that is invariant to re-parametrization, we replace  $\tilde{d}_{\rm L}$  by an asymptotically equivalent version that behaves like a standard LRT statistic. Specifically, let

(2.1) 
$$\overline{L}(\psi) = \frac{1}{m} \sum_{\ell=1}^{m} L^{\ell}(\psi), \quad \text{where} \quad L^{\ell}(\psi) = \log f(X^{\ell} \mid \psi).$$

We emphasize that  $\overline{L}(\psi)$  is not a real log-likelihood (even if we drop the divider m), because it does not properly model the completed datasets:  $\mathbb{X} = \{X^1, \dots, X^m\}$  (e.g., addressing the issue that all  $X^{\ell}$ s share the same  $X_{\text{obs}}$ ). Nevertheless,  $\overline{L}(\psi)$  can be treated as a log-likelihood for computational purposes. In particular, we can maximize it to obtain

$$(2.2) \quad \hat{\psi}_0^* = \hat{\psi}_0^*(\mathbb{X}) = \underset{\psi \in \Psi}{\arg\max} \ \overline{L}(\psi), \qquad \hat{\psi}^* = \hat{\psi}^*(\mathbb{X}) = \underset{\psi \in \Psi}{\arg\max} \ \overline{L}(\psi).$$

The corresponding log-likelihood ratio test statistic is given by

(2.3) 
$$\widehat{d}_{\mathcal{L}} = 2\left\{\overline{L}(\widehat{\psi}^*) - \overline{L}(\widehat{\psi}_0^*)\right\} = \frac{1}{m} \sum_{\ell=1}^m d_{\mathcal{L}}(\widehat{\psi}_0^*, \widehat{\psi}^* \mid X^{\ell}).$$

Thus, in contrast to  $\widetilde{d}_{\rm L}$  of (1.8),  $\widehat{d}_{\rm L}$  aggregates MI datasets through averaging MI LRT functions as in (2.1), rather than averaging MI test statistics and moments, as in (1.7). Although  $\sqrt{n}(\widehat{\psi}_0^* - \overline{\psi}_0) \stackrel{\rm pr}{\to} \mathbf{0}$  and  $\sqrt{n}(\widehat{\psi}^* - \overline{\psi}) \stackrel{\rm pr}{\to} \mathbf{0}$  as  $n \to \infty$  for each m, only  $\widehat{d}_{\rm L}$ , not  $\widetilde{d}_{\rm L}$ , is guaranteed to be non-negative and invariant to parametrization of  $\psi$  for all m, n. Indeed, the likelihood principle guides us to consider averaging individual log-likelihoods than individual MLEs, since the former has a much better chance to capture functional features of the real log-likelihood than any of their (local) maximizers can.

To derive properties of  $\hat{d}_{L}$ , we need the usual RCs on MLE and MI.

Assumption 2. The sampling model  $f(X \mid \psi)$  satisfies the following:

- (a)  $\psi \mapsto \underline{L}(\psi) = n^{-1} \log f(X \mid \psi)$  is twice continuously differentiable;
- (b) the complete-data MLE  $\hat{\psi}(X)$  is the unique solution of  $\partial \underline{L}(\psi)/\partial \psi = \mathbf{0}$ ;
- (c) if  $\underline{I}(\psi) = -\partial^2 \underline{L}(\psi)/\partial \psi \partial \psi^{\mathsf{T}}$ , then for each  $\psi$ , there exists a positive definite matrix  $\underline{\mathcal{I}}(\psi) = \mathcal{U}_{\psi}^{-1}$  such that  $\underline{I}(\psi) \xrightarrow{\mathrm{pr}} \underline{\mathcal{I}}(\psi)$  as  $n \to \infty$ ; and
- (d) the observed-data MLE  $\hat{\psi}_{obs}$  of  $\psi$  obeys

(2.4) 
$$\left[ \mathcal{T}_{\psi}^{-1/2} \left( \widehat{\psi}_{\text{obs}} - \psi^{\star} \right) \mid \psi^{\star} \right] \Rightarrow \mathcal{N}_{h}(\mathbf{0}, I_{h})$$

as  $n \to \infty$ , where  $I_h$  is the  $h \times h$  identity matrix.

Assumption 3. The imputation model is proper (Rubin, 2004):

(2.5) 
$$\left[ \mathscr{B}_{\psi}^{-1/2} \left( \hat{\psi}^{\ell} - \hat{\psi}_{\text{obs}} \right) \mid X_{\text{obs}} \right] \Rightarrow \mathscr{N}_{h}(\mathbf{0}, I_{h}),$$

$$(2.6) \quad \left[ \mathcal{T}_{\psi}^{-1} \left( U_{\psi}^{\ell} - \mathcal{U}_{\psi} \right) \mid X_{\text{obs}} \right] \stackrel{\text{pr}}{\to} \mathbf{0}, \qquad \left[ \mathcal{T}_{\psi}^{-1} \left( B_{\psi} - \mathcal{B}_{\psi} \right) \mid X_{\text{obs}} \right] \stackrel{\text{pr}}{\to} \mathbf{0}$$

independently for  $\ell = 1, ..., m$ , as  $n \to \infty$ , provided that  $\mathscr{B}_{\psi}^{-1}$  is well-defined.

Assumption 2 holds under the usual RCs that guarantee normality and consistency of MLEs. When the imputations  $X_{\rm mis}^1,\ldots,X_{\rm mis}^m$  are drawn independently from (correctly specified) posterior predictive distribution  $f(X_{\rm mis} | X_{\rm obs})$ , Assumption 3 is typically satisfied. Clearly, we can replace  $\psi$  by its sub-vector  $\theta$  in Assumptions 2 and 3. These  $\theta$ -version assumptions are sufficient to guarantee the validity of the following Theorem 2.4 and Corollary 2.3. For simplicity, Assumption 1, the  $\theta$ -version of Assumptions 2 and 3, and conditions that are strong enough to guarantee Property 1.1 are collectively written as  $RC_{\theta}$ , which are commonly assumed for MI inference.

THEOREM 2.1. Assume  $RC_{\theta}$ . Under  $H_0$ , we have (i)  $\hat{d}_L \ge 0$  for all m, n; (ii)  $\hat{d}_L$  is invariant to parametrization of  $\psi$  for all m, n; and (iii)  $\hat{d}_L = \tilde{d}_L$  as  $n \to \infty$  for each m.

Consequently, an improved combining rule is defined as

(2.7) 
$$\widehat{D}_{L}(\boldsymbol{r}_{m}) = \frac{\widehat{d}_{L}}{k(1+\boldsymbol{r}_{m})},$$

for a given value of  $r_m$ . It follows the forms of (1.4) and of (1.9). The issue is then how to estimate  $r_m$  that avoids the defects of  $\tilde{r}_L$  of (1.8).

2.2. An Improved Estimator of  $r_m$ . Using  $\hat{d}_L$  in (2.3), we can modify  $\tilde{r}_L$  in (1.8) to a potentially better estimator:

(2.8) 
$$\hat{r}_{L} = \frac{m+1}{k(m-1)} (\bar{d}_{L} - \hat{d}_{L}).$$

Although  $\hat{d}_L \ge 0$  is guaranteed by our construction,  $\hat{r}_L \ge 0$  does not hold in general for a finite m. However, it is guaranteed in the following situation.

PROPOSITION 2.2. Write  $\psi = (\theta^{\intercal}, \eta^{\intercal})^{\intercal}$ , where  $\eta$  represents a nuisance parameter that is distinct from  $\theta$ . If there exist functions  $L_{\uparrow}$  and  $L_{\downarrow}$  such that, for all X, the log-likelihood function  $L(\psi \mid X) = \log f(X \mid \psi)$  is of the form  $L(\psi \mid X) = L_{\uparrow}(\theta \mid X) + L_{\downarrow}(\eta \mid X)$ , then  $\hat{r}_{L} \geq 0$  for all m, n.

The condition in Proposition 2.2 means that the likelihood function of  $\psi$  is separable. Then, the profile likelihood estimator of  $\eta$  given  $\theta$ , i.e.,  $\hat{\eta}_{\theta} = \arg\max_{\eta} L(\theta, \eta \mid X)$ , does not depend on  $\theta$ . Trivially, if there is no nuisance parameter  $\eta$ , the separation condition is satisfied. More generally, we have

COROLLARY 2.3. Assume  $RC_{\theta}$ . We have (i) under  $H_0$ ,  $\hat{r}_L \stackrel{pr}{\to} r$  as  $m, n \to \infty$ ; and (ii) under  $H_1$ ,  $\hat{r}_L \stackrel{pr}{\to} r_0$  as  $m, n \to \infty$ , where  $r_0 \ge 0$  is some finite value depending on  $\theta_0$  and  $\theta^*$ .

Corollary 2.3 ensures that, under  $H_0$ ,  $\hat{r}_L$  is non-negative asymptotically and converges in probability to the true  $\mathcal{F}$ . But it also reveals another fundamental defect of  $\hat{r}_L$ : under  $H_1$ , the limit  $\mathcal{F}_0$  may not equal to  $\mathcal{F}$ , a problem we will address in § 2.2. Fortunately, since  $\hat{d}_L \stackrel{\text{pr}}{\to} \infty$  under  $H_1$ , the LRT statistic  $\hat{D}_L(\hat{r}_L)$  is still powerful, albeit the power may be somewhat reduced. Similarly,  $\tilde{r}_L$  of (1.8) has the same asymptotic properties and defects, but  $\hat{r}_L$  behaves more nicely than  $\tilde{r}_L$  for finite m. This hinges closely on the high sensitivity of  $\tilde{r}_L$  to the parametrization of  $\psi$  for small m, e.g., in some cases,  $\tilde{r}_L$  becomes more negative as  $H_1$  moves away from  $H_1$ ; see § 4.1.

Whereas we can fix the occasional negativeness of  $\hat{r}_L$  by using  $\hat{r}_L^+ = \max(0, \hat{r}_L)$ , such an ad hoc fix misses the opportunity to improve upon  $\hat{r}_L$ , and indeed it cannot fix the inconsistency of  $\hat{r}_L$  under  $H_1$ .

2.3. A Complication Caused by Nuisance Parameter. To better understand the source of the negativity of  $\hat{r}_{\rm L}$ , we extend  $\overline{L}(\psi)$  in (2.1) to allow it take m different arguments:

(2.9) 
$$\overline{L}(\psi^1, \dots, \psi^m) = \frac{1}{m} \sum_{\ell=1}^m L^{\ell}(\psi^{\ell}).$$

Table 1 The definitions of hypotheses  $H_0^0$ ,  $H_0^1$ ,  $H_1^0$ ,  $H_1^1$ .

	$\mathscr{C}^0: \psi^1 = \dots = \psi^m \in \Psi$ (i.e., $\mathscr{V} = 0$ )	$\mathscr{C}^1:\psi^1,\ldots,\psi^m\in\Psi$
	(i.e., $\mathcal{F} = 0$ )	(i.e., $r \geq 0$ )
$\mathscr{C}_0: \theta^1 = \dots = \theta^m = \theta_0 \in \Theta$ (i.e., $H_0$ -constrained)	$H_0^0 = \mathscr{C}_0 \cap \mathscr{C}^0$	$H^1_0=\mathscr{C}_0\cap\mathscr{C}^1$
$\mathscr{C}_1: \theta^1, \dots, \theta^m \in \Theta$ (i.e., not $H_0$ -constrained)	$H_1^0 = \mathscr{C}_1 \cap \mathscr{C}^0$	$H_1^1=\mathscr{C}_1\cap\mathscr{C}^1$



FIG 1. The relationships between the four hypotheses  $H_0^0$ ,  $H_0^1$ ,  $H_1^0$ ,  $H_1^1$ . Each arrow denotes an implication, e.g.,  $H_0^0 \Rightarrow H_0^1$  means that  $H_0^0$  implies  $H_0^1$ .

Using the "log-likelihood"  $\overline{L}(\psi^1, \ldots, \psi^m)$ , we can construct, at least conceptually, four hypotheses  $H_0^0$ ,  $H_0^1$ ,  $H_1^0$ ,  $H_1^1$  defined in Table 1. Each of them consists of zero, one or two of the constraints

$$\mathscr{C}_0: \theta^1 = \cdots = \theta^m = \theta_0$$
 and  $\mathscr{C}^0: \psi^1 = \cdots = \psi^m$ .

The constraint  $\mathscr{C}_0$  is equivalent to  $H_0$ , and the constraint  $\mathscr{C}^0$  means that all  $\psi^{\ell}$ s are equal, and hence it is effectively equivalent to r = 0, i.e., no missing information. The relationships among  $H_0^0$ ,  $H_0^1$ ,  $H_1^0$ ,  $H_1^1$  can be visualized in Figure 1. Define the maximized value of  $\overline{L}(\psi^1, \ldots, \psi^m)$  under hypothesis  $H \in \{H_0^0, H_0^1, H_1^0, H_1^1\}$  by  $\mathbb{L}(H)$ . Then we can re-express  $(\overline{d}_L - \widehat{d}_L)/2$  as

(2.10) 
$$(\overline{d}_{L} - \widehat{d}_{L})/2 = \{ \mathbb{L}(H_{1}^{1}) - \mathbb{L}(H_{1}^{0}) \} - \{ \mathbb{L}(H_{0}^{1}) - \mathbb{L}(H_{0}^{0}) \}.$$

Whereas the two bracketed terms in (2.10) are non-negative because they correspond to two LRT statistics, the difference between these two terms is not guaranteed to be non-negative. A simple example illustrates this well. For the regression model  $[Y \mid X_1, X_2] \sim \mathcal{N}(\beta_0 + \beta_1 X_1 + \beta_2 X_2, \sigma^2)$ , the LRT statistic for testing  $H_1^0: \beta_1 = 0, \beta_2 \in \mathbb{R}$  against  $H_1^1: \beta_1, \beta_2 \in \mathbb{R}$  is not necessarily larger (or smaller) than that for testing  $H_0^0: \beta_1 = \beta_2 = 0$  against  $H_0^1: \beta_1 \in \mathbb{R}, \beta_2 = 0$ . A schematic illustration is provided in Figure 2.

The decomposition (2.10) provides another interpretation of  $\hat{r}_L$ . The test statistic  $\mathbb{L}(H_1^1) - \mathbb{L}(H_1^0)$  seeks evidence for detecting the falsity of r = 0 in

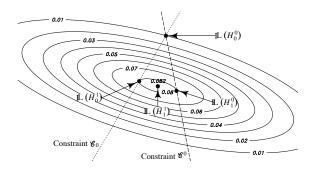


FIG 2. A schematic illustration of the sign of (2.10). The contour lines of  $\overline{L}(\psi^1,\ldots,\psi^m)$  are plotted. The two straight lines refer to constraints  $\mathscr{C}_0$  and  $\mathscr{C}^0$ . Since  $\mathbb{L}(H_1^1) = 0.082$ ,  $\mathbb{L}(H_0^1) = \mathbb{L}(H_1^0) = 0.08$ , and  $\mathbb{L}(H_0^0) = 0.01$ , we have  $\{\mathbb{L}(H_1^1) - \mathbb{L}(H_1^0)\} - \{\mathbb{L}(H_0^1) - \mathbb{L}(H_0^0)\} = 0.002 - 0.007 < 0$ . Note that the function  $\overline{L}(\psi^1,\ldots,\psi^m)$  in (2.9) is at least 4-dimensional (i.e.,  $\theta^1, \theta^2, \eta^1, \eta^2$ ) generally, so the above illustration in a 2-dimension space is just conceptual.

both  $\theta$  and  $\eta$ , whereas  $\mathbb{L}(H_0^1) - \mathbb{L}(H_0^0)$  seeks evidence only in  $\eta$ . For cases where  $\theta$  and  $\eta$  are orthogonal (at least locally), the left-hand side of (2.10) can be viewed as a measure of evidence against r = 0 solely from  $\theta$ ; Proposition 2.2 already hinted this possibility. However, the "test statistic" (2.10) has a very serious problem apart from being possibly negative. Because  $\mathcal{C}_0$  requires all  $\theta^{\ell}$ s to coincide with a specific  $\theta_0$ ,  $\mathcal{C}_0$  is not nested within  $\mathcal{C}^0$ , i.e.,  $\mathcal{C}^0 \Rightarrow \mathcal{C}_0$ . Hence  $\hat{r}_L$  is guaranteed to consistently estimate  $r_m$  only under  $H_0$ . This explains Corollary 2.3, and leads to an improvement below.

2.4. A Consistent and Non-negative Estimator of  $r_m$ . Our new estimator simply drops the second term in (2.10), that is, we estimate  $r_m$  by

(2.11) 
$$\widehat{r}_{L}^{\Diamond} = \frac{m+1}{h(m-1)} (\overline{\delta}_{L} - \widehat{\delta}_{L}), \text{ where}$$

(2.12) 
$$\overline{\delta}_{L} = 2\overline{L}(\hat{\psi}^{1}, \dots, \hat{\psi}^{m}), \qquad \hat{\delta}_{L} = 2\overline{L}(\hat{\psi}^{*}, \dots, \hat{\psi}^{*}),$$

where h is the dimension of  $\psi$ , and the rhombus " $\diamondsuit$ " symbolizes a robust estimator. It is robust, because it is consistent under either  $H_0$  or  $H_1$ , as long as we are willing to impose the EFMI assumption on the entire parameter  $\psi$ , a stronger requirement than Assumption 1. This expansion from  $\theta$  to  $\psi$  is inevitable because the LRT must handle the entire  $\psi$ , not just  $\theta$ . The collection of Assumptions 2–4 will be referred to as  $RC_{\psi}$ .

Assumption 4 (EFMI of  $\psi$ ). There is  $r \ge 0$  such that  $\mathcal{T}_{\psi} = (1+r)\mathcal{U}_{\psi}$ .

THEOREM 2.4. Assume  $RC_{\psi}$ . Then for any value of  $\psi$ , we have (i)  $\hat{r}_L^{\Diamond} \geqslant 0$  for all m, n; (ii)  $\hat{r}_L^{\Diamond}$  is invariant to parametrization of  $\psi$  for all m, n; and (iii)  $\hat{r}_L^{\Diamond} \stackrel{pr}{\to} r$  as  $m, n \to \infty$ , where r is given in Assumption 4.

With the improved combining rule  $\hat{D}_{L}(r_{m})$  of (2.7) and improved estimators for  $r_{m}$ , we are ready to propose two MI LRT statistics:

(2.13) 
$$\hat{D}_{L}^{+} = \hat{D}_{L}(\hat{r}_{L}^{+})$$
 and  $\hat{D}_{L}^{\Diamond} = \hat{D}_{L}(\hat{r}_{L}^{\Diamond}).$ 

For comparison, we also study the test statistic  $\hat{D}_{L} = \hat{D}_{L}(\hat{r}_{L})$ .

2.5. Reference Null Distributions. The estimators  $\hat{r}_{\rm L}^+$  and  $\tilde{r}_{\rm L}$  have the same functional form asymptotically  $(n \to \infty)$  and rely on the same set of assumptions, hence they have the same asymptotic distribution.

LEMMA 2.5. Suppose  $RC_{\theta}$  and m > 1. Under  $H_0$ , we have, jointly,

$$(2.14) \qquad \frac{\hat{r}_{\rm L}^+}{r_m} \Rightarrow M_2 \quad and \quad \hat{D}_{\rm L}^+ \Rightarrow \frac{(1+r_m)\,M_1}{1+r_mM_2}$$

as  $n \to \infty$ , where  $M_1 \sim \chi_k^2/k$  and  $M_2 \sim \chi_{k(m-1)}^2/\{k(m-1)\}$  are independent.

Consequently, the null distribution of  $\hat{D}_{\rm L}^+ = \hat{D}_{\rm L}(\hat{r}_{\rm L}^+)$  can be approximated by  $F_{k, {\rm df}(\hat{r}_{\rm L}^+, k)}$ , but a better approximation will be provided shortly.

For the other proposal, although  $\hat{r}_{L}^{+} - \hat{r}_{L}^{\Diamond} \xrightarrow{\text{pr}} 0$  as  $n \to \infty$  under  $H_0$ , their non-degenerated distributions (after proper scaling) are different because  $\hat{r}_{L}^{\Diamond}$  relies on an average FMI in  $\psi$ , but  $\hat{r}_{L}^{+}$  only on an average FMI in  $\theta$ .

Theorem 2.6. Suppose  $RC_{\psi}$  and m > 1. Then for any value of  $\psi$ ,

(2.15) 
$$\frac{\hat{r}_{\rm L}^{\Diamond}}{r_m} \Rightarrow M_3 \sim \frac{\chi_{h(m-1)}^2}{h(m-1)}$$

as  $n \to \infty$ , where  $M_3$  is independent of the  $M_1$  defined in (2.14).

Theorem 2.6 implies that, if n can be regarded as infinity and  $\hat{r}_{\rm L}^{\Diamond}$  is uniformly integrable in  $\mathcal{L}^2$ , then

$$\mathsf{Bias}(\widehat{r}_{\mathsf{L}}^{\Diamond}) = \mathsf{E}(\widehat{r}_{\mathsf{L}}^{\Diamond}) - \mathscr{V}_m = 0 \qquad \text{and} \qquad \mathsf{Var}(\widehat{r}_{\mathsf{L}}^{\Diamond}) = \frac{2\mathscr{V}_m^2}{h(m-1)} = O(m^{-1})$$

as  $m \to \infty$ . Therefore  $\hat{r}_{\rm L}^{\lozenge}$  is a  $\sqrt{m}$ -consistent estimator of  $\nu$  in  $\mathscr{L}^2$ . Moreover, for each m > 1 and as  $n \to \infty$ , we have

$$\frac{\mathsf{Bias}(\widehat{r}_{\mathtt{L}}^+)}{\mathsf{Bias}(\widehat{r}_{\mathtt{L}}^{\lozenge})} \to 1 \qquad \text{and} \qquad \frac{\mathsf{Var}(\widehat{r}_{\mathtt{L}}^+)}{\mathsf{Var}(\widehat{r}_{\mathtt{L}}^{\lozenge})} \to \frac{h}{k} \geqslant 1,$$

which implies that  $\hat{r}_{\rm L}^{\Diamond}$  is no less efficient than  $\hat{r}_{\rm L}^{+}$  when  ${\rm RC}_{\psi}$  holds. This is of no surprise because of the extra information in  $\hat{r}_{\rm L}^{\Diamond}$  brought in by the stronger Assumption 4. Result (2.15) also gives us the exact (i.e., for any m>1, but assuming  $n\to\infty$ ) reference null distribution of  $\hat{D}_{\rm L}^{\Diamond}$ , as given below.

Theorem 2.7. Assume  $RC_{\psi}$  and m > 1. Under  $H_0$ , we have

(2.16) 
$$\widehat{D}_{L}^{\Diamond} \Rightarrow \frac{(1 + r_{m}) M_{1}}{1 + r_{m} M_{3}} \equiv D$$

as  $n \to \infty$ , where  $M_1 \sim \chi_k^2/k$  and  $M_3 \sim \chi_{h(m-1)}^2/\{h(m-1)\}$  are independent.

The impact of the nuisance parameter  $r_m$  on the null distribution diminishes with m, because  $\hat{D}_{\rm L}^{\Diamond}$  and  $\hat{D}_{\rm L}^{+}$  converge in distribution to  $M_1 = \chi_k^2/k$  as  $m, n \to \infty$ . Since  $M_3 \stackrel{\rm pr}{\to} 1$  faster than  $M_2 \stackrel{\rm pr}{\to} 1$ ,  $\hat{D}_{\rm L}^{\Diamond}$  is expected to be more robust to  $r_m$ . Nevertheless, because m typically is small in practice (e.g.,  $m \le 10$ ), we cannot ignore the impact of  $r_m$ . This issue has been largely dealt with in the literature by seeking an  $F_{k,\rm df}$  distribution as an approximate null distribution, as in Li et al. (1991b). However, directly adopting their df of (1.6) leads to poorer approximation for our purposes; see below. A better approximation is to match the first two moments of the denominator of (2.16),  $1 + r_m M_3$ , with that of a scaled  $\chi^2$ :  $a\chi_b^2/b$ . This yields  $a = 1 + r_m$  and  $b = (1 + r_m^{-1})^2 h(m-1)$ , and the approximated  $F_{k, {\rm df}(r_m, h)}$ , where

(2.17) 
$$\widehat{\mathrm{df}}(r_m, h) = \left\{ \frac{1 + r_m}{r_m} \right\}^2 h(m - 1) = \frac{h(m - 1)}{f_m^2}.$$

This degrees of freedom is appealing because it simply inflates the denominator degrees of freedom  $M_3$  by dividing it by  $\mathcal{F}_m^2$ , where  $\mathcal{F}_m = \mathcal{F}_m/(1+\mathcal{F}_m)$  the finite imputation corrected FMI. Intuitively, the less missing information, the closer  $F_{k,\hat{\text{df}}(\mathcal{F}_m,h)}$  should be to  $\chi_k^2/k$ , the usual large-n asymptotic  $\chi^2$  test; as mentioned earlier, for small n, see Barnard and Rubin (1999).

To compare  $F_{k,\widehat{\mathrm{df}}(r_m,h)}$  with  $F_{k,\widetilde{\mathrm{df}}(r_m,h)}$  as approximations to the distribution of D given in (2.16), we compute via simulations

$$\widetilde{\alpha} = \mathsf{P}\left\{D > F_{k, \widehat{\mathrm{df}}(r_m, h)}^{-1}(1 - \alpha)\right\} \quad \text{and} \quad \widehat{\alpha} = \mathsf{P}\left\{D > F_{k, \widehat{\mathrm{df}}(r_m, h)}^{-1}(1 - \alpha)\right\},$$

where  $F_{k,\mathrm{df}}^{-1}(q)$  denotes the q-quantile of  $F_{k,\mathrm{df}}$ . We draw  $N=2^{18}$  independent dent copies D for each of the following possible combinations:  $m \in \{3, 5, 7\}$ ,  $k \in \{1, 2, 4, 8\}, \ \tau = h/k \in \{1, 2, 3\}, \ \not f_m \in \{0, 0.1, \dots, 0.9\}, \ \text{and following Ben-}$ jamin et al. (2018)'s recommendation, we use both  $\alpha \in \{0.5\%, 5\%\}$ . The results for  $\alpha = 0.5\%$  and for  $\alpha = 5\%$  are shown respectively in Figure 3 and in the Appendix. In general,  $\hat{\alpha}$  approximates  $\alpha$  much better than  $\tilde{\alpha}$ , especially when m, k, h are small. When m, h are larger, their performances are similar because both  $F_{k,\widetilde{\mathrm{df}}(r_m,h)}$  and  $F_{k,\widehat{\mathrm{df}}(r_m,h)}$  get close to  $\chi_k^2/k$ . But the performances of  $\widetilde{\alpha}$  and  $\widehat{\alpha}$  are not monotonic in  $\mathscr{S}_m$ . Generally speaking, the performance of  $F_{k,\hat{\mathrm{df}}(r_m,h)}$  is particularly good for  $0\% \lesssim f_m \lesssim 30\%$ . Consequently, we recommend using  $F_{k,\hat{\mathrm{df}}(\hat{r}_L^\lozenge,h)}$  as an approximate null distribution for  $\widehat{D}_{L}^{\Diamond}$ , and  $F_{k,\widehat{\mathrm{df}}(\widehat{r}_{L}^{+},k)}$  for  $\widehat{D}_{L}^{+}$ , as employed in the rest of this paper. However, these approximations obviously suffer from the usual "plug-in problem" by ignoring the uncertainties in estimating  $r_m$ . Since the  $F_{k,df}$  is not too sensitive to the value of df once it is reasonably large (df  $\geq 20$ ), the "plug-in problem" is less an issue here than in many other context, leading to acceptable approximations as empirically demonstrated in Section 4. Nevertheless, further improvements are likely and should be sought.

**3.** Computational Considerations and Comparisons. The statistic  $\overline{d}_L$  of (1.7) is easy to compute because only the standard complete-data procedure  $\mathscr{D}_L: X \mapsto d_L(\widehat{\psi}_0(X), \widehat{\psi}(X) \mid X)$  is needed. However,  $\widehat{d}_L$  of (2.3) and  $\widehat{r}_L^{\Diamond}$  of (2.8) in general cannot be computed solely by  $\mathscr{D}_L$ , e.g.,  $\widehat{d}_L$  requires

$$\overline{\mathscr{D}}_{\mathrm{L}}: \mathbb{X} \mapsto \frac{1}{m} \sum_{\ell=1}^{m} d_{\mathrm{L}}(\widehat{\psi}_{0}^{*}(\mathbb{X}), \widehat{\psi}^{*}(\mathbb{X}) \mid X^{\ell}).$$

Creating a subroutine for this computation requires additional effort and information that may beyond a user's capacity. Here we show how to compute or approximate  $\hat{d}_L$  and  $\hat{r}_L^{\Diamond}$  solely by  $\mathcal{D}_L$  or a trivial modification of  $\mathcal{D}_L$ .

3.1. Computationally Feasible Combining Rule. We begin with precise notation for our complete data X and its sampling model  $f(X|\psi)$ . For the vast majority of real-world datasets, X is of the form of an  $n \times p$  matrix, with rows indicating subjects and columns variables/attributes. We then write  $X = (X_1, \ldots, X_n)^{\mathsf{T}}$ , and its sampling model by  $f_n(X \mid \psi)$ . Correspondingly, the  $\ell$ th completed-dataset by MI is  $X^{\ell} = (X_1^{\ell}, \ldots, X_n^{\ell})^{\mathsf{T}}$ . Define the stacked dataset by  $X^{\mathsf{S}} = [(X^1)^{\mathsf{T}}, \ldots, (X^m)^{\mathsf{T}}]^{\mathsf{T}}$ , a matrix having mn rows, which is conceptually different from the collection of datasets  $\mathbb{X} = \{X^1, \ldots, X^m\}$ .

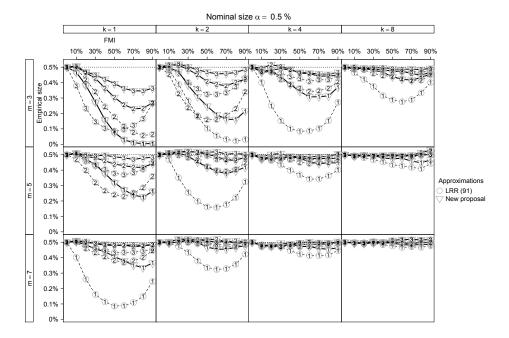


FIG 3. The performances of two different approximated null distributions when the nominal size is  $\alpha=0.5\%$ . The vertical axis denotes  $\widehat{\alpha}$  or  $\widetilde{\alpha}$ , and the horizontal axis denotes the value of  $f_m$ . The number attached to each line denotes the value of  $\tau=h/k$ . The proposed approximation  $\widehat{\alpha}$  is denoted by thick solid lines with triangles, and the existing approximation  $\widetilde{\alpha}$  is denoted by thin dashed lines with circles.

Treating  $X^{S}$  as a dataset with size mn, we can define

(3.1) 
$$\overline{L}^{S}(\psi) = \frac{1}{m} \log f_{mn}(X^{S} \mid \psi),$$

which, other than the scaling factor 1/m, is just the ordinary log-likelihood function of  $\psi$  based on the dataset  $X^{\rm S}$  (for computation purposes). Consequently, as long as the user's complete-data procedure can handle size mn instead of just n, the user can apply it to  $X^{\rm S}$  to obtain

(3.2) 
$$\hat{\psi}_0^{\mathrm{S}} = \underset{\psi \in \Psi}{\operatorname{arg\,max}} \, \overline{L}^{\mathrm{S}}(\psi) \quad \text{and} \quad \hat{\psi}^{\mathrm{S}} = \underset{\psi \in \Psi}{\operatorname{arg\,max}} \, \overline{L}^{\mathrm{S}}(\psi).$$

Consequently, the quantities

(3.3) 
$$\hat{\delta}_{0,S} = 2\overline{L}^{S}(\hat{\psi}_{0}^{S})$$
 and  $\hat{\delta}_{S} = 2\overline{L}^{S}(\hat{\psi}^{S})$ 

are readily available from the user's complete-data procedure. It is then desirable if we can replace  $\overline{L}(\psi)$  by  $\overline{L}^{S}(\psi)$  in the proposed test statistics.

Precisely, in parallel to (2.7), (2.8) and (2.11), we define

(3.4) 
$$\widehat{D}_{S}(r_{m}) = \frac{\widehat{d}_{S}}{k(1+r_{m})},$$
 with  $\widehat{d}_{S} = \widehat{\delta}_{S} - \widehat{\delta}_{0,S}$  of (3.3);

(3.5) 
$$\hat{r}_{S} = \frac{m+1}{k(m-1)} (\bar{d}_{S} - \hat{d}_{S}), \text{ with } \bar{d}_{S} = \bar{d}_{L} \text{ of } (1.7);$$

(3.6) 
$$\widehat{r}_{S}^{\Diamond} = \frac{m+1}{h(m-1)} (\overline{\delta}_{S} - \widehat{\delta}_{S}), \quad \text{with } \overline{\delta}_{S} = \overline{\delta}_{L} \text{ of } (2.12);$$

and  $\hat{r}_{\rm S}^+ = \max(0,\hat{r}_{\rm S})$ . The "stacked" counterparts of  $\hat{D}_{\rm L}^{\Diamond}$  and its existing counterparts  $\hat{D}_{\rm L}$  and  $\hat{D}_{\rm L}^+$  (see (2.13)) then are given by

$$(3.7) \qquad \hat{D}_{S}^{\Diamond} = \hat{D}_{S}(\hat{r}_{S}^{\Diamond}), \qquad \hat{D}_{S} = \hat{D}_{S}(\hat{r}_{S}), \qquad \hat{D}_{S}^{+} = \hat{D}_{S}(\hat{r}_{S}^{+}).$$

PROPOSITION 3.1. If  $X = (X_1, \dots, X_n)^{\mathsf{T}}$  is row-independent for arbitrary n, i.e.,  $f(X \mid \psi) = \prod_{i=1}^n f(X_i \mid \psi)$ , then (2.1) and (3.1) are the same:  $\overline{L}(\psi) \equiv \overline{L}^{\mathsf{S}}(\psi)$ . Consequently,  $\hat{D}_{\mathsf{S}} \equiv \hat{D}_{\mathsf{L}}$  and  $\hat{D}_{\mathsf{S}}^{\Diamond} \equiv \hat{D}_{\mathsf{L}}^{\Diamond}$ .

Since for many applications, the rows correspond to individual subjects, the row-independence assumption typically holds for arbitrary n. Hence we can extend from n to mn, assuming the user's complete-data procedure is not size-limited. Even if it does not hold, we can still have  $\hat{d}_L \simeq \hat{d}_S$  under some RCs that guarantee  $\overline{L}(\psi)$  and  $\overline{L}^S(\psi)$  are close; see Appendix A, where we also reveal a subtle but important difference between  $\hat{r}_L^{\Diamond}$  and  $\hat{r}_S^{\Diamond}$ .

Similar to  $\mathcal{D}_L$  in (1.11), we define complete-data functions (with data X being the only input)

(3.8) 
$$\mathscr{D}_{L,0}(X) = 2\log f(X \mid \hat{\psi}_0(X)), \quad \mathscr{D}_{L,1}(X) = 2\log f(X \mid \hat{\psi}(X)).$$

The subroutine for evaluating the complete-data LRT function  $X \mapsto \mathcal{D}_{L}(X)$  is usually available, as is the subroutine for  $X \mapsto \mathcal{D}_{L,1}(X)$ ; for example, the function logLik in R extracts the maximum of complete data log-likelihood for objects belonging to classes "glm", "lm", "nls" and "Arima".

The Algorithms 1 and 2 listed below compute tests given  $\widehat{D}_{S}^{\Diamond}$  and  $\widehat{D}_{S}^{+}$ , respectively. Whenever possible, we recommend the use of the robust MI LRT given by Algorithm 1, since it has the best theoretical guarantee. The second test can be useful when  $\mathscr{D}_{L}$  but not  $\mathscr{D}_{L,1}$  is available.

3.2. Computational Comparison with Existing Tests. First, we list some existing estimators of  $r_m$  and their computation. Let  $s_{\mathrm{W},1}^2$  and  $s_{\mathrm{W},1/2}^2$  be the sample variances of  $\{d_{\mathrm{W}}^\ell\}_{\ell=1}^m$  and  $\{\sqrt{d_{\mathrm{W}}^\ell}\}_{\ell=1}^m$ , respectively. By the method of

# **Algorithm 1:** (Robust) MI LRT statistic $\hat{D}_{S}^{\Diamond}$

```
Input: Datasets X^1, \ldots, X^m; dimensions h, k; functions \mathcal{D}_{L,1}, \mathcal{D}_L in (3.8), (1.11). begin

Compute \overline{\delta}_S = m^{-1} \{ \mathcal{D}_{L,1}(X^1) + \cdots + \mathcal{D}_{L,1}(X^m) \}.

(i) Stack the datasets to form X^S = [(X^1)^\intercal, \ldots, (X^m)^\intercal]^\intercal.

(ii) Compute \hat{d}_S = m^{-1} \mathcal{D}_L(X^S) and \hat{\delta}_S = m^{-1} \mathcal{D}_{L,1}(X^S).

Calculate \hat{r}_S^{\Diamond} according to (3.6), and \hat{D}_S^{\Diamond} according to (3.4) and (3.7).

Calculate \hat{d}f(\hat{r}_S^{\Diamond}, h) according to (2.17).

Compute the p-value as 1 - F_{k,\widehat{d}f}(\hat{r}_S^{\Diamond}, h) (\widehat{D}_S^{\Diamond}).
```

## **Algorithm 2:** MI LRT statistic $\hat{D}_{S}^{+}$

```
Input: Datasets X^1, \ldots, X^m; dimension k; function \mathcal{D}_L in (1.11). begin

Compute \overline{d}_L = m^{-1} \{ \mathcal{D}_L(X^1) + \cdots + \mathcal{D}_L(X^m) \}.

(i) Stack the datasets to form X^S = [(X^1)^\intercal, \ldots, (X^m)^\intercal]^\intercal.

(ii) Compute \widehat{d}_S = m^{-1} \mathcal{D}_L(X^S).

Calculate \widehat{r}_S^+ according to (3.5), and \widehat{D}_S^+ according to (3.4) and (3.7).

Calculate \widehat{df}(\widehat{r}_S^+, k) according to (2.17).

Compute the p-value as 1 - F_{k,\widehat{df}(\widehat{r}_S^+, k)(\widehat{D}_S^+)}.
```

moments concerning  $s_{\mathrm{W},1}^2$  and  $s_{\mathrm{W},1/2}^2$ , Rubin (2004) and Li *et al.* (1991a) respectively proposed estimating  $r_m$  by

(3.9) 
$$\widetilde{r}_{W,1} = \frac{(1+1/m)s_{W,1}^2}{2\overline{d}_W + \sqrt{\left\{4\overline{d}_W^2 - 2ks_{W,1}^2\right\}^+}}, \qquad \widetilde{r}_{W,1/2} = (1+1/m)s_{W,1/2}^2,$$

where  $\{a\}^+ = \max(0, a)$ . Note that when k is large and m is small, using (3.9) may lead to power loss, although the users have no choice when they are given only  $\{d_{\mathbf{W}}^{\ell}\}_{\ell=1}^{m}$ . A trivial modification of  $\widetilde{r}_{\mathbf{L}}$  of (1.8), i.e.,  $\widetilde{r}_{\mathbf{L}}^+ = \max(0, \widetilde{r}_{\mathbf{L}})$ , is a better alternative if the user is able to compute  $\widetilde{r}_{\mathbf{L}}$ .

Second, we list some alternative MI combining rules. Having the above estimators of  $r_m$ , we can insert them into the following combining rules:

$$\widetilde{D}'_{\mathrm{W}}(\boldsymbol{r}_{m}) = \frac{\widetilde{d}'_{\mathrm{W}}}{k(1+\boldsymbol{r}_{m})}, \quad \widetilde{D}_{\mathrm{L}}(\boldsymbol{r}_{m}) = \frac{\widetilde{d}_{\mathrm{L}}}{k(1+\boldsymbol{r}_{m})}, \quad \widetilde{D}_{\mathrm{L}}^{+}(\boldsymbol{r}_{m}) = \left\{\widetilde{D}_{\mathrm{L}}(\boldsymbol{r}_{m})\right\}^{+}.$$

Using (1.3) and (1.8), we can also define the following combining rules:

(3.10) 
$$\overline{D}'_{W}(\boldsymbol{r}_{m}) = \frac{\overline{d}'_{W} - \frac{k(m-1)}{m+1}\boldsymbol{r}_{m}}{k(1+\boldsymbol{r}_{m})}, \qquad \overline{D}_{L}(\boldsymbol{r}_{m}) = \frac{\overline{d}_{L} - \frac{k(m-1)}{m+1}\boldsymbol{r}_{m}}{k(1+\boldsymbol{r}_{m})}.$$

The combining rule  $\overline{D}'_W(r_m)$  is useful when computing  $\widetilde{d}'_W$  is difficult but computing  $\overline{d}'_W$  and estimating  $r_m$  are simple. However, the resulting power may deteriorate if the estimator of  $r_m$  is inefficient or inaccurate. This type of test statistics is also mentioned in Li et al. (1991a). Indeed, there are infinitely many asymptotically equivalent test statistics, e.g., any convex combination of  $\widetilde{D}_W(r_m)$  and  $\overline{D}_W(r_m)$ , i.e.,  $\phi \widetilde{D}_W(r_m) + (1-\phi)\overline{D}_W(r_m)$ , for  $\phi \in [0,1]$ . When  $\widetilde{r}_{W,1}$  or  $\widetilde{r}_{W,1/2}$  is used for estimating  $r_m$ , the null distributions of the resulting MI test statistics can be approximated by  $F_{k,\widetilde{\text{df}}'(r_m,k)}$ , where  $\widetilde{\text{df}}'(r_m,k) = (m-1)(1+r_m^{-1})^2k^{-3/m}$ ; see Li et al. (1991a).

Next, we introduce and recall some notation to facilitate the comparison: (a) standard complete-data moments estimation ( $\mathcal{M}_W$  and  $\mathcal{M}_L$ ) and testing ( $\mathcal{D}_W$  and  $\mathcal{D}_L$ ) procedures, and (b) non-standard complete-data procedures ( $\widetilde{\mathcal{D}}_L$ ,  $\overline{\mathcal{D}}_L$ ,  $\mathcal{D}_{L,1}$  and  $\overline{\mathcal{D}}_{L,1}$ ), where

$$\mathcal{M}_{\mathrm{W}}(X) = \left\{ \widehat{\theta}(X), U(X) \right\}, \quad \mathcal{M}_{\mathrm{L}}(X) = \left\{ \widehat{\psi}(X), \widehat{\psi}_{0}(X) \right\},$$

$$\mathcal{D}_{\mathrm{W}}(X) = d_{\mathrm{W}}(\widehat{\theta}(X), U(X)), \quad \overline{\mathcal{D}}_{\mathrm{L},1}(\mathbb{X}) = \frac{2}{m} \sum_{\ell=1}^{m} \log f(X^{\ell} \mid \widehat{\psi}^{*}(\mathbb{X})).$$

Clearly,  $\mathcal{M}_W$  produces  $\mathcal{D}_W$ , and  $\{\mathcal{M}_L, \widetilde{\mathcal{D}}_L\}$  produces  $\mathcal{D}_L$ . If users can perform optimization,  $\widetilde{\mathcal{D}}_L$  produces  $\{\mathcal{M}_L, \mathcal{D}_L, \mathcal{D}_{L,1}\}$ . Note that an un-normalized density can be used in  $\mathcal{D}_{L,1}$  and  $\overline{\mathcal{D}}_{L,1}$ .

Table 2 summarizes whether a particular pair of  $D(\cdot)$  and r, resulting the statistic D(r), has the following statistical or computational properties.

- (Inv) D(r) and r are invariant to re-parametrization of  $\psi$ ;
- (Rob) r is robust against  $\theta_0$ , i.e., consistent under both  $H_0$  and  $H_1$ ;
- $(\geq 0)$  D(r) and r are non-negative for all m and n;
- (Pow) the test has high power to reject  $H_0$  under  $H_1$ ;
- (Def) D(r) and r are always well-defined and numerically well-conditioned;
- (Sca) the MI procedure requires users only to deal with scalars;
- (Dep)  $X_1, \ldots, X_n$  can be dependent; and
- (EFMI) whether EFMI is assumed for  $\theta$  or for  $\psi$ .

In summary,  $\hat{D}_{S}(\hat{r}_{S}^{+})$  is the most computationally attractive test statistic. If the user is willing to make stronger assumptions,  $\hat{D}_{S}(\hat{r}_{S}^{\Diamond})$  has better statistical properties, and is still computationally feasible. Nevertheless,  $\hat{D}_{L}(\hat{r}_{L}^{\Diamond})$  is the most general test statistic and has the best statistical properties.

Computational requirements and statistical properties of MI test statistics, their associated combining rules and estimators of FMI  $r_m$ . The symbols "+" and "-" mean that the test statistic (or estimator) is equipped and not equipped with the indicated property, respectively; see the end of § 3.2 for heading descriptions. The reference papers/book are abbreviated as follows: Rubin (2004) (R04), Li et al. (1991a) (LMRR91) and Meng and Rubin (1992) (MR92).

		Combini	ng Rule	Estimate	or of $r_m$	Approx. null	${\rm distribution}^{\it a}$					Proj	pertie	s		
Test	No.	Formula	Routine	Formula	Routine	Original	Proposed	Reference	Inv	Rob	≥ 0	Pow	Def	Sca	Dep	EFMI
WT	W-1	$\widetilde{D}'_{\mathrm{W}}(r_m)$	$\mathscr{M}_{\mathrm{W}}$	$\widetilde{r}_{\mathrm{W}}'$	$\mathscr{M}_{\mathrm{W}}$	$F_{k,\widetilde{\operatorname{df}}(r_m,k)}^{b}$	$F_{k,\widehat{\mathrm{df}}(r_m,k)}$	R04	_	+ c	+	_	_	_	+	$\theta$
	W-2	$\widetilde{D}_{\mathrm{W}}'(\boldsymbol{r}_m)^{\boldsymbol{d}}$	$\mathscr{M}_{\mathrm{W}}$	$\widetilde{r}'_{\mathrm{W},1}$	$\mathscr{D}_{\mathrm{W}}$	$F_{k,\widetilde{\mathrm{df}}'(r_m,k)}$	NA	R04	_	_	+	_	_	_	+	$\theta$
	W-3	$\widetilde{D}_{\mathrm{W}}'(\boldsymbol{r}_m)$	$\mathscr{M}_{\mathrm{W}}$	$\widetilde{r}'_{\mathrm{W},1/2}$	$\mathscr{D}_{\mathrm{W}}$	$F_{k,\widetilde{\mathrm{df}}'(r_m,k)}$	NA	LMRR91	_	_	+	_	_	_	+	$\theta$
	W-4	$D_{\mathrm{W}}(T)^{e}$	$\mathscr{M}_{\mathrm{W}}$	$\widetilde{r}_{\mathrm{W}}'$	$\mathscr{D}_{\mathrm{W}}$	$F_{k,\widetilde{\mathrm{df}}(r_m,k)}$	$F_{k,\widehat{\mathrm{df}}(r_m,k)}$	R04	_	+	+	_	_	_	+	$ heta^{f}$
	W-5	$\overline{D}'_{\mathrm{W}}(r_m)$	$\mathscr{D}_{\mathrm{W}}$	$\widetilde{r}_{\mathrm{W},1}'$	$\mathscr{D}_{\mathrm{W}}$	$F_{k,\widetilde{\mathrm{df}}'(r_m,k)}$	NA	R04	_	_	_	_	_	+	+	$\theta$
	W-6	$\overline{D}'_{\mathrm{W}}(r_m)$	$\mathscr{D}_{\mathrm{W}}$	$\widetilde{r}'_{\mathrm{W},1/2}$	$\mathscr{D}_{\mathrm{W}}$	$F_{k,\widetilde{\mathrm{df}}'(r_m,k)}$	NA	LMRR91	-	_	-	_	-	+	+	$\theta$
LRT	L-1	$\widetilde{D}_{\mathrm{L}}(\boldsymbol{r}_{m})$	$\mathscr{M}_{\mathrm{L}},\widetilde{\mathscr{D}}_{\mathrm{L}}$	$\widetilde{r}_{ m L}$	$\mathscr{M}_{\mathrm{L}},\widetilde{\mathscr{D}}_{\mathrm{L}}$	$F_{k,\widetilde{\mathrm{df}}(r_m,k)}$	$F_{k,\widehat{\mathrm{df}}(r_m,k)}$	MR92	_	_	_	_	+	g	+	$\theta$
	L-2	$\widetilde{D}_{\mathrm{L}}^{+}(\boldsymbol{r}_{m})$	$\mathscr{M}_{\mathrm{L}},\widetilde{\mathscr{D}}_{\mathrm{L}}$	$\widetilde{r}_{ m L}^+$	$\mathscr{M}_{\mathrm{L}},\widetilde{\mathscr{D}}_{\mathrm{L}}$	$F_{k,\widetilde{\mathrm{df}}(r_m,k)}$	$F_{k,\widehat{\mathrm{df}}(r_m,k)}$	$MR92^h$	-	_	+	_	+	-	+	$\theta$
	L-3	$\hat{D}_{\mathrm{S}}(r_m)$	$\mathscr{D}_{\mathrm{L}}$	$\widehat{r}_{ ext{S}}$	$\mathscr{D}_{\mathrm{L}}$	$F_{k,\widetilde{\mathrm{df}}(r_m,k)}$	$F_{k,\widehat{\mathrm{df}}(r_m,k)}$	Proposal	+	_	_	_	+	+	_	$\theta$
	L-4	$\widehat{D}_{\mathrm{S}}(m{r}_m)$	$\mathscr{D}_{\mathrm{L}}$	$\widehat{r}_{ ext{S}}^{+}$	$\mathscr{D}_{\mathrm{L}}$	$F_{k,\widetilde{\mathrm{df}}(r_m,k)}$	$F_{k,\widehat{\mathrm{df}}(r_m,k)}$	Proposal	+	_	+	_	+	+	_	$\theta$
	L-5	$\hat{D}_{\mathrm{S}}(m{r}_m)$	$\mathscr{D}_{\mathrm{L}}$	$\widehat{r}_{\mathrm{S}}^{\Diamond}$	$\mathscr{D}_{\mathrm{L},1}$	$F_{k,\widetilde{\mathrm{df}}(r_m,h)}$	$F_{k,\widehat{\mathrm{df}}(r_m,h)}$	Proposal	+	+	+	+	+	+	_	$\psi$
	L-6 <sup>i</sup>	$\hat{D}_{\mathrm{L}}(\boldsymbol{r}_{m})$	$\overline{\mathscr{D}}_{\mathrm{L}}$	$\widehat{r}_{ m L}^+$	$\overline{\mathscr{D}}_{\mathrm{L}}$	$F_{k,\widetilde{\mathrm{df}}(r_m,k)}$	$F_{k,\widehat{\mathrm{df}}(r_m,k)}$	Proposal	+	_	+	_	+	+	+	$\theta$
	$L-7^{j}$	$\hat{D}_{\mathrm{L}}(m{r}_m)$	$\overline{\mathscr{D}}_{\mathrm{L}}$	$\widehat{r}_{\rm L}^{\Diamond}$	$\overline{\mathscr{D}}_{\mathrm{L},1}$	$F_{k,\widetilde{\mathrm{df}}(r_m,h)}$	$F_{k,\widehat{\mathrm{df}}(r_m,h)}$	Proposal	+	+	+	+	+	+	+	$\psi$

<sup>&</sup>lt;sup>a</sup>In actual computation, the  $r_m$  in the denominator degree of freedom of F is replaced by its corresponding estimator given in the previous column.

<sup>&</sup>lt;sup>b</sup>The original approximate null distribution documented in Rubin (2004) was modified by Li *et al.* (1991a). This footnote also applies to W-2,4,5.

<sup>&</sup>lt;sup>c</sup>The estimator  $\widetilde{r}'_{W}$  does not depend on  $\theta_{0}$ , but its MSE may be inflated under  $H_{1}$  if a bad parametrization of  $\theta$  is used.

<sup>&</sup>lt;sup>d</sup>The originally proposed combining rule is  $\overline{D}'_{W}(r_{m})$ ; see (3.10). Although  $\overline{D}'_{W}(r_{m})$  is more computational feasible, the power loss is more significant than  $\widetilde{D}'_{W}(r_{m})$  after inserting an inefficient estimator  $\widetilde{r}'_{W,1}$  for  $r_{m}$ . This footnote also applies to W-3.

<sup>&</sup>lt;sup>e</sup>Computing the test statistic  $D_{\mathrm{W}}(T) = d_{\mathrm{W}}(\overline{\theta}, T)/k$  does not require estimating  $r_m$ .

<sup>&</sup>lt;sup>f</sup>EFMI is not required for the test statistic  $D_{\rm W}(T)$ , but it is required for its approximate null distribution.

<sup>&</sup>lt;sup>g</sup>Averaging and processing vector estimators of  $\psi$ , but not their covariance matrixes, is needed. This footnote also applies to L-2.

<sup>&</sup>lt;sup>h</sup>It is a trivial modification of the original proposal in MR92 by replacing  $\tilde{r}_L$  with  $\tilde{r}_L^+ = \max\{0, \tilde{r}_L\}$ .

<sup>&</sup>lt;sup>i</sup>L-6 is equivalent to L-4 when the rows of X are independent.

 $<sup>^{</sup>j}$ L-7 is equivalent to L-5 when the rows of X are independent.

		•						• •	
						Paramet	ers		
$\mathbf{E}$	experiment	I	Fixe	d			Variable	;	
No.	Varying	$\rho$	p	f	Case 1	Case 2	Case 3	Case 4	Case 5
I	Correlation $\rho$	-	2	0.5	-0.8	-0.4	0	0.4	0.8
II	Dimension $p$	0.4	_	0.5	2	3	4	5	6
III	FMI €	0.4	2	_	0.1	0.3	0.5	0.7	0.9

Table 3

The values of parameters used in the simulation experiment in § 4.1.

### 4. Empirical Investigation and Findings.

4.1. Simulation Studies. Suppose that  $X_1, \ldots, X_n \sim \mathcal{N}_p(\mu, \Sigma)$  independently, where  $\Sigma = \sigma^2\{(1-\rho)I_p + \rho \mathbf{1}_p \mathbf{1}_p^{\mathsf{T}}\}$ , and  $\mathbf{1}_p$  is the *p*-vector of ones. The values of p,  $\sigma^2$ ,  $\rho$  and  $\mu$  are specified below. Further assume that only  $n_{\text{obs}} = \lfloor (1-f)n \rfloor$  data points are observed, where  $f \in (0,1)$  is the FMI. Let  $X_{\text{obs}} = \{X_i : i = 1, \ldots, n_{\text{obs}}\}$  and  $X_{\text{mis}} = \{X_i : i = n_{\text{obs}} + 1, \ldots, n\}$ . Suppose that we want to test whether the means of all components are equal, i.e.,  $H_0 : \mu = \mu_0 \mathbf{1}_p$ , where  $\mu_0 \in \mathbb{R}$  is an unknown constant.

Obviously, one may directly use the observed dataset to construct the LRT statistic  $D_{\rm L}$  without MI. The test  $D_{\rm L}$  (denoted by L-0) is regarded as a benchmark for comparison. Throughout this subsection, W-1,2,3,4 and L-1,2,3,4,5 listed in Table 2 are compared. In the imputation step, a Bayesian model is employed for imputation. Assume a multivariate Jeffreys prior on  $(\mu, \Sigma)$ , i.e.,  $f(\mu, \Sigma) \propto |\Sigma|^{-(p+1)/2}$ . Let  $\overline{X}_{\rm obs}$  and  $S_{\rm obs}$  be the sample mean and sample covariance matrix based on  $X_{\rm obs}$ . Then, the  $\ell$ th imputed missing dataset can be produced by the following procedure, for  $\ell = 1, \ldots, m$ .

- 1. Draw a posterior sample  $\Sigma^{\ell}$  from the inverse-Wishart distribution with  $(n_{\text{obs}} 1)$  degrees of freedom and scale matrix  $S_{\text{obs}}^{-1}$ .
- 2. Draw one posterior sample  $\mu^{\ell}$  from  $\mathcal{N}_p(\overline{X}_{\text{obs}}, \Sigma^{\ell}/n_{\text{obs}})$ .
- 3. Draw  $(n-n_{\text{obs}})$  imputed missing values  $\{X_i^{\ell}: i=n_{\text{obs}}+1,\ldots,n\}$  from  $\mathcal{N}_p(\mu^{\ell}, \Sigma^{\ell})$  independently. Also, denote  $X_i^{\ell}=X_i$  for  $i=1,\ldots,n_{\text{obs}}$ .

With the  $\ell$ th completed dataset, the unconstrained MLEs for  $\mu$  and  $\Sigma$  are

$$\widehat{\mu}^{\ell} = \frac{1}{n} \sum_{i=1}^{n} X_i^{\ell}, \qquad \widehat{\Sigma}^{\ell} = \frac{1}{n} \sum_{i=1}^{n} \left( X_i^{\ell} - \widehat{\mu}^{\ell} \right) \left( X_i^{\ell} - \widehat{\mu}^{\ell} \right)^{\mathsf{T}}.$$

Whereas we generate data using a covariance matrix with common variance and correlation, our model does not assume any structure for  $\Sigma$ . The only restriction we can impose is the common-mean assumption under the null,

for which the constrained MLEs are

$$\widehat{\mu}_0^{\ell} = \left\{ \frac{\mathbf{1}_p^{\mathsf{T}}(\widehat{\Sigma}^{\ell})^{-1}\widehat{\mu}^{\ell}}{\mathbf{1}_p^{\mathsf{T}}(\widehat{\Sigma}^{\ell})^{-1}\mathbf{1}_p} \right\} \mathbf{1}_p, \qquad \widehat{\Sigma}_0^{\ell} = \widehat{\Sigma}^{\ell} + \left(\widehat{\mu}^{\ell} - \widehat{\mu}_0^{\ell}\right) \left(\widehat{\mu}^{\ell} - \widehat{\mu}_0^{\ell}\right)^{\mathsf{T}}.$$

In the experiment, we study the impact of parametrization on different test statistics. For the Wald tests, three parametrizations of  $\theta$  are examined:

- (i)  $\theta = (\mu_2 \mu_1, \dots, \mu_p \mu_{p-1})^{\mathsf{T}}$  differences of means, (ii)  $\theta = (\mu_2/\mu_1, \dots, \mu_p/\mu_{p-1})^{\mathsf{T}} \mathbf{1}_{p-1}$  relative differences of means, and (iii)  $\theta = (\mu_2^3 \mu_1^3, \dots, \mu_p^3 \mu_{p-1}^3)^{\mathsf{T}}$  differences of cubic means.

(iii) 
$$\theta = (\mu_2^3 - \mu_1^3, \dots, \mu_p^3 - \mu_{p-1}^3)^\mathsf{T}$$
 — differences of cubic means.

For any case above,  $H_0$  can be re-expressed as  $\theta_0 = \mathbf{0}_{p-1}$ , an (p-1)-vector of zeros. For LRTs, the following parametrizations of  $\psi$  are used:

- (i)  $\psi = {\mu; \Sigma}$  means and covariances,
- (ii)  $\psi = \{\sqrt{\sigma_{ii}}/\mu_i, 1 \leq i \leq p; \Sigma\}$  noise-to-signal and covariances, and (iii)  $\psi = \{\mu^{\mathsf{T}}\Sigma^{-1/2}; \Sigma^{-1}\}$  standardized means and precisions,

where  $\Sigma = (\sigma_{ij})$  and  $\Sigma^{1/2}$  is the symmetric square root of  $\Sigma$ . The dimension of  $\psi$  is  $h = (p^2 + 3p)/2$ .

In the first part of the experiment, we study the distribution of p-values derived from each test under  $H_0$ . In particular, we use n = 100, m = 3,  $\sigma^2 = 5$  and  $\mu = \mathbf{1}_p$ , with various values of  $\rho$ , p and f specified in Table 3. All simulations are repeated 2<sup>12</sup> times. The comparison under parametrization (ii) is shown in Figure 4; whereas those under parametrizations (i) and (iii) are deferred to Appendix C. Note that, for Wald tests under parametrization (ii), the matrix  $U^{\ell}$  is singular in less than 0.25% of the replications, and those cases are removed from the analysis (which should favor the Wald tests).

The empirical sizes (i.e., type-I errors) of the MI Wald tests generally deviate from the nominal size  $\alpha$  under parametrization (ii). In contrast, the sizes of all LRTs are closer to  $\alpha$ . However, the original L-1 and its trivial modification L-2 do not have accurate sizes when  $|\rho|$  or  $\not$  is large. They can be over-sized or under-sized depending on which parametrization is used. Moreover, the trivial modification L-2 does not help to correct the size, and it may even worsen the test. For our test statistics L-3 and L-4, they are invariant to parametrizations and have quite accurate sizes, although they are under-sized in challenging cases where both p and  $\ell$  are large. Moreover, they are identical throughout our simulation experiments, i.e., we never observed  $\hat{r}_{\rm L} < 0$ . For our recommended statistic L-5, it gives the most satisfactory overall results. It generally has very accurate size, except that it is slightly over-sized for large p, a problem that should diminish when we use m beyond the smallest recommended m = 3.

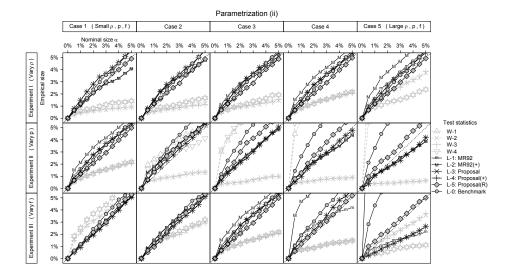


FIG 4. The comparison between empirical size and nominal size  $\alpha$  under parametrization (ii) for  $\alpha \in (0, 5\%]$ . The Wald tests (W-1,2,3,4) and LRTs (L-0,1,2,3,4,5) are represented by grey dashed and black solid lines, respectively. The LRT statistic  $\widetilde{D}_L$  (L-1: MR92) and its modification  $\widetilde{D}_L^+$  (L-2: MR92(+)) are the tests that greatly improved upon by our proposals  $\widehat{D}_S$  (L-3: Proposal),  $\widehat{D}_S^+$  (L-4: Proposal(+)) and  $\widehat{D}_S^{\diamondsuit}$  (L-5: Proposal(R)).

Interestingly, as seen clearly in Figure 4, the benchmark L-0 performs very badly for large p and f. The sample size per parameter, n/h, is small; for  $p \ge 4$ ,  $n/h \le 100/14 < 8$ . The asymptotic null distribution  $\chi_k^2/k$  then can fail badly under arbitrary or even all parametrizations; (ii) apparently falls into this category. An F approximation would be more appropriate. But this is exactly what is being used for MI tests, albeit with different choices of the denominator degrees of freedom. Table 4 documents how often  $\tilde{r}_L$ ,  $\tilde{D}_L$  and  $\hat{r}_S$  are negative. In some cases, nearly half of the simulated values of  $\tilde{r}_L$  and  $\tilde{D}_L$  are negative. In contrast,  $\hat{r}_S$  is always non-negative in our simulation, despite the fact that it can be negative in theory.

To study the power of each test, we set f = 0.5, p = 2,  $\rho = 0.8$ ,  $\sigma^2 = 5$  and  $\mu = (-2 + \delta, -2 + 2\delta)^{\mathsf{T}}$  for different values of  $m \in \{3, 10, 30\}$ ,  $n \in \{100, 400, 1600\}$  and  $\delta = \mu_2 - \mu_1 \in [0, 4]$ . The empirical power functions for size 0.5% tests under parametrizations (i), (ii) and (iii) are plotted in Figure 5. The results for size 5% tests are deferred to Table 12 of the Appendix. Generally, none of the Wald tests exhibits monotonically increasing power as  $\delta$  increases, and their performance is affected significantly by parametrization. In particular, the powers can be as low as zero when  $1 \lesssim \delta \lesssim 2$  under

Table 4 The empirical proportions of negative  $\widetilde{r}_L$  and  $\widetilde{D}_L$ . The results under parametrizations (ii) and (iii) are shown. For parametrization (i),  $\widetilde{r}_L \geqslant 0$  and  $\widetilde{D}_L \geqslant 0$  in the experiments.

		Case									
		1	2	3	4	5	1	2	3	4	5
Experiment	Parametrization		% of	f $\widetilde{r}_{ m L}$	< 0			% c	of $\widetilde{D}_{\mathrm{L}}$	< 0	
I	(ii)	1	2	3	4	5	26	16	13	12	12
	(iii)	6	6	7	7	7	1	1	1	1	2
II	(ii)	4	1	0	0	0	12	5	3	4	3
	(iii)	7	3	1	1	1	1	0	0	0	0
III	(ii)	13	6	4	4	3	55	25	12	5	2
	(iii)	18	9	7	5	4	20	5	1	1	0

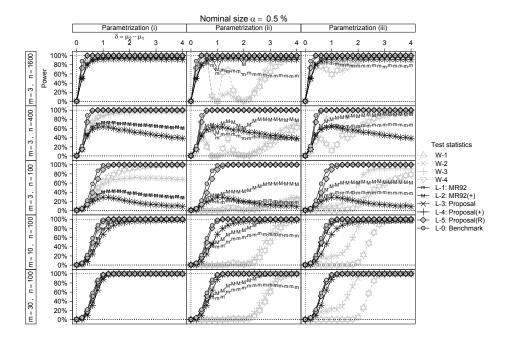


Fig 5. The power curves under different parametrizations. The nominal size is  $\alpha = 0.5\%$ . In each plot, the vertical axis denotes the power, whereas the horizontal axis denotes the value of  $\delta = \mu_2 - \mu_1$ . The legend in Figure 4 also applies here.

parametrizations (ii) and (iii). Under parametrization (ii), L-1 is not powerful even for large  $\delta$ . Moreover, its trivial modifications L-2 cannot retrieve all the power it should have. On the other hand, our first proposed test

Table 5 The range of empirical size  $[\min \hat{\alpha}, \max \hat{\alpha}]$  in percentage, where  $\max$  and  $\min$  are taken over the three parametrizations. Only one value is recorded for those tests that are invariant to parametrization. The nominal size is  $\alpha = 0.5\%$ .

	Range of empirical size: $[\min \hat{\alpha}, \max \hat{\alpha}]/\%$										
(n,m)	(1600, 3)	(400, 3)	(100, 3)	(100, 10)	(100, 30)						
W-1	[0.90, 1.05]	[0.76, 1.05]	[0.20, 1.22]	[0.07, 0.56]	[0.02, 0.49]						
W-2	[0.90, 1.05]	[0.98, 1.22]	[0.93, 1.25]	[0.32, 0.73]	[0.20, 0.85]						
W-3	[0.98, 1.05]	[0.98, 1.25]	[0.90, 1.29]	[0.34, 0.71]	[0.22, 0.73]						
W-4	[0.90, 1.05]	[0.76, 1.05]	[0.20, 1.22]	[0.07, 0.56]	[0.02, 0.49]						
L-1	[0.90, 1.03]	[1.10, 1.64]	[1.15, 1.49]	[0.37, 1.05]	[0.10, 0.46]						
L-2	[0.90, 1.05]	[1.10, 1.76]	[1.15, 2.37]	[0.37, 0.98]	[0.10, 0.49]						
L-3	0.90	1.10	0.83	0.24	0.07						
L-4	0.90	1.10	0.83	0.24	0.07						
L-5	0.46	0.44	0.68	0.46	0.42						
L-0	0.39	0.66	0.66	0.66	0.66						

statistics L-3 and L-4 perform better than L-1 and L-2 at least for large m, however, they also lose a significant amount of power when m is small.

Compared with all these, our recommended test statistic L-5 performs extremely well for all m and n, with power very close to the benchmark L-0 even for small m. To ensure the comparisons of power are fair, we also investigate the empirical (actual) size,  $\hat{\alpha}$ , in comparison to the nominal type-I error  $\alpha$ . Table 5 shows the minimum and maximum of the empirical sizes over the three parametrizations considered in each test — and only one value is needed for those tests that are invariant to parametrization — when the nominal size  $\alpha = 0.5\%$ . We see the deviations from the nominal  $\alpha$  can be noticeable, especially when m=3. To take that into account, we report the empirical size adjusted power, that is,  $O = \text{power}/\hat{\alpha}$ , which also has the interpretation as (an approximated) posterior odds of  $H_1$  to  $H_0$  (Bayarri et al., 2016). Figure 6 plots the result for size 0.5% tests. Compared with the benchmark L-0, the odds O of the proposed robust MI test (L-5) is closer to the nominal value  $1/\alpha$  as  $\delta \to \infty$ . Nevertheless, the finite sample performances of all size 0.5% tests are less satisfactory than those for size 5% tests (given in Appendix) because a very large sample size n is required in order to approximate the tail behavior of test statistics satisfactorily.

We also compare the performance of estimators of  $\mathcal{F}_m$  for different  $\delta$  and parametrizations. In our experiment, we have  $\mathcal{F}_m = 1 + 1/m$  because we have set  $\mathcal{F} = 1$ . The MSEs of estimators  $\hat{f} = \hat{r}/(1+\hat{r})$  of  $f_m = r_m/(1+r_m)$  are

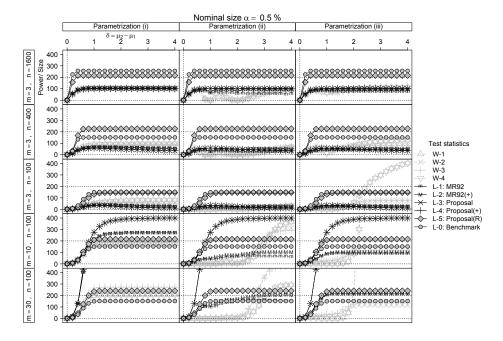


FIG 6. The ratios of empirical power to empirical size under different parametrizations. The nominal size is  $\alpha=0.5\%$ . In each plot, the vertical axis denotes the ratio, whereas the horizontal axis denotes  $\delta=\mu_2-\mu_1$ . The legend in Figure 4 also applies here.

shown in Figure 7, in log scale. Clearly, the only estimator that is consistent, invariant to parametrization and robust against  $\delta$  is our proposal  $\hat{f}_{\rm L}^{\Diamond} = \hat{r}_{\rm L}^{\Diamond}/(1+\hat{r}_{\rm L}^{\Diamond})$ . It concentrates at the true value  $\not f_m$  quite closely even for small m and n. Since  $\hat{f}_{\rm L}^{\Diamond}$  is the only reliable estimator of  $\not f_m$ , it verifies why L-5 has the greatest power. On the other hand, the estimator  $\hat{f}_{\rm L} = \hat{r}_{\rm L}/(1+\hat{r}_{\rm L})$  has very large MSE when  $\delta \neq 0$ . It also explains why L-1 is not powerful.

4.2. Monte Carlo Experiments Without EFMI. To check how robust various tests are to the assumption of EFMI, we simulate  $X_i = (X_{i1}, \ldots, X_{ip})^{\intercal} \stackrel{\text{iid}}{\sim} \mathcal{N}_p(\mu, \Sigma)$  for  $i = 1, \ldots, n$ . Let  $R_{ij}$  be defined by  $R_{ij} = 1$  if  $X_{ij}$  is observed, otherwise  $R_{ij} = 0$ . Suppose that the first variable  $X_{\cdot 1}$  is always observed, and the rest form a monotone missing pattern as defined by a logistic model on the missing propensity:

$$P(R_{ij} = 0 \mid R_{i,j-1} = a) = \begin{cases} [1 + \exp(\alpha_0 + \alpha_1 X_{i,j-1})]^{-1} & \text{if } a = 1; \\ 1 & \text{if } a = 0, \end{cases}$$

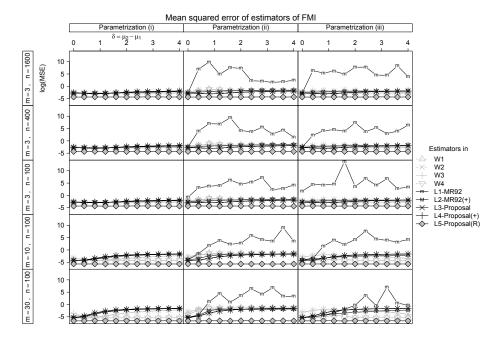


FIG 7. The MSEs of estimators of  $f_m$  used in the test statistics. In each plot, the vertical axis denotes the log of MSE, whereas the horizontal axis denotes the value of  $\delta = \mu_2 - \mu_1$ . The legend in Figure 4 also applies here.

for  $j=2,\ldots,p$ , where  $\alpha_0,\alpha_1\in\mathbb{R}$ . If  $\alpha_1=0$ , then the data are missing completely at random (MCAR); otherwise they are missing at random (MAR), as defined in Rubin (1976). Let  $n_j=\sum_{i=1}^n R_{ij}$  be the number of observed jth component. Without loss of generality, assume  $X_{\rm obs}$  is arranged in such a way that  $R_{ij}\geqslant R_{i'j}$  for all i< i' and j.

To impute the missing data, it is useful to represent  $X_i$  by

$$\begin{bmatrix} X_{i1} \mid \beta_1, \tau_1^2 \end{bmatrix} \sim \mathcal{N}(\beta_1, \tau_1^2),$$
$$\begin{bmatrix} X_{ij} \mid X_{i,1:(j-1)}, \beta_j, \tau_j^2 \end{bmatrix} \sim \mathcal{N}(\beta_j^{\mathsf{T}} Z_{ij}, \tau_j^2), \qquad j = 2, \dots, p,$$

where  $\tau_1^2, \ldots, \tau_p^2 \in \mathbb{R}^+$ ,  $\beta_j \in \mathbb{R}^j$ ,  $X_{i,1:(j-1)} = (X_{i1}, \ldots, X_{i,j-1})^\intercal$  and  $Z_{ij} = (1, X_{i,1:(j-1)}^\intercal)^\intercal$  for  $j \geq 2$ . Denote the (complete-case) least squares estimators of  $\beta_j$  and  $\tau_j^2$  by  $\widehat{\beta}_j = (Z_j^\intercal Z_j)^{-1} Z_j^\intercal W_j$  and  $\widehat{\tau}_j^2 = \frac{1}{n_j - j} (W_j - Z_j \widehat{\beta}_j)^\intercal (W_j - Z_j \widehat{\beta}_j)$ , where  $Z_j = (Z_{1j}, \ldots, Z_{n_jj})^\intercal$  and  $W_j = (X_{1j}, \ldots, X_{n_jj})^\intercal$ .

To perform MI, we assume a Bayesian model with the non-informative prior  $f(\beta_1, \ldots, \beta_p, \tau_1^2, \ldots, \tau_p^2) \propto 1/(\tau_1^2 \cdots \tau_p^2)$ . For each  $\ell = 1, \ldots, m$ , the  $\ell$ th imputed dataset  $X^{\ell}$ , whose (i, j)th element is  $X_{ij}^{\ell}$ , is produced as follows.

- 1. Let  $X_{ij}^{\ell} = X_{ij}$  for all  $1 \leq j \leq p$  and  $i \leq n_j$ .

- 1. Let  $X_{ij} = X_{ij}$  for all  $1 \le j \le p$  and  $i \le n_j$ . 2. For each j = 2, ..., p, repeat Step 3 to Step 5. 3. Draw a sample  $(\tau_j^{\ell})^2$  from  $\hat{\tau}_j^2(n_j j)/\chi_{n_j j}^2$ . 4. Draw a sample  $\beta_j^{\ell}$  from  $\mathcal{N}_j(\hat{\beta}_j, (\tau_j^{\ell})^2 (Z_j^{\mathsf{T}} Z_j)^{-1})$ . 5. Draw a sample  $X_{ij}^{\ell}$  from  $\mathcal{N}((\beta_j^{\ell})^{\mathsf{T}} Z_{ij}^{\ell}, (\tau_j^{\ell})^2)$  for  $i = n_j + 1, ..., n$ , where  $Z_{ij}^\ell = (1, (X_{i,1:(j-1)}^\ell)^{\mathsf{T}})^{\mathsf{T}}.$

We test  $H_0: \mu = \mathbf{0}_p$  against  $H_1: \mu \neq \mathbf{0}_p$ . In the experiments, we set  $\mu = \delta \mathbf{1}_p$ , where  $\delta \in [0, 0.6]$ ; the (i, j)th element of  $\Sigma$  to be  $\Sigma_{ij} = 0.5^{|i-j|}$  for  $i, j = 1, \dots, p; n = 500; m \in \{3, 5\}; p = 5; (\alpha_0, \alpha_1) \in \{(2, -1), (1, 0)\}.$  Our model treats  $\Sigma$  unknown, and hence k=p and  $h=(3p+p^2)/2$ . With the  $\ell$ th imputed dataset, the  $H_0$ -constrained MLEs of  $\mu$  and  $\hat{\Sigma}$  are  $\hat{\mu}_0^{\ell} = \mathbf{0}_p$  and  $\hat{\Sigma}_0^{\ell} = (X^{\ell})^{\mathsf{T}}(X^{\ell})/n$ ; whereas the unconstrained counterparts are  $\hat{\mu}^{\ell} = \mathbf{1}_n^{\mathsf{T}} X^{\ell}/n$  and  $\hat{\Sigma}^{\ell} = (X^{\ell} - \hat{\mu}^{\ell})^{\mathsf{T}}(X^{\ell} - \hat{\mu}^{\ell})/n$ . Under  $H_0$  and MAR, the fractions of missing observations of the five variables are (0, 16%, 28%, 38%, 47%), whereas the average fractions of missing information, i.e, the eigenvalues of  $\mathscr{B}_{\theta}\mathscr{T}_{\theta}^{-1}$ , are (0, 19%, 34%, 45%, 55%). So, the assumption of EFMI does not hold.

We compare (L4)  $\hat{D}_{L}^{+} \approx F_{k,\hat{\mathrm{df}}(\hat{r}_{L}^{+},k)}$ , (L5)  $\hat{D}_{L}^{\Diamond} \approx F_{k,\hat{\mathrm{df}}(\hat{r}_{L}^{\Diamond},h)}$ , (C1) completedata (asymptotic) LRT using  $\{X_i : i = 1, ..., n\}$ , and (C2) complete-case (asymptotic) LRT using  $\{X_i : i = 1, \dots, n_p\}$ . The results are shown in Figure 8. The size of  $\hat{D}_{L}^{\Diamond}$  is quite accurate when the nominal size is small. If the data are MCAR, complete-case test C2 is valid, however, with slightly less power. (Test C2 is typically invalid without MCAR.) In terms of power-to-size ratio, the performance of  $\widehat{D}_{\rm L}^{\Diamond}$  is the best among the three implementable tests L4, L5 and C2. Its performance is comparable to the (unavailable) completedata test C1. Note also that the power-to-size ratio of  $\widehat{D}_{L}^{+}$  and  $\widehat{D}_{L}^{\Diamond}$  become closer to the nominal value 1/0.5% when m increases. All these indicate that the performance of our proposed tests are acceptable despite of the serious violation of the EFMI assumption.

4.3. Applications to a Care-Survival Data. Meng and Rubin (1992) applied their test to the data given in Table 6, where i, j and k index, respectively, clinic (A or B), amount of parental care (more or less) and survival status (died or survived). However, the clinic label k is missing for some of the observations (and the missing-data mechanism was assumed to be ignorable). Two hypotheses were tested in Meng and Rubin (1992). The first is whether the clinic and parental care are conditionally independent given the survival status, and the second is whether all three variables are mutually independent. The MI datasets are generated from a Bayesian model in § 4.2 of Meng and Rubin (1992). Our aim here is to investigate the impact on the

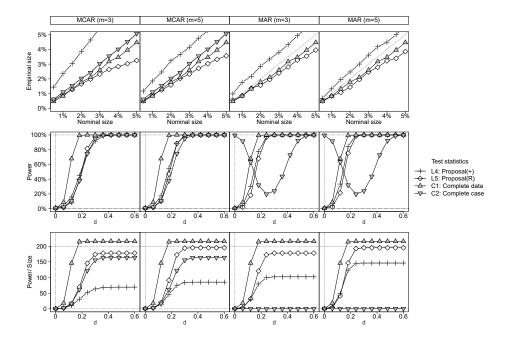


FIG 8. The empirical size, empirical power, and their ratio. The first row of plots show the empirical sizes. The size of the complete-case test (C2) under MAR is off the chat (always equals to one) in the experiment because it is invalid. The second and third rows of plots show the powers and the power-to-size ratios, respectively, where the nominal size is 0.5%.

test statistics  $\widetilde{D}_{S}$ ,  $\widehat{D}_{S}^{+}$  and  $\widehat{D}_{S}^{\Diamond}$  by different parametrizations of  $\{\pi_{ijk}\}$ ; and the impact on the estimators  $\widetilde{r}_{L}$ ,  $\widehat{r}_{S}^{+}$  and  $\widehat{r}_{S}^{\Diamond}$  under different null hypotheses.

Table 6
Data from Meng and Rubin (1992). The notation "?" indicates missing label.

		Survival Status (j)		
Clinic $(k)$	$ {\bf Parental \ care} \ (i)$	Died	Survived	
A	Less	3	176	
	More	4	293	
В	Less	17	197	
	More	2	23	
?	Less	10	150	
	More	5	90	

Specifically, the  $\ell$ th imputed log-likelihood function is  $\log f(X^{\ell} \mid \pi) = \sum_{c} n_{c}^{\ell} \log \pi_{c}$ , where  $X^{\ell}$  are the cell counts  $n_{c}^{\ell}$  in the  $\ell$ th imputed dataset.

Hence the unconstrained MLE of  $\pi_c$  is  $\hat{\pi}_c^\ell = n_c^\ell/n_+^\ell$ , where  $n_+^\ell = \sum_c n_c^\ell$ . Consequently, the joint log-likelihood based on the stacked data is

(4.1) 
$$\log f(X^{S} \mid \pi) = \sum_{\ell=1}^{m} \sum_{c} n_{c}^{\ell} \log \pi_{c} = \sum_{c} n_{c}^{+} \log \pi_{c},$$

where  $n_c^+ = \sum_{\ell=1}^m n_c^\ell$ . Thus the unconstrained MLE with respect to (4.1) is  $\hat{\pi}_c^{\rm S} = n_c^+/n_+^+$ , where  $n_+^+ = \sum_c n_c^+$ . Similarly, we can find the constrained MLEs under a given null. We consider the following parametrizations:

- (i)  $\psi_{ijk} = \pi_{ijk}$  the identity map,
- (ii)  $\psi_{ijk} = \log\{\pi_{ijk}/(1-\pi_{ijk})\}$  the logit transformation, and (iii)  $\psi_{ij1} = \pi_{ij1}$  and  $\psi_{ij2} = \pi_{ij2}/\pi_{ij1}$  ratios of probabilities.

The *p*-values  $\widetilde{p}_{\rm L}$ ,  $\widehat{p}_{\rm S}^+$  and  $\widehat{p}_{\rm S}^{\Diamond}$  of the tests  $\widetilde{D}_{\rm L}$ ,  $\widehat{D}_{\rm S}^+$  and  $\widehat{D}_{\rm S}^{\Diamond}$  and the associated estimates of  $r_m$ , i.e.,  $\widetilde{r}_L$ ,  $\widehat{r}_S^+$  and  $\widehat{r}_S^{\Diamond}$ , are shown in Table 7. A more detailed comparison is deferred to Table 9 of the Appendix.

The LRTs using  $\widetilde{D}_L$ ,  $\widehat{D}_S^+$  and  $\widehat{D}_S^{\diamondsuit}$  under different parametrizations in § 4.3.

	Parametrization (i)										
ĺ	<i>H</i> <sub>0</sub> : Co	nditional indepen	dence	H <sub>0</sub> : Full in							
m	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\Diamond}$	$\widetilde{D}_{\mathrm{L}},\widehat{D}_{\mathrm{S}}^{+},\widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{ m L}, \widehat{p}_{ m S}^+, \widehat{p}_{ m S}^{\diamondsuit}$	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\Diamond}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{\mathrm{L}}, \widehat{p}_{\mathrm{S}}^{+}, \widehat{p}_{\mathrm{S}}^{\diamondsuit}$					
3	0.54, 0.54, 0.38	0.08, 0.08, 0.09	0.93, 0.93, 0.92	0.31, 0.31, 0.38	54.2, 54.2, 51.4	0,0,0					
10	0.50, 0.50, 0.70	0.14, 0.14, 0.12	0.87, 0.87, 0.88	0.56, 0.56, 0.70	45.4, 45.4, 41.7	0, 0, 0					
50	0.31, 0.31, 0.45	0.11, 0.11, 0.10	0.90, 0.90, 0.91	0.33, 0.33, 0.45	51.5, 51.5, 47.3	0, 0, 0					
	Parametrization (ii)										
ĺ	<i>H</i> <sub>0</sub> : Co	nditional indepen	idence	$H_0$ : Full in							
m	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\lozenge}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{ m L}, \widehat{p}_{ m S}^+, \widehat{p}_{ m S}^{\diamondsuit}$	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\lozenge}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{\mathrm{L}},\widehat{p}_{\mathrm{S}}^{+},\widehat{p}_{\mathrm{S}}^{\diamondsuit}$					
3	1.08, 0.54, 0.38	-0.07, 0.08, 0.09	1.00, 0.93, 0.92	0.61, 0.31, 0.38	43.9, 54.2, 51.4	0,0,0					
10	0.99, 0.50, 0.70	-0.10, 0.14, 0.12	1.00, 0.87, 0.88	1.09, 0.56, 0.70	33.7, 45.4, 41.7	0, 0, 0					
50	0.63, 0.31, 0.45	-0.10, 0.11, 0.10	1.00, 0.90, 0.91	0.65, 0.33, 0.45	41.3, 51.5, 47.3	0, 0, 0					
		Par	ametrization	(iii)							
ĺ	<i>H</i> <sub>0</sub> : Co	nditional indepen	dence	$H_0$ : Full in	dependence						
m	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\Diamond}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{\mathrm{L}}, \widehat{p}_{\mathrm{S}}^{+}, \widehat{p}_{\mathrm{S}}^{\Diamond}$	$\widetilde{r}_{ m L}, \widehat{r}_{ m S}^+, \widehat{r}_{ m S}^{\lozenge}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{\mathrm{L}},\widehat{p}_{\mathrm{S}}^{+},\widehat{p}_{\mathrm{S}}^{\diamondsuit}$					
3	-2.35, 0.54, 0.38	-1.16, 0.08, 0.09	1.00, 0.93, 0.92	-1.22, 0.31, 0.38	-321, 54.2, 51.4	1,0,0					
10	-2.04, 0.50, 0.70	-2.20, 0.14, 0.12	1.00, 0.87, 0.88	-1.85, 0.56, 0.70	-86, 45.4, 41.7	1, 0, 0					
50	-1.22, 0.31, 0.45	-7.39, 0.11, 0.10	1.00, 0.90, 0.91	$\left  -1.06, 0.33, 0.45 \right $	-1136, 51.5, 47.3	1, 0, 0					

The simulation outputs demonstrate that  $\hat{D}_{\mathrm{S}}^{+}$  and  $\hat{D}_{\mathrm{S}}^{\Diamond}$  are invariant to parametrizations, whereas  $\widetilde{D}_{L}$  is not. Moreover, the impact on  $\widetilde{D}_{L}$  is large under the parametrization (iii). In particular, the value of  $\tilde{r}_{\rm L}$  is inflated; and some of the values of  $\tilde{r}_{\rm L}$  and  $\tilde{D}_{\rm L}$  are negative, leading to the meaningless  $\tilde{p}_{\rm L}=1$ , especially under parametrization (iii). In contrast,  $\hat{r}_{\rm S}\geqslant 0$  for all cases in this example (and hence  $\hat{r}_{\rm S}^+=\hat{r}_{\rm S}$ ). In addition,  $\hat{D}_{\rm S}^+\approx\hat{D}_{\rm S}^{\Diamond}$  for testing the conditional independence, a hypothesis that is not rejected by either test. In contrast, for testing the full independence,  $\hat{D}_{\rm S}^+$  and  $\hat{D}_{\rm S}^{\Diamond}$  are not very close to each other, but they both lead to essentially zero p-value, and hence both reject the null hypothesis. These results reconfirm the conclusions in Meng and Rubin (1992). Last but not least, the estimator  $\hat{r}_{\rm S}^{\Diamond}$  does not change under different null hypotheses, however it is not true for  $\tilde{r}_{\rm L}$  and  $\hat{r}_{\rm S}^+$ .

- 5. Conclusions, Limitations and Future Work. In addition to conducting a general comparative study of MI tests, we proposed two particularly promising MI LRT based on  $\hat{D}_{\rm S}^{\Diamond} = \hat{D}_{\rm S}(\hat{r}_{\rm S}^{\Diamond})$  and  $\hat{D}_{\rm S}^{+} = \hat{D}_{\rm S}(\hat{r}_{\rm S}^{+})$ . Both test statistics are non-negative, invariant to re-parametrizations, and powerful to reject a false null hypothesis (at least for large enough m). Test  $\hat{D}_{\rm S}^{\Diamond}$  is most principled, and the resulting test has the desirable monotonically increasing power as  $H_1$  departs from  $H_0$ . However, it is derived under the stronger assumption of EFMI for  $\psi$ , not just for  $\theta$ ; and row independence of  $X_{\rm com}$  is needed for the ease of computation. (The computationally more demanding test based on  $\hat{D}_{\rm L}(\hat{r}_{\rm L}^{\Diamond})$  relaxes the independence assumption.) The main advantage of  $\hat{D}_{\rm S}^{+}$  is that it is easier to compute, as it requires only standard complete-data computer subroutines for likelihood ratio tests. One drawback is that the ad hoc fix  $\hat{r}_{\rm S}^{+} = \max(0, \hat{r}_{\rm S})$  is inconsistent in general. However, the inconsistency does not appear to significantly affect the asymptotic power, at least in our experiments. Whereas  $\hat{D}_{\rm S}^{+}$  and  $\hat{D}_{\rm S}^{\Diamond}$  significantly improve over existing counterparts, more studies are needed, for reasons listed below.
  - When the missing data mechanism is not ignorable but the imputers fail to fully take that into account, the issue of uncongeniality becomes critical (Meng, 1994a). Xie and Meng (2017) provides theoretical tools for addressing such an issue in the context of estimation, and research is needed to extend their findings to the setting of hypothesis testing.
  - Although the violation of the EFMI assumption may not (seriously) invalidate a test, it will affect its power. It is therefore desirable to explore MI tests without this assumption.
  - The robust  $\widehat{D}_{S}^{\Diamond}$  relies on a stronger assumption of EFMI on  $\psi$ . We can modify it so only EFMI on  $\theta$  is required, but the modification may be very difficult to compute and may require users to have access to non-trivial complete-data procedures. Hence a computational feasible robust test that only assumes EFMI on  $\theta$  needs to be developed.

• Because the FMI is a fundamental nuisance parameter here and there is no (known) pivotal quantity, all MI tests are approximate in nature. In particular, they all have the potential of doing poorly when FMI is large and/or m is small. It is therefore of both theoretical and practical interest to seek powerful MI tests that are least affected by FMI.

#### References.

- Barnard, J. and Rubin, D. B. (1999) Small-sample degrees of freedom with multiple imputation. *Biometrika*, **86**, 948–955.
- Bayarri, M. J., Benjamin, D. J., Berger, J. O. and Sellke, T. M. (2016) Rejection odds and rejection ratios: A proposal for statistical practice in testing hypotheses. *Journal of Mathematical Psychology*, **72**, 90–103.
- Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E.-J., Berk, R., Bollen, K. A., Brembs, B., Brown, L., Camerer, C. et al. (2018) Redefine statistical significance. Nature Human Behaviour, 2, 6–10.
- Berglund, P. and Heeringa, S. G. (2014) Multiple imputation of missing data using SAS. SAS Institute.
- Blocker, A. W. and Meng, X.-L. (2013) The potential and perils of preprocessing: Building new foundations. *Bernoulli*, 19, 1176–1211.
- Carlin, J. B., Galati, J. C. and Royston, P. (2008) A new framework for managing and analyzing multiply imputed data in stata. *The Stata Journal*, **8**, 49–67.
- Cox, D. R. and Reid, N. (1987) Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society B*, **49**, 1–39.
- Grund, S., Robitzsch, A. and Luedtke, O. (2017) Tools for Multiple Imputation in Multilevel Modeling.
- Harel, O. and Zhou, X.-H. (2007) Multiple imputation review of theory, implementation and software. *Statistics in medicine*, **26**, 3057–3077.
- Holan, S. H., Toth, D., Ferreira, M. A. R. and F., K. A. (2010) Bayesian multiscale multiple imputation with implications for data confidentiality. *Journal of the American Statistical Association*, 105, 564–577.
- Horton, N. and Kleinman, K. P. (2007) Much ado about nothing: A comparison of missing data methods and software to fit incomplete data regression models. *The American Statistician*, **61**, 79–90.
- Kenward, M. G. and Carpenter, J. R. (2007) Multiple imputation: current perspectives. Statistical Methods in Medical Research, 16, 199–218.
- Kim, J. K. and Shao, J. (2013) Statistical Methods for Handling Incomplete Data. Chapman and Hall/CRC.
- Kim, J. K. and Yang, S. (2017) A note on multiple imputation under complex sampling. *Biometrika*, **104**, 221–228.
- King, G., Honaker, J., Joseph, A. and Scheve, K. (2001) Analyzing incomplete political science data: An alternative algorithm for multiple imputation. *American Political Science Review*, **95**, 49–69.
- Li, K. H., Meng, X.-L., Raghunathan, T. E. and Rubin, D. B. (1991a) Significance levels from repeated *p*-values with multiply-imputed data. *Statistica Sinica*, **1**, 65–92.
- Li, K. H., Raghunathan, T. E. and Rubin, D. B. (1991b) Large-sample significance levels from multiply imputed data using moment-based statistics and an F reference distribution. Journal of the American Statistical Association, 86, 1065–1073.
- Medeiros, R. (2008) Likelihood ratio tests for multiply imputed datasets: Introducing milrtest.

- Meng, X.-L. (1994a) Multiple-imputation inferences with uncongenial sources of input. Statistical Science, 9, 538–573.
- Meng, X.-L. (1994b) Posterior predictive p-values. The Annals of Statistics, 22, 1142–1160. Meng, X.-L. (2002) Discussion of "Bayesian measures of model complexity and fit" by
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P. and Van Der Linde, A. *Journal of the Royal Statistical Society B*, **64**, 633.
- Meng, X.-L. and Rubin, D. B. (1992) Performing likelihood ratio tests with multiply-imputed data sets. *Biometrika*, **79**, 103–111.
- Peugh, J. L. and Enders, C. K. (2004) Missing data in educational research: A review of reporting practices and suggestions for improvement. Review of Educational Research, 74, 525–556.
- Rose, R. A. and Fraser, M. W. (2008) A simplified framework for using multiple imputation in social work research. *Social Work Research*, **32**, 171–178.
- Royston, P. and White, I. R. (2011) Multiple imputation by chained equations (mice): Implementation in stata. *Journal of Statistical Software*, **45**, 1–20.
- Rubin, D. B. (1976) Inference and missing data. Biometrika, 63, 581–592.
- Rubin, D. B. (1978) Multiple imputations in sample surveys a phenomenological Bayesian approach to nonresponse. *Proceedings of the Survey Research Methods Section of the American Statistical Association*, 20–34.
- Rubin, D. B. (1996) Multiple imputation after 18+ years. Journal of the American statistical Association, 91, 473–489.
- Rubin, D. B. (2004) Multiple Imputation for Nonresponse in Surveys. John Wiley & Sons. Rubin, D. B. and Schenker, N. (1986) Multiple imputation for interval estimation from simple random samples with ignorable nonresponse. Journal of the American Statistical Association, 81, 366–374.
- Schafer, J. L. (1999) Multiple imputation: A primer. Statistical Methods in Medical Research, 8, 3–15.
- Serfling, R. J. (2001) Approximation Theorems of Mathematical Statistics. Wiley-Interscience.
- Su, Y.-S., Gelman, A., Hill, J. and Yajima, M. (2011) Multiple imputation with diagnostics (mi) in R: Opening windows into the black box. *Journal of Statistical Software*, **45**, 1–31.
- Tu, X. M., Meng, X.-L. and Pagano, M. (1993) The AIDS epidemic: Estimating survival after AIDS diagnosis from surveillance data. *Journal of the American Statistical Asso*ciation, 88, 26–36.
- van Buuren, S. and Groothuis-Oudshoorn, K. (2011) Mice: Multivariate imputation by chained equations in R. *Journal of Statistical Software*, **45**, 1–67.
- van Buuren S (2012) Flexible Imputation of Missing Data. Chapman and Hall/CRC.
- van der Vaart, A. W. (2000) Asymptotic Statistics. Cambridge University Press.
- Wallace, D. L. (1980) The Behrens-Fisher and Fieller-Creasy problems. In R. A. Fisher: An Appreciation (eds. S. E. Fienberg and D. V. Hinkley), 119–147. Springer New York.
- Wang, N. and Robins, J. M. (1998) Large-sample theory for parametric multiple imputation procedures. *Biometrika*, 85, 935–948.
- Xie, X. and Meng, X.-L. (2017) Dissecting multiple imputation from a multi-phase inference perspective: What happens when God's, imputer's and analyst's models are uncongenial? (with discussion). *Statistica Sinica*, **27**, 1485–1594.

### APPENDIX A: SUPPLEMENTARY RESULTS

**A.1.** Another Motivation for  $\hat{r}_L^{\Diamond}$ . The definition of  $\hat{r}_L^{\Diamond}$  can also be motivated by the following observation. First, observe that one simple method to construct an always non-negative estimator of  $r_m$  is to perturb  $\hat{\psi}_0^*$  and  $\hat{\psi}_0^{\ell}$  by a suitable amount, say  $\Delta$ , so that the perturbed version of  $\hat{r}_L$  is always non-negative, and is still asymptotically equivalent to the original  $\hat{r}_L$ . We show, in Theorem A.1 below, that the right amount of  $\Delta$  is  $\Delta = \hat{\psi}^* - \hat{\psi}_0^*$ . Using the perturbed version of  $\hat{r}_L$ , we obtain

$$\widehat{r}_{\mathrm{L}}^{\triangle} = \frac{m+1}{k(m-1)}\widehat{\delta}_{\mathrm{L}}^{\triangle},$$

where

$$\widehat{\delta}_{\mathrm{L}}^{\triangle} = \frac{2}{m} \sum_{\ell=1}^{m} \log \left\{ \frac{f(X^{\ell} \mid \widehat{\psi}^{\ell})}{f(X^{\ell} \mid \widehat{\psi}^{*})} \frac{f(X^{\ell} \mid \widehat{\psi}^{*}_{0} + \Delta)}{f(X^{\ell} \mid \widehat{\psi}^{\ell}_{0} + \Delta)} \right\} = \frac{1}{m} \sum_{\ell=1}^{m} d_{\mathrm{L}}(\widehat{\psi}^{\ell}_{0} + \Delta, \widehat{\psi}^{\ell} \mid X^{\ell}).$$

Then we have the following result.

THEOREM A.1. Suppose  $RC_{\theta}$ . Under  $H_0$ , we have (i)  $\hat{r}_L^{\triangle} \ge 0$  for all m, n; and (ii)  $\hat{r}_L^{\triangle} = \hat{r}_L$  as  $n \to \infty$  for each m.

Although  $\hat{r}_{\rm L}^{\triangle} \geqslant 0$ , it is only invariant to affine transformations, and not robust against  $\theta_0$ , and less computational feasible than  $\hat{r}_{\rm L}$ ; see § 3. However, it gives us some insights on how to construct a potentially better estimator. Note that, in (A.1), the constrained MLE is not used in  $d_{\rm L}(\cdot,\cdot\mid X^{\ell})$ , but it is still always non-negative. We call this a "pseudo" LRT statistics. Then,  $\hat{\delta}_{\rm L}^{\triangle}$  is just a multiple of an average of many "pseudo" LRT statistics. In order to find a good estimator of  $r_m$ , we may seek for an estimator which admits this form. Indeed, our proposed estimator  $\hat{r}_{\rm L}^{\triangle}$  also takes the same form:

$$\widehat{r}_{\mathbf{L}}^{\Diamond} = \frac{m+1}{h(m-1)} \frac{1}{m} \sum_{\ell=1}^{m} d_{\mathbf{L}}(\widehat{\psi}^*, \widehat{\psi}^{\ell} \mid X^{\ell}).$$

**A.2. Results for Dependent Data.** This is a supplement for § 3.1. If the data are not independent,  $\hat{d}_{\rm L} \simeq \hat{d}_{\rm S}$  is still true under the following conditions.

Assumption 5. (a) Define  $R(\psi) = \overline{\underline{L}}^{S}(\psi) - \overline{\underline{L}}(\psi)$ , where  $\overline{\underline{L}}(\psi) = (mn)^{-1} \sum_{\ell=1}^{m} \log f(X^{\ell} \mid \psi)$  and  $\overline{\underline{L}}^{S}(\psi) = (mn)^{-1} \log f(X^{S} \mid \psi)$ . For each m, as  $n \to \infty$ ,

$$\sup_{\psi \in \Psi} |R(\psi)| = O_p(1/n), \qquad \sup_{\psi \in \Psi} \left| \frac{\partial}{\partial \psi} R(\psi) \right| = O_p(1/n).$$

(b) For each m, there exists a continuous function  $\psi \mapsto \overline{\mathcal{Z}}(\psi)$ , which is free of n but may depend on m, such that, as  $n \to \infty$ ,

$$\sup_{\psi \in \Psi} \left| \underline{\underline{L}}(\psi) - \underline{\underline{\mathscr{L}}}(\psi) \right| = o_p(1).$$

(c) Let  $\psi_0^* = \arg\max_{\psi \in \Psi} \frac{\overline{\mathcal{L}}(\psi)}{\psi(\theta) = \theta_0} \frac{\overline{\mathcal{L}}(\psi)}{\overline{\mathcal{L}}(\psi)}$  and  $\psi^* = \arg\max_{\psi \in \Psi} \frac{\overline{\mathcal{L}}(\psi)}{\overline{\mathcal{L}}(\psi)}$ . For any fixed m, and for all  $\varepsilon > 0$ , there exists  $\delta > 0$  such that

$$\sup_{\substack{\psi \in \Psi : |\psi_0^* - \psi| > \varepsilon \\ \theta(\psi) = \theta_0}} \left\{ \underline{\underline{\mathcal{Z}}}(\psi_0^*) - \underline{\underline{\mathcal{Z}}}(\psi) \right\} \geqslant \delta, \quad \sup_{\substack{\psi \in \Psi : |\psi^* - \psi| > \varepsilon \\ }} \left\{ \underline{\underline{\mathcal{Z}}}(\psi^*) - \underline{\underline{\mathcal{Z}}}(\psi) \right\} \geqslant \delta.$$

Conditions (b) and (c) in Assumption 5 are standard RCs that are usually assumed for M-estimators (see § 5 of van der Vaart (2000)); whereas condition (a) is satisfied by many models (see Example A.1 below).

THEOREM A.2. Suppose  $RC_{\theta}$  and Assumption 5. Under both  $H_0$  and  $H_1$ , we have (i)  $\hat{d}_S, \hat{r}_S \ge 0$  for all m, n; (ii)  $\hat{d}_S, \hat{r}_S$  are invariant to the parametrization of  $\psi$  for all m, n; and (iii)  $\hat{d}_L = \hat{d}_S$  and  $\hat{r}_L = \hat{r}_S$  as  $n \to \infty$  for each m.

Theorem A.2 implies that the handy test statistics  $\hat{D}_{S}$  and  $\hat{D}_{S}^{+}$  approximate  $\hat{D}_{L}$  and  $\hat{D}_{L}^{+}$  for dependent data, provided that Assumption 5 holds.

EXAMPLE A.1. Consider a stationary autoregressive model of order one. Suppose the complete data  $X = (X_1, \ldots, X_n)^{\mathsf{T}}$  is generated as following:  $X_1 \sim \mathcal{N}(0, v^2)$  and  $[X_i | X_{i-1}] \sim \mathcal{N}(\phi X_{i-1}, \sigma^2)$  for  $i \geq 2$ , where  $v^2 = \sigma^2(1+\phi)/(1-\phi)$ . Then  $\psi = (\phi, \sigma^2)^{\mathsf{T}}$ , and

$$\begin{split} \underline{\overline{L}}(\psi) &= -\frac{1}{2}\log(2\pi) - \frac{1}{2n}\log v^2 - \frac{1}{mn}\sum_{\ell=1}^m \frac{X_1^\ell}{2v^2} - \frac{n-1}{2n}\log\sigma^2 \\ &- \frac{1}{mn}\sum_{\ell=1}^m \sum_{i=2}^n \frac{(X_i^\ell - \phi X_{i-1}^\ell)^2}{2\sigma^2}, \\ \underline{\overline{L}}^{\mathrm{S}}(\psi) &= -\frac{1}{2}\log(2\pi) - \frac{1}{2mn}\log v^2 - \frac{(X_1^1)^2}{2mnv^2} - \frac{mn-1}{2mn}\log\sigma^2 \\ &- \frac{1}{mn}\sum_{\ell=1}^m \sum_{i=2}^n \frac{(X_i^\ell - \phi X_{i-1}^\ell)^2}{2\sigma^2} - \frac{1}{mn}\sum_{\ell=2}^m \frac{(X_1^\ell - \phi X_n^{\ell-1})^2}{2\sigma^2}. \end{split}$$

Then, it is easy to see that condition (a) of Assumption 5 is satisfied.

### APPENDIX B: PROOFS

PROOF OF THEOREM 2.1. (i, ii) From (2.3), we know  $\hat{d}_L \ge 0$  is invariant to parametrization  $\psi$ . (iii) Since  $\hat{d}_L$  is invariant to transformation of  $\psi$ , we assume, without loss of generality, that  $\psi$  admits a parameterization such that  $\operatorname{Cov}(\hat{\theta}^\ell, \hat{\eta}^\ell) \simeq \mathbf{0}$  by taking suitable linear transformation of  $\psi$ . Also write  $U_{\eta}^\ell$  as an efficient estimator of  $\operatorname{Var}(\hat{\eta})$  based on  $X^\ell$ ; and recall that  $U_{\theta}^\ell = U^\ell$  is an efficient estimator of  $\operatorname{Var}(\hat{\theta})$  based on  $X^\ell$ .

Using Taylor's expansion on  $\psi \mapsto \overline{L}(\psi) = m^{-1} \sum_{\ell=1}^m \log f(X^\ell \mid \psi)$  around  $\widehat{\psi}^* = ((\widehat{\theta}^*)^\intercal, (\widehat{\eta}^*)^\intercal)^\intercal$ , we know that for  $\psi \simeq \widehat{\psi}^*$ ,

(B.1) 
$$\overline{L}(\psi) \simeq \overline{L}(\widehat{\psi}^*) - \frac{1}{2} \left( \psi - \widehat{\psi}^* \right)^{\mathsf{T}} \overline{I}(\widehat{\psi}^*) \left( \psi - \widehat{\psi}^* \right),$$

where  $\overline{I}(\psi) = -\partial^2 \overline{L}(\psi)/\partial \psi \partial \psi^{\dagger}$ , which satisfies

(B.2) 
$$\overline{I}(\widehat{\psi}^*) \simeq \begin{pmatrix} \overline{U}_{\theta}^{-1} & \mathbf{0} \\ \mathbf{0} & \overline{U}_{\eta}^{-1} \end{pmatrix}$$

with  $\overline{U}_{\eta} = m^{-1} \sum_{i=1}^{m} U_{\eta}^{\ell}$ . Under the null,  $\hat{\psi}^* = \hat{\psi}_0^*$ . So, using (B.1), we have

$$\widehat{d}_{L} \simeq \left(\widehat{\psi}_{0}^{*} - \widehat{\psi}^{*}\right)^{\mathsf{T}} \overline{I}(\widehat{\psi}^{*}) \left(\widehat{\psi}_{0}^{*} - \widehat{\psi}^{*}\right), 
\simeq \left(\begin{array}{c} \theta_{0} - \widehat{\theta}^{*} \\ \widehat{\eta}(\theta_{0}) - \widehat{\eta}(\widehat{\theta}^{*}) \end{array}\right)^{\mathsf{T}} \left(\begin{array}{c} \overline{U}_{\theta}^{-1} & \mathbf{0} \\ \mathbf{0} & \overline{U}_{\eta}^{-1} \end{array}\right) \left(\begin{array}{c} \theta_{0} - \widehat{\theta}^{*} \\ \widehat{\eta}(\theta_{0}) - \widehat{\eta}(\widehat{\theta}^{*}) \end{array}\right) 
(B.3) \simeq (\overline{\theta}^{\mathsf{T}} - \theta_{0}) \overline{U}_{\theta}^{-1} (\overline{\theta}^{\mathsf{T}} - \theta_{0}) = \widetilde{d}'_{W},$$

where we have used (a)  $\hat{\theta}^* = \overline{\theta}$ ; see, e.g., Lemma 1 of Wang and Robins (1998), and (b)  $\hat{\eta}(\theta_0) - \hat{\eta}(\hat{\theta}^*) = O_p(1/n)$  if  $\theta_0 - \hat{\theta}^* = O_p(1/\sqrt{n})$ ; see Cox and Reid (1987). Since  $\tilde{d}_W' = \tilde{d}_L$  (Meng and Rubin, 1992), we have  $\hat{d}_L = \tilde{d}_L$ .

PROOF OF PROPOSITION 2.2. The given condition implies that

$$\begin{split} \widehat{\psi}^{\ell} &= ((\widehat{\theta}^{\ell})^{\mathsf{T}}, (\widehat{\eta}^{\ell})^{\mathsf{T}})^{\mathsf{T}}, \qquad \widehat{\psi}^{\ell}_{0} &= (\theta^{\mathsf{T}}_{0}, (\widehat{\eta}^{\ell})^{\mathsf{T}})^{\mathsf{T}}, \\ \widehat{\psi}^{*} &= ((\widehat{\theta}^{*})^{\mathsf{T}}, (\widehat{\eta}^{*})^{\mathsf{T}})^{\mathsf{T}}, \qquad \widehat{\psi}^{*}_{0} &= (\theta^{\mathsf{T}}_{0}, (\widehat{\eta}^{*})^{\mathsf{T}})^{\mathsf{T}}. \end{split}$$

Clearly, we also have the decomposition:  $L^{\ell}(\psi) = L^{\ell}_{\dagger}(\theta) + L^{\ell}_{\ddagger}(\eta)$  for all  $\ell$ , where  $L^{\ell}_{\dagger}(\theta) = L_{\dagger}(\theta \mid X^{\ell})$  and  $L^{\ell}_{\ddagger}(\eta) = L_{\ddagger}(\eta \mid X^{\ell})$ . Then,

$$\overline{d}_{L} - \widehat{d}_{L} = \frac{2}{m} \sum_{\ell=1}^{m} \left\{ L^{\ell}(\widehat{\psi}^{\ell}) - L^{\ell}(\widehat{\psi}^{\ell}_{0}) - L^{\ell}(\widehat{\psi}^{*}) + L^{\ell}(\widehat{\psi}^{*}_{0}) \right\}$$

$$= \frac{2}{m} \sum_{\ell=1}^{m} \left\{ L^{\ell}_{\uparrow}(\widehat{\theta}^{\ell}) - L^{\ell}_{\uparrow}(\widehat{\theta}^{*}) \right\} \geqslant 0$$

since 
$$L_{\dagger}^{\ell}(\widehat{\theta}^{\ell}) \geqslant L_{\dagger}^{\ell}(\widehat{\theta}^{*})$$
 for all  $\ell$ .

PROOF OF COROLLARY 2.3. Applying Taylor's expansion on  $\psi \mapsto L^{\ell}(\psi)$ , we can find  $\check{\psi}^{\ell}$  lying on the line segment joining  $\hat{\psi}^{\ell}$  and  $\hat{\psi}^{\ell}_0$  such that

$$L^{\ell}(\widehat{\psi}_{0}^{\ell}) = L^{\ell}(\widehat{\psi}^{\ell}) - \frac{1}{2} \left( \widehat{\psi}_{0}^{\ell} - \widehat{\psi}^{\ell} \right)^{\mathsf{T}} I^{\ell}(\widecheck{\psi}^{\ell}) \left( \widehat{\psi}_{0}^{\ell} - \widehat{\psi}^{\ell} \right),$$

where  $I^{\ell}(\psi) = -\partial^2 L^{\ell}(\psi)/\partial\psi\partial\psi^{\dagger}$ . By the lower order variability of  $I^{\ell}(\check{\psi}^{\ell})$ , we can find  $\check{\psi}^*$  such that  $I^{\ell}(\check{\psi}^{\ell}) \cong I^{\ell}(\check{\psi}^*)$  for all  $\ell$ . Then, using similar techniques as in (B.2) and (B.3), we have

$$(B.4) L^{\ell}(\widehat{\psi}^{\ell}) - L^{\ell}(\widehat{\psi}^{\ell}_{0}) \cong \frac{1}{2} \left( \widehat{\psi}^{\ell}_{0} - \widehat{\psi}^{\ell} \right)^{\mathsf{T}} I^{\ell}(\widecheck{\psi}^{*}) \left( \widehat{\psi}^{\ell}_{0} - \widehat{\psi}^{\ell} \right)$$

$$\cong \frac{1}{2} \left( \theta_{0} - \widehat{\theta}^{\ell} \right)^{\mathsf{T}} \widecheck{U}^{-1} \left( \theta_{0} - \widehat{\theta}^{\ell} \right)$$

for some matrix  $\check{U}$ . Similarly, we have

(B.5) 
$$L^{\ell}(\widehat{\psi}^*) - L^{\ell}(\widehat{\psi}_0^*) \simeq \frac{1}{2} \left( \theta_0 - \widehat{\theta}^* \right)^{\mathsf{T}} \widecheck{U}^{-1} \left( \theta_0 - \widehat{\theta}^* \right).$$

Write  $A^{\otimes 2} = AA^{\dagger}$  for any appropriate matrix A. Using (B.4), (B.5) and the cyclic property of trace, we have

$$\overline{d}_{L} - \widehat{d}_{L} \simeq \frac{1}{m} \sum_{\ell=1}^{m} \left\{ \left( \theta_{0} - \widehat{\theta}^{\ell} \right)^{\mathsf{T}} \widecheck{U}^{-1} \left( \theta_{0} - \widehat{\theta}^{\ell} \right) - \left( \theta_{0} - \widehat{\theta}^{*} \right)^{\mathsf{T}} \widecheck{U}^{-1} \left( \theta_{0} - \widehat{\theta}^{*} \right) \right\}$$

$$= \operatorname{tr} \left[ \widecheck{U}^{-1} \left\{ \frac{1}{m} \sum_{\ell=1}^{m} \left( \theta_{0} - \widehat{\theta}^{\ell} \right)^{\otimes 2} - \left( \theta_{0} - \widehat{\theta}^{*} \right)^{\otimes 2} \right\} \right]$$

$$\simeq \operatorname{tr} \left[ \widecheck{U}^{-1} \frac{1}{m} \sum_{\ell=1}^{m} \left\{ (\widehat{\theta}^{\ell})^{\otimes 2} - \overline{\theta}^{\otimes 2} \right\} \right] \simeq \operatorname{tr} \left( \widecheck{U}^{-1} B \right) \simeq \operatorname{tr} \left( \mathcal{U}_{\theta, 0}^{-1} \mathscr{B}_{\theta} \right)$$

as  $m, n \to \infty$ , where  $\mathcal{U}_{\theta,0}$  is a deterministic matrix that depends on both  $\theta_0$  and  $\theta^*$ , and satisfies  $n(\check{U} - \mathcal{U}_{\theta,0}) \stackrel{\mathrm{pr}}{\to} 0$ . Note that  $\mathrm{tr}(\mathcal{U}_{\theta,0}^{-1}\mathcal{B}_{\theta}) = k_{\ell} r_0$ , for some finite  $r_0$  by Assumption 2. Then  $\hat{r}_L \stackrel{\mathrm{pr}}{\to} r_0 = \mathrm{tr}(\mathcal{U}_{\theta,0}^{-1}\mathcal{B}_{\theta})/k$ , proving (ii). (But  $\mathcal{U}_{\theta,0}$  may not equal to  $\mathcal{U}_{\theta}$ , and hence  $\hat{r}_L$  may not be consistent for  $r_m$ .)

If  $H_0$  is true, then  $\overline{\theta} \stackrel{\text{pr}}{\to} \theta_0 = \theta^*$  and  $\widecheck{U} \simeq \overline{U} \simeq \mathcal{U}_{\theta} = \mathcal{U}_{\theta,0}$ . Then,  $\widehat{r}_L \stackrel{\text{pr}}{\to} r$  as  $m, n \to \infty$ . So, (i) follows.

PROOF OF THEOREM 2.4. (i, ii) It is trivial by the definition of  $\hat{r}_{\rm L}^{\Diamond}$ . (iii) Applying Taylor's expansion to  $\psi \mapsto L^{\ell}(\psi)$  again, we know there is  $\check{\psi}^{\ell}$  lying

on the line segment joining  $\widehat{\psi}^{\ell}$  and  $\widehat{\psi}^*$  such that

(B.6) 
$$L^{\ell}(\widehat{\psi}^*) = L^{\ell}(\widehat{\psi}^{\ell}) - \frac{1}{2} \left( \widehat{\psi}^* - \widehat{\psi}^{\ell} \right)^{\mathsf{T}} I^{\ell}(\widecheck{\psi}^{\ell}) \left( \widehat{\psi}^* - \widehat{\psi}^{\ell} \right).$$

By the lower order variability of  $I^{\ell}(\check{\psi}^{\ell})$ , we know that  $I^{\ell}(\check{\psi}^{\ell}) \cong \overline{I}(\hat{\psi}^*)$  for all  $\ell$ , where  $\overline{I}(\psi) = m^{-1} \sum_{\ell=1}^m I^{\ell}(\psi)$ . We also know that  $\hat{\psi}^* \cong \overline{\psi}$ . Thus

$$\overline{\delta}_{L} - \widehat{\delta}_{L} \simeq \frac{1}{m} \sum_{\ell=1}^{m} \left( \widehat{\psi}^{*} - \widehat{\psi}^{\ell} \right)^{\mathsf{T}} \overline{I}(\widehat{\psi}^{*}) \left( \widehat{\psi}^{*} - \widehat{\psi}^{\ell} \right) \\
= \operatorname{tr} \left\{ \overline{I}(\widehat{\psi}^{*}) \frac{1}{m} \sum_{\ell=1}^{m} \left( \widehat{\psi}^{*} - \widehat{\psi}^{\ell} \right)^{\otimes 2} \right\} \\
\simeq \operatorname{tr} \left\{ \overline{I}(\widehat{\psi}^{*}) \frac{1}{m} \sum_{\ell=1}^{m} \left( \widehat{\psi}^{\ell} - \overline{\psi} \right)^{\otimes 2} \right\} \simeq \operatorname{tr} \left( \mathcal{U}_{\psi}^{-1} \mathcal{B}_{\psi} \right)$$
(B.7)

as  $m, n \to \infty$ . By the assumption of EFMI of  $\psi$ , we have  $\widehat{r}_L^{\lozenge} \xrightarrow{\operatorname{pr}} r$ .

PROOF OF LEMMA 2.5. First, recall that, as  $n \to \infty$ , the observed data MLE  $\hat{\theta}_{\text{obs}}$  of  $\theta$  satisfies (2.4), which can be written as  $[\hat{\theta}_{\text{obs}} \mid \theta^{\star}] \stackrel{\mathfrak{D}}{\approx} \mathcal{N}_k(\theta^{\star}, \mathcal{T}_{\theta})$ , where  $A_{1,n} \stackrel{\mathfrak{D}}{\approx} A_{2,n}$  means that  $A_{1,n}$  and  $A_{2,n}$  have the same asymptotic distribution, i.e., there exist deterministic sequences  $\mu_n$  and  $\Sigma_n$  such that  $(A_{1,n} - \mu_n)\Sigma_n^{-1/2} \Rightarrow A$  and  $(A_{2,n} - \mu_n)\Sigma_n^{-1/2} \Rightarrow A$  for some non-degenerate random variable A. From Assumption 3, a proper imputation model is used. So, we have (2.5), which is equivalent to say that, as  $n \to \infty$ ,

(B.8) 
$$\left[\widehat{\theta}^{\ell} \mid X_{\text{obs}}\right] \stackrel{\mathcal{D}}{\approx} \mathcal{N}_{k}(\widehat{\theta}_{\text{obs}}, \mathscr{B}_{\theta}),$$

independently for for  $\ell = 1, ..., m$ . Therefore we can represent

(B.9) 
$$\widehat{\theta}_{\text{obs}} \stackrel{\mathfrak{D}}{\approx} \theta^{\star} + \mathcal{T}_{\theta}^{1/2} W,$$

(B.10) 
$$\widehat{\theta}^{\ell} \stackrel{\mathcal{D}}{\approx} \widehat{\theta}_{\text{obs}} + \mathcal{B}_{\theta}^{1/2} Z_{\ell}, \qquad \ell = 1, \dots, m$$

where  $Z_1, \ldots, Z_m, W \stackrel{\text{iid}}{\sim} \mathcal{N}_k(0, I_k)$ . Also write  $Z_\ell = (Z_{1\ell}, \ldots, Z_{k\ell})^{\intercal}$ , for  $\ell = 1, 2, \ldots, m$ , and  $W = (W_1, \ldots, W_k)^{\intercal}$ . Averaging (B.10) over  $\ell$ , we have  $\overline{\theta} \stackrel{\mathcal{D}}{\approx} \widehat{\theta}_{\text{obs}} + \mathcal{B}_{\theta}^{1/2} \overline{Z}_{\bullet}$ , where  $\overline{Z}_{\bullet} = m^{-1} \sum_{\ell=1}^{m} Z_{\ell}$ . Since  $\mathcal{B}_{\theta} = \mathcal{P} \mathcal{U}_{\theta}$ , we have

$$\mathcal{U}_{\theta}^{-1/2}(\widehat{\theta}^{\ell} - \theta^{\star}) \stackrel{\mathcal{D}}{\approx} (1 + r)^{1/2}W + r^{1/2}Z_{\ell},$$

$$\mathcal{U}_{\theta}^{-1/2}(\overline{\theta} - \theta^{\star}) \stackrel{\mathcal{D}}{\approx} (1 + r)^{1/2}W + r^{1/2}\overline{Z}_{\bullet}.$$

Note that (2.6) implies  $\mathcal{U}_{\theta} \cong \overline{U}$ . Under  $H_0$ , we have  $\theta^* = \theta_0$  and

$$\begin{split} \overline{d}_{\mathrm{L}} & \simeq & \overline{d}'_{\mathrm{W}} \stackrel{\mathcal{D}}{\approx} \sum_{i=1}^{k} \left\{ (1+r)^{1/2} W_i + r^{1/2} Z_{i\ell} \right\}^2, \\ \\ \widehat{d}_{\mathrm{L}} & \simeq & \widetilde{d}_{\mathrm{L}} \simeq & \widetilde{d}'_{\mathrm{W}} \stackrel{\mathcal{D}}{\approx} \sum_{i=1}^{k} \left\{ (1+r)^{1/2} W_i + r^{1/2} \overline{Z}_i \right\}^2. \end{split}$$

After some simple algebra, we obtain

$$\widehat{r}_{\mathrm{L}}^{+} \stackrel{\mathcal{D}}{\approx} \frac{(m+1)\boldsymbol{r}}{mk} \sum_{i=1}^{k} s_{Z_{i}}^{2} \quad \text{and} \quad \widehat{D}_{\mathrm{L}}^{+} \stackrel{\mathcal{D}}{\approx} \frac{m \sum_{i=1}^{k} \left\{ (1+\boldsymbol{r})^{1/2} W_{i} + \boldsymbol{r}^{1/2} \overline{Z}_{i\bullet} \right\}^{2}}{mk + (m+1)\boldsymbol{r} \sum_{i=1}^{k} s_{Z_{i}}^{2}},$$

where  $s_{Z_i}^2 = (m-1)^{-1} \sum_{\ell=1}^m (Z_{i\ell} - \overline{Z}_{i\bullet})^2$  is the sample variance of  $\{Z_{i\ell}\}_{\ell=1}^m$ . Since  $W_i$ ,  $\overline{Z}_{i\bullet}$  and  $s_{Z_i}^2$  are mutually independent for each fixed i, we can simplify the representation of  $\widehat{D}_{L}^+$  to

$$\widehat{r}_{\mathrm{L}}^{+} \overset{\mathfrak{D}}{\approx} \frac{(m+1) \cancel{r}}{m(m-1) k} \sum_{i=1}^{k} H_{i}^{2} \quad \text{and} \quad \widehat{D}_{\mathrm{L}}^{+} \overset{\mathfrak{D}}{\approx} \frac{(m-1) \{m+(m+1) \cancel{r}\} \sum_{i=1}^{k} G_{i}^{2}}{m(m-1) k + (m+1) \cancel{r} \sum_{i=1}^{k} H_{i}^{2}},$$

where  $G_i^2 \stackrel{\text{iid}}{\sim} \chi_1^2$  and  $H_i^2 \stackrel{\text{iid}}{\sim} \chi_{m-1}^2$ , for  $i = 1, \dots, k$ , are all mutually independent. Clearly, they can be further simplified to (2.14).

PROOF OF THEOREM 2.6. Similar to (B.9) and (B.10), we can have a more general representation:

$$\hat{\psi}_{\text{obs}} \stackrel{\mathcal{D}}{\approx} \psi^{\star} + \mathcal{T}_{\psi}^{1/2}W; \quad \hat{\psi}^{\ell} \stackrel{\mathcal{D}}{\approx} \hat{\psi}_{\text{obs}} + \mathcal{B}_{\psi}^{1/2}Z_{\ell}, \qquad \ell = 1, \dots, m,$$

where  $Z_1, \ldots, Z_h, W \stackrel{\text{iid}}{\sim} \mathcal{N}_h(0, I_h)$ . Also write  $Z_\ell = (Z_{1\ell}, \ldots, Z_{h\ell})^{\intercal}$ , for  $\ell = 1, 2, \ldots, m$ , and  $W = (W_1, \ldots, W_h)^{\intercal}$ . Using (B.7), we have

$$\overline{\delta}_{L} - \widehat{\delta}_{L} = \operatorname{tr} \left\{ \overline{I}(\widehat{\psi}^{*}) \frac{1}{m} \sum_{\ell=1}^{m} \left( \widehat{\psi}^{\ell} - \overline{\psi} \right) \left( \widehat{\psi}^{\ell} - \overline{\psi} \right)^{\mathsf{T}} \right\}$$

$$\overset{\mathcal{D}}{\approx} \operatorname{tr} \left[ \mathcal{U}_{\psi}^{-1} \frac{1}{m} \sum_{\ell=1}^{m} \left\{ (\mathcal{T}_{\psi} - \mathcal{U}_{\psi})^{1/2} \left( Z_{\ell} - \overline{Z}_{\bullet} \right) \right\}^{\otimes 2} \right]$$

$$= \frac{1}{m} \sum_{\ell=1}^{m} \operatorname{tr} \left\{ r I_{h} \left( Z_{\ell} - \overline{Z}_{\bullet} \right)^{\otimes 2} \right\} = \frac{r}{m} \sum_{\ell=1}^{m} \sum_{i=1}^{h} (Z_{i\ell} - \overline{Z}_{i\bullet})^{2}.$$

Equivalently, we can say  $\overline{\delta}_{L} - \hat{\delta}_{L} \Rightarrow r \chi^{2}_{h(m-1)}/m$  as  $n \to \infty$ . Hence

$$\hat{r}_{\rm L}^{\Diamond} \Rightarrow r \cdot \frac{m+1}{hm(m-1)} \cdot \chi_{h(m-1)}^2,$$

which is equivalent to (2.15). Note that it is true under both  $H_0$  and  $H_1$ .  $\square$ 

PROOF OF THEOREM 2.7. From the representations of  $\widehat{d}_{L}^{\Diamond}$  and  $\widehat{r}_{L}^{\Diamond}$  in Lemma 2.5 and Theorem 2.6, we know that they are asymptotically  $(n \to \infty)$  independent. The proof then follows the derivation for Lemma 2.5.

PROOF OF PROPOSITION 3.1. It is trivial.

PROOF OF THEOREM A.1. (i) Using the representation (A.1), we can easily see that  $\hat{r}_L^{\triangle} \ge 0$ . (ii) It suffices to show

$$m^{-1} \sum_{\ell=1}^{m} d_{\mathcal{L}}(\widehat{\psi}_{0}^{\ell} + \Delta_{m}, \widehat{\psi}^{\ell} \mid X^{\ell}) \simeq \overline{d}_{\mathcal{L}} - \widetilde{d}_{\mathcal{L}},$$

where  $\Delta_m = \hat{\psi}^* - \hat{\psi}_0^*$ . Under  $H_0$ ,  $\Delta_m \simeq 0$  and  $\hat{\psi}_0^{\ell} \simeq \hat{\psi}^{\ell}$ , so  $\hat{\psi}_0^{\ell} + \Delta_m \simeq \hat{\psi}^{\ell}$ . Using Taylor's expansion on  $\psi \mapsto L^{\ell}(\psi)$  around its maximizer  $\hat{\psi}^{\ell}$ , we have for  $\psi \simeq \hat{\psi}^{\ell}$  that

$$L^{\ell}(\psi) \simeq L^{\ell}(\widehat{\psi}^{\ell}) - \frac{1}{2} \left( \psi - \widehat{\psi}^{\ell} \right)^{\mathsf{T}} I^{\ell}(\widehat{\psi}^{\ell}) \left( \psi - \widehat{\psi}^{\ell} \right).$$

Under the parametrization of  $\psi$  in the proof of Theorem 2.1, we know that the upper  $k \times k$  sub-matrix of  $I^{\ell}(\widehat{\psi}^{\ell})$  is  $(U^{\ell})^{-1}$ . Using the lower order variability of  $U^{\ell}$ , we have  $(U^{\ell})^{-1} \simeq \overline{U}^{-1}$  and

$$\frac{1}{m} \sum_{\ell=1}^{m} d_{L}(\hat{\psi}_{0}^{\ell} + \Delta_{m}, \hat{\psi}^{\ell} \mid X^{\ell})$$

$$\stackrel{=}{=} \frac{1}{m} \sum_{\ell=1}^{m} \left( \hat{\psi}_{0}^{\ell} + \Delta_{m} - \hat{\psi}^{\ell} \right)^{\mathsf{T}} I^{\ell}(\hat{\psi}^{\ell}) \left( \hat{\psi}_{0}^{\ell} + \Delta_{m} - \hat{\psi}^{\ell} \right)$$

$$\stackrel{=}{=} \frac{1}{m} \sum_{\ell=1}^{m} (\hat{\theta}^{\ell} - \overline{\theta})^{\mathsf{T}} \overline{U}^{-1} (\hat{\theta}^{\ell} - \overline{\theta}) = \overline{d}'_{W} - \widetilde{d}'_{W} \stackrel{=}{=} \overline{d}_{L} - \widehat{d}_{L}.$$

Therefore, the desired result follows.

PROOF OF THEOREM A.2. Throughout this proof, conditions (a), (b) and (c) refer to the list given in Assumption 5. (i, ii) It trivially follows

from the definitions of  $\hat{d}_S$  and  $\hat{r}_S$ . (iii) First, by the definition of maximizer and condition (a), we have

$$\begin{split} \underline{\overline{L}}(\hat{\psi}^*) - \underline{\overline{L}}(\hat{\psi}^{\mathrm{S}}) &= \underline{\overline{L}}(\hat{\psi}^*) - \underline{\overline{L}}^{\mathrm{S}}(\hat{\psi}^{\mathrm{S}}) + \underline{\overline{L}}^{\mathrm{S}}(\hat{\psi}^{\mathrm{S}}) - \underline{\overline{L}}(\hat{\psi}^{\mathrm{S}}) \\ &\leqslant \underline{\overline{L}}(\hat{\psi}^*) - \underline{\overline{L}}^{\mathrm{S}}(\hat{\psi}^*) + \underline{\overline{L}}^{\mathrm{S}}(\hat{\psi}^{\mathrm{S}}) - \underline{\overline{L}}(\hat{\psi}^{\mathrm{S}}) \\ &\leqslant 2\sup_{\psi \in \Psi} \left| \underline{\overline{L}}(\psi) - \underline{\overline{L}}^{\mathrm{S}}(\psi) \right| = O_p(1/n), \end{split}$$

which, together with condition (b), imply that

$$\underline{\underline{\mathscr{Z}}}(\psi^{*}) - \underline{\underline{\mathscr{Z}}}(\widehat{\psi}^{S}) 
= \{\underline{\underline{\mathscr{Z}}}(\psi^{*}) - \underline{\underline{L}}(\psi^{*})\} + \{\underline{\underline{L}}(\psi^{*}) - \underline{\underline{L}}(\widehat{\psi}^{S})\} + \{\underline{\underline{L}}(\widehat{\psi}^{S}) - \underline{\underline{\mathscr{Z}}}(\widehat{\psi}^{S})\} 
(B.11) 
$$\leqslant 2 \sup_{\psi \in \Psi} |\underline{\underline{L}}(\psi) - \underline{\underline{\mathscr{Z}}}(\psi)| + \{\underline{\underline{L}}(\widehat{\psi}^{*}) - \underline{\underline{L}}(\widehat{\psi}^{S})\} = o_{p}(1).$$$$

Using (B.11) and (c), we have  $\hat{\psi}^S \xrightarrow{pr} \psi^*$ . By (b) and (c), we also have  $\hat{\psi}^* \xrightarrow{pr} \psi^*$ . So,  $|\hat{\psi}^S - \hat{\psi}^*| \xrightarrow{pr} \mathbf{0}$  as  $n \to \infty$ . By the definition of maximizer,

(B.12) 
$$\mathbf{0} = \nabla \overline{L}^{S}(\hat{\psi}^{S}) = \nabla \overline{L}(\hat{\psi}^{S}) + \nabla R(\hat{\psi}^{S}).$$

where  $\nabla g(\psi) = \partial g(\psi)/\partial \psi$  is the gradient of  $\psi \mapsto g(\psi)$ . By condition (a), we know  $\nabla R(\widehat{\psi}^{S}) = O_{p}(1/n)$ . Thus, together with (B.12), we have  $\nabla \overline{\underline{L}}(\widehat{\psi}^{S}) = O_{p}(1/n)$ . Also, by the definition of MLE, we have  $\nabla \overline{\underline{L}}(\widehat{\psi}^{*}) = \mathbf{0}$ .

By Taylor's expansion, there exists  $\check{\psi}$  such that

(B.13) 
$$\underline{\overline{L}}(\widehat{\psi}^*) - \underline{\overline{L}}(\widehat{\psi}^S) = \left\{ \nabla \underline{\overline{L}}(\widecheck{\psi}) \right\}^{\mathsf{T}} \left( \widehat{\psi}^* - \widehat{\psi}^S \right) = o_p(1/n),$$

where we have used the continuity of  $\psi \mapsto \nabla \underline{\underline{L}}(\psi)$  to yield  $\nabla \underline{\underline{L}}(\check{\psi}) = O_p(1/n)$ . Rewriting (B.13), we have

(B.14) 
$$\underline{\overline{L}}(\widehat{\psi}^*) - \underline{\overline{L}}^{S}(\widehat{\psi}^{S}) = R(\widehat{\psi}^{S}) + o_p(1/n).$$

Similar to (B.14), we have

$$(B.15) \overline{\underline{L}}(\hat{\psi}_0^*) - \overline{\underline{L}}^S(\hat{\psi}_0^S) = R(\hat{\psi}_0^S) + o_p(1/n).$$

Then, using (B.14) and (B.15), we have

$$\begin{aligned} \left| \hat{d}_{\mathcal{L}} - \hat{d}_{\mathcal{S}} \right| &= 2n \left| \left\{ \underline{\underline{L}}(\hat{\psi}^*) - \underline{\underline{L}}^{\mathcal{S}}(\hat{\psi}^{\mathcal{S}}) \right\} - \left\{ \underline{\underline{L}}(\hat{\psi}^*_0) - \underline{\underline{L}}^{\mathcal{S}}(\hat{\psi}^{\mathcal{S}}_0) \right\} \right| \\ &= 2n \left| R(\hat{\psi}^{\mathcal{S}}) - R(\hat{\psi}^{\mathcal{S}}_0) + o_p(1/n) \right|. \end{aligned}$$

Now consider two cases.

- (i) Under  $H_0$ , we have  $\hat{d}_L = O_p(1)$  and  $\hat{\psi}_0^S \simeq \hat{\psi}^S$ . Thus condition (a) implies  $R(\hat{\psi}^{S}) - R(\hat{\psi}_{0}^{S}) = o_{p}(1/n)$ . Then, we have  $|\hat{d}_{L} - \hat{d}_{S}| = o_{p}(\hat{d}_{L})$ .
- (ii) Under  $H_1$ , we have  $\hat{d}_L \stackrel{\text{pr}}{\to} \infty$ . Condition (a) and (B.11) imply that  $\underline{\overline{L}}(\hat{\psi}^*) \underline{\overline{L}}^S(\hat{\psi}^S) = O_p(1/n)$ . Similarly, we also have  $\underline{\overline{L}}(\hat{\psi}^*) \underline{\overline{L}}^S(\hat{\psi}^S) = O_p(1/n)$ .  $O_p(1/n)$ . Hence  $\left|\hat{d}_L - \hat{d}_S\right| = O_p(1)$ . Thus we have  $\left|\hat{d}_L - \hat{d}_S\right| = o_p(\hat{d}_L)$ .

Therefore, under either  $H_0$  or  $H_1$ , we also have  $\left|\hat{d}_L - \hat{d}_S\right| = o_p(\hat{d}_L)$ . Since  $\hat{d}_{L} = \hat{d}_{S}$  and  $\bar{d}_{L} = \bar{d}_{S}$ , we know  $\hat{r}_{L} = \hat{r}_{S}$ .

Note that, even under the assumption of this theorem,  $\hat{r}_S$  and  $\hat{r}_S^{\Diamond}$  are not equivalent. From (3.5) and (3.6),  $\hat{r}_S$  and  $\hat{r}_S^{\Diamond}$  are a "difference" estimator and a "difference" estimator, respectively. So, the "bias" of using  $\overline{L}^{\rm S}(\psi)$  cannot be canceled out in  $\widehat{r}_{\rm S}^{\diamond}$ .

### APPENDIX C: ADDITIONAL FIGURES AND TABLES

This section presents additional figures and tables in  $\S 2.5$  and  $\S 4$ 

- Figure 9: the performance of different approximations to the reference null distribution when  $\alpha = 5\%$ ; see § 2.5.
- Figures 10 and 11: the empirical distributions of the p-values under  $H_0$  and parametrizations (i) and (iii), respectively; see § 4.1.
- Figures 12 and 13: the empirical power functions and the empirical ratio of power-to-size for size 5% tests, respectively; see § 4.1.
- Table 8: the ranges of empirical sizes over different parametrizations for size 5% tests; see § 4.1.
- Table 9: detailed results of the care-survival example in § 4.3.

Department of Statistics THE CHINESE UNIVERSITY OF HONG KONG SHATIN, N.T., HONG KONG

E-MAIL: kinwaichan@sta.cuhk.edu.hk

Department of Statistics HARVARD UNIVERSITY

Cambridge, Massachusetts 02138, U.S.A.

E-MAIL: meng@stat.harvard.edu

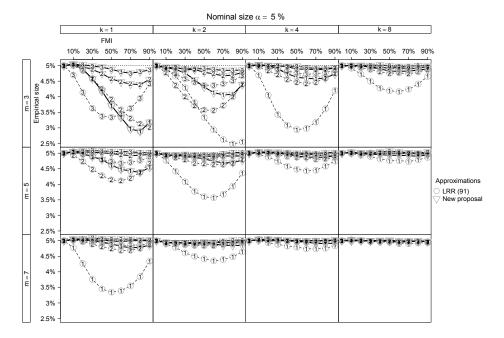


FIG 9. The performances of two different approximate null distributions when the nominal size is  $\alpha=5\%$ . The vertical axis denotes  $\hat{\alpha}$  or  $\tilde{\alpha}$ , and the horizontal axis denotes the value of  $f_m$ . The number attached to each line denotes the value of  $\tau=h/k$ . The proposed approximation  $\hat{\alpha}$  is denoted by thick solid lines with triangles, and the existing approximation  $\tilde{\alpha}$  is denoted by thin dashed lines with circles.

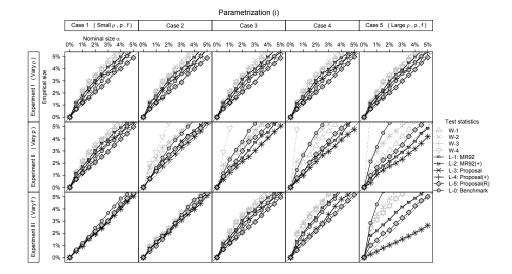


FIG 10. The comparison between empirical size and nominal size  $\alpha$  under parametrization (i) for  $\alpha \in (0,5\%]$ . The Wald tests (W-1,2,3,4) and LRTs (L-0,1,2,3,4,5) are represented by grey dashed and black solid lines, respectively. The LRT statistic  $\widetilde{D}_{\rm L}$  (L-1: MR92) (Meng and Rubin, 1992) and its modification  $\widetilde{D}_{\rm L}^+$  (L-2: MR92(+)) are the existing counterparts of our proposals  $\widehat{D}_{\rm S}$  (L-3: Proposal),  $\widehat{D}_{\rm S}^+$  (L-4: Proposal(+)) and  $\widehat{D}_{\rm S}^{\diamond}$  (L-5: Proposal(R)).

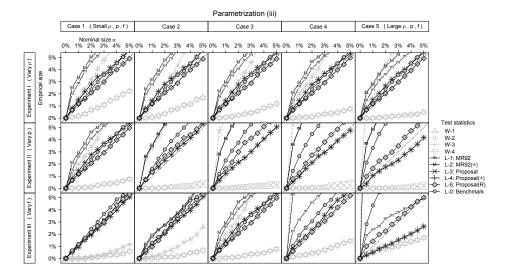


FIG 11. The comparison between empirical size and nominal size  $\alpha$  under parametrization (iii) for  $\alpha \in (0, 5\%]$ . The legend in Figure 10 also applies here.

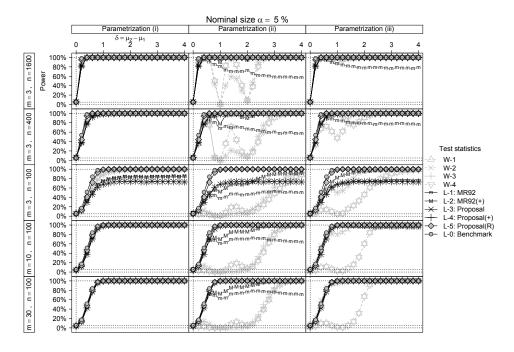


FIG 12. The power curves under different parametrizations. The nominal size is  $\alpha = 5\%$ . In each plot, the vertical axis denotes the power, whereas the horizontal axis denotes the value of  $\delta = \mu_2 - \mu_1$ . The legend in Figure 10 also applies here.

Table 8 The range of empirical size  $[\min \hat{\alpha}, \max \hat{\alpha}]$  in percentage, where  $\max$  and  $\min$  are taken over the three parametrizations. Only one value is recorded for those tests that are invariant to parametrization. The nominal size is  $\alpha = 5\%$ .

Range of empirical size: $[\min \hat{\alpha}, \max \hat{\alpha}]/\%$										
(n,m)	(1600, 3)	(400, 3)	(100, 3)	(100, 10)	(100, 30)					
W-1	[5.62, 5.71]	[5.30, 6.03]	[3.22, 6.20]	[1.64, 4.81]	[1.37, 5.00]					
W-2	[5.93, 6.05]	[6.08, 7.18]	[5.52, 8.69]	[4.42, 8.47]	[4.20, 8.50]					
W-3	[5.81, 6.03]	[6.01, 6.98]	[5.37, 8.28]	[4.20, 7.67]	[4.10, 7.50]					
W-4	[5.62, 5.71]	[5.30, 6.03]	[3.22, 6.20]	[1.64, 4.81]	[1.37, 5.00]					
L-1	[5.57, 6.15]	[6.37, 6.57]	[5.88, 6.47]	[4.39, 5.66]	[4.22, 5.32]					
L-2	[5.52, 6.10]	[6.37, 6.52]	[5.88, 7.47]	[4.39, 5.66]	[4.22, 5.32]					
L-3	5.76	6.37	5.42	3.78	3.71					
L-4	5.76	6.37	5.42	3.78	3.71					
L-5	4.96	5.32	4.93	4.79	4.54					
L-0	5.03	5.03	5.57	5.57	5.57					

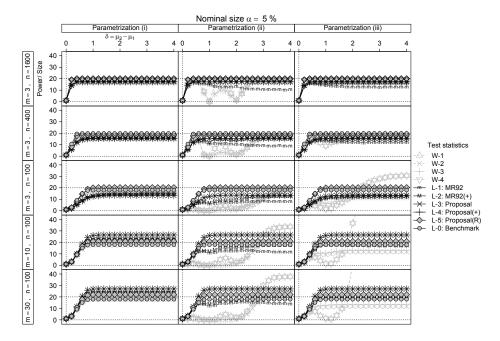


FIG 13. The ratios of empirical power to empirical size under different parametrizations. The nominal size is  $\alpha=5\%$ . In each plot, the vertical axis denotes the ratio, whereas the horizontal axis denotes  $\delta=\mu_2-\mu_1$ . The legend in Figure 10 also applies here.

Table 9
The LRTs using  $\widetilde{D}_L$ ,  $\widehat{D}_S^+$  and  $\widehat{D}_S^{\Diamond}$  under different parametrizations in § 4.3.

	Parametrization (i)											
	<i>H</i> <sub>0</sub> : Co	nditional indepen	idence	$H_0$ : Full in	dependence							
m	$\widetilde{r}_{\mathrm{L}}, \widehat{r}_{\mathrm{S}}^{+}, \widehat{r}_{\mathrm{S}}^{\Diamond}$	$\tilde{D}_{\mathrm{L}}, \hat{D}_{\mathrm{S}}^{+}, \hat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{ m L}, \widehat{p}_{ m S}^+, \widehat{p}_{ m S}^{\diamondsuit}$	$\widetilde{r}_{ m L}, \widehat{r}_{ m S}^+, \widehat{r}_{ m S}^{\Diamond}$	$\tilde{D}_{\mathrm{L}}, \hat{D}_{\mathrm{S}}^{+}, \hat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{ m L},\widehat{p}_{ m S}^+,\widehat{p}_{ m S}^{\diamondsuit}$						
2	0.63, 0.64, 0.83	0.14, 0.14, 0.12	0.87, 0.87, 0.89	0.53, 0.53, 0.83	44.4, 44.4, 37.1	0,0,0						
3	0.54, 0.54, 0.38	0.08, 0.08, 0.09	0.93, 0.93, 0.92	0.31, 0.31, 0.38	54.2, 54.2, 51.4	0, 0, 0						
5	0.49, 0.48, 0.89	0.12, 0.12, 0.10	0.89, 0.89, 0.91	0.72, 0.72, 0.89	40.8, 40.8, 37.1	0, 0, 0						
7	0.23, 0.23, 0.47	0.06, 0.06, 0.05	0.94, 0.94, 0.95	0.31, 0.31, 0.47	53.2, 53.2, 47.6	0, 0, 0						
10	0.50, 0.50, 0.70	0.14, 0.14, 0.12	0.87, 0.87, 0.88	0.56, 0.56, 0.70	45.4, 45.4, 41.7	0, 0, 0						
25	0.35, 0.35, 0.47	0.06, 0.06, 0.06	0.94, 0.94, 0.95	0.35, 0.35, 0.47	51.4, 51.4, 47.0	0, 0, 0						
50	0.31, 0.31, 0.45	0.11, 0.11, 0.10	0.90, 0.90, 0.91	0.33, 0.33, 0.45	51.5, 51.5, 47.3	0, 0, 0						
	Parametrization (ii)											
	$H_0$ : Co	nditional indepen	dence	$H_0$ : Full in	dependence							
m	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\lozenge}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{ m L}, \widehat{p}_{ m S}^+, \widehat{p}_{ m S}^{\diamondsuit}$	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\lozenge}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{\mathrm{L}},\widehat{p}_{\mathrm{S}}^{+},\widehat{p}_{\mathrm{S}}^{\Diamond}$						
2	1.23, 0.64, 0.83	0.01, 0.14, 0.12	0.99, 0.87, 0.89	0.98, 0.53, 0.83	34.2, 44.4, 37.1	0, 0, 0						
3	1.08, 0.54, 0.38	-0.07, 0.08, 0.09	1.00, 0.93, 0.92	0.61, 0.31, 0.38	43.9, 54.2, 51.4	0, 0, 0						
5	1.02, 0.48, 0.89	-0.09, 0.12, 0.10	1.00, 0.89, 0.91	1.40, 0.72, 0.89	29.0, 40.8, 37.1	0, 0, 0						
7	0.45, 0.23, 0.47	-0.07, 0.06, 0.05	1.00, 0.94, 0.95	0.58, 0.31, 0.47	43.9, 53.2, 47.6	0, 0, 0						
10	0.99, 0.50, 0.70	-0.10, 0.14, 0.12	1.00, 0.87, 0.88	1.09, 0.56, 0.70	33.7, 45.4, 41.7	0, 0, 0						
25	0.71, 0.35, 0.47	-0.14, 0.06, 0.06	1.00, 0.94, 0.95	0.68, 0.35, 0.47	41.0, 51.4, 47.0	0, 0, 0						
50	0.63, 0.31, 0.45	-0.10, 0.11, 0.10	1.00, 0.90, 0.91	0.65, 0.33, 0.45	41.3, 51.5, 47.3	0, 0, 0						
		Par	ametrization	(iii)								
	$H_0$ : Co	nditional indepen	dence	$H_0$ : Full in	dependence							
m	$\widetilde{r}_{ m L},\widehat{r}_{ m S}^+,\widehat{r}_{ m S}^{\lozenge}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\widetilde{p}_{ m L}, \widehat{p}_{ m S}^+, \widehat{p}_{ m S}^{\diamondsuit}$	$\widetilde{r}_{ m L}, \widehat{r}_{ m S}^+, \widehat{r}_{ m S}^{\lozenge}$	$\widetilde{D}_{\mathrm{L}}, \widehat{D}_{\mathrm{S}}^{+}, \widehat{D}_{\mathrm{S}}^{\Diamond}$	$\underbrace{\widetilde{p}_{\mathrm{L}},\widehat{p}_{\mathrm{S}}^{+},\widehat{p}_{\mathrm{S}}^{\Diamond}}_{\mathrm{L}}$						
2	1.06, 0.64, 0.83	0.04, 0.14, 0.12	0.96, 0.87, 0.88	-0.38, 0.53, 0.83	109, 44.4, 37.1	0, 0, 0						
3	-2.35, 0.54, 0.38	-1.16, 0.08, 0.09	1.00, 0.93, 0.92	-1.22, 0.31, 0.38	-321, 54.2, 51.4	1, 0, 0						
5	-2.64, 0.48, 0.89	-1.38, 0.12, 0.10	1.00, 0.89, 0.91	-2.24, 0.72, 0.89	-58, 40.8, 37.1	1, 0, 0						
7	-0.01, 0.23, 0.47	0.25, 0.06, 0.05	0.78, 0.94, 0.95	-0.34, 0.31, 0.47	107, 53.2, 47.6	0, 0, 0						
10	-2.04, 0.50, 0.70	-2.20, 0.14, 0.12	1.00, 0.87, 0.88	-1.85, 0.56, 0.70	-86, 45.4, 41.7	1, 0, 0						
25	-1.39, 0.35, 0.47	-4.30, 0.06, 0.06	1.00, 0.94, 0.95	-1.12, 0.35, 0.47	-603, 51.4, 47.0	1, 0, 0						
50	-1.22, 0.31, 0.45	-7.39, 0.11, 0.10	1.00, 0.90, 0.91	-1.06, 0.33, 0.45	-1136, 51.5, 47.3	1,0,0						