

MULTIPLE IMPROVEMENTS OF MULTIPLE IMPUTATION LIKELIHOOD RATIO TESTS

Kin Wai Chan¹ and Xiao-Li Meng²

Department of Statistics, The Chinese University of Hong Kong¹

Department of Statistics, Harvard University²

Supplementary Material

A Supplementary Results

A.1 A Complication Caused by Nuisance Parameter

This section supplement the discussion of Section 2.2 in the main article. Recall that the likelihood function $L^{(\ell)}(\cdot)$ is based on both observed data X_{obs} and imputed data $X_{\text{mis}}^{(\ell)}$, which varies across ℓ . Hence, each imputed likelihood $L^{(\ell)}(\cdot)$ is associated with a (imputation-specific) pseudo parameter $\psi^{(\ell)}$, may vary across $\ell = 1, \dots, m$.

To see the source of the negativity of $\hat{\tau}_L$, we extend $\bar{L}(\psi)$ in (2.1) to

$$\bar{L}(\psi^{(1)}, \dots, \psi^{(m)}) = \frac{1}{m} \sum_{\ell=1}^m L^{(\ell)}(\psi^{(\ell)}). \quad (\text{A.1})$$

Using the “log-likelihood” $\bar{L}(\psi^{(1)}, \dots, \psi^{(m)})$, we can construct, at least conceptually, four hypotheses $H_0^0, H_0^1, H_1^0, H_1^1$ defined in Table A.1. Each of them consists of zero, one or two of the constraints $\mathcal{E}_0 : \theta^{(1)} = \dots = \theta^{(m)} = \theta_0$ and $\mathcal{E}^0 : \psi^{(1)} = \dots = \psi^{(m)}$, where $\theta^{(\ell)} = \theta(\psi^{(\ell)})$ is the interested part of $\psi^{(\ell)}$ for each ℓ . The constraint \mathcal{E}_0 is equivalent to H_0 , and the constraint \mathcal{E}^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 A.2. Define the maximized value of $\bar{L}(\psi^{(1)}, \dots, \psi^{(m)})$

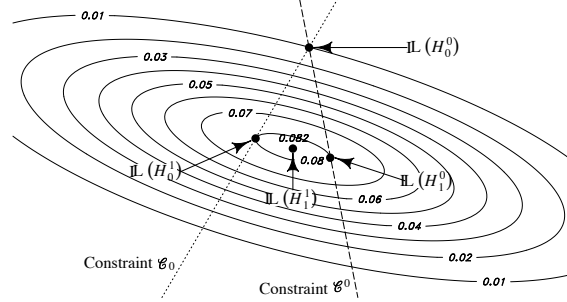


Figure A.1: A schematic illustration of the sign of (A.2). The contour lines of $\bar{L}(\psi^{(1)}, \dots, \psi^{(m)})$ are plotted. The two straight lines refer to constraints \mathcal{E}_0 and \mathcal{E}^0 . Since $\mathbb{L}(H_1^1) = 0.082$, $\mathbb{L}(H_0^0) = \mathbb{L}(H_1^0) = 0.08$, and $\mathbb{L}(H_0^0)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

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

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

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 $(\bar{d}_L - \hat{d}_L)/2$ as

$$(\bar{d}_L - \hat{d}_L)/2 = \{\mathbb{L}(H_1^1) - \mathbb{L}(H_1^0)\} - \{\mathbb{L}(H_0^1) - \mathbb{L}(H_0^0)\}. \quad (\text{A.2})$$

Whereas the two bracketed terms in (A.2) are non-negative as they correspond to two LRT statistics, their difference can be negative.

A simple example illustrates this well. For the regression model $[Y | X_1, X_2] \sim \mathcal{N}(\beta_0 + \beta_1 X_1 + \beta_2 X_2, \sigma^2)$, 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 $H_0^1 : \beta_1 \in \mathbb{R}, \beta_2 = 0$; see Figure A.1 for a schematic illustration.

The decomposition (A.2) 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 $\boldsymbol{r} = 0$ in

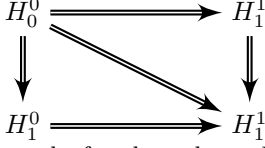


Figure A.2: The relationships between the four hypotheses $H_0^0, H_1^0, H_0^1, H_1^1$. Each arrow denotes an implication, e.g., $H_0^0 \Rightarrow H_1^0$ means that H_0^0 implies H_1^0 .

both θ and η , whereas $\mathbb{L}(H_0^1) - \mathbb{L}(H_0^0)$ seeks evidence only in η . For cases where θ and η are orthogonal (at least locally), the left-hand side of (A.2) can be viewed as a measure of evidence against $\nu = 0$ solely from θ ; Proposition 1 already hinted this possibility. However, the “test statistic” (A.2) has a fatal flaw. Because \mathcal{E}_0 requires all $\theta^{(\ell)}$ s to coincide with a specific θ_0 , \mathcal{E}_0 is not nested within \mathcal{E}^0 , i.e., $\mathcal{E}^0 \not\Rightarrow \mathcal{E}_0$. Hence \hat{r}_L is guaranteed to consistently estimate ν_m only under H_0 . This explains Corollary 1, and leads to an improvement in Section 2.2. In it not hard to see that our new estimator \hat{r}_L^\diamond simply drops the second term in (A.2).

A.2 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 ν_m is to perturb $\hat{\psi}_0^*$ and $\hat{\psi}_0^{(\ell)}$ by a suitable amount, say Δ , 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 Δ is $\Delta = \hat{\psi}^* - \hat{\psi}_0^*$. Using the perturbed version of \hat{r}_L , we obtain

$$\hat{r}_L^\Delta = \frac{m+1}{k(m-1)} \hat{\delta}_L^\Delta,$$

where

$$\hat{\delta}_L^\Delta = \frac{2}{m} \sum_{\ell=1}^m \log \left\{ \frac{f(X^{(\ell)} | \hat{\psi}^{(\ell)}) f(X^{(\ell)} | \hat{\psi}_0^* + \Delta)}{f(X^{(\ell)} | \hat{\psi}^*) f(X^{(\ell)} | \hat{\psi}_0^{(\ell)} + \Delta)} \right\} = \frac{1}{m} \sum_{\ell=1}^m d_L(\hat{\psi}_0^{(\ell)} + \Delta, \hat{\psi}^{(\ell)} | X^{(\ell)}).$$

Then we have the following result.

Theorem A.1. *Suppose RC_θ . Under H_0 , we have (i) $\hat{r}_L^\Delta \geq 0$ for all m, n ; and (ii) $\hat{r}_L^\Delta \simeq \hat{r}_L$ as $n \rightarrow \infty$ for each m .*

Although $\widehat{r}_L^\Delta \geq 0$, it is only invariant to affine transformations, and not robust against θ_0 , and less computational feasible than \widehat{r}_L ; see Section 3. However, it gives us some insights on how to construct a potentially better estimator. Note that, in (A.3), the constrained MLE is not used in $d_L(\cdot, \cdot | X^{(\ell)})$, but it is still always non-negative. We call this a “pseudo” LRT statistics. Then, $\widehat{\delta}_L^\Delta$ 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 estimator \widehat{r}_L^\diamond also takes the same form:

$$\widehat{r}_L^\diamond = \frac{m+1}{h(m-1)} \frac{1}{m} \sum_{\ell=1}^m d_L(\widehat{\psi}^*, \widehat{\psi}^{(\ell)} | X^{(\ell)}).$$

A.3 Additional result for Section 2.3

This section presents the additional simulation result for Section 2.3. The performance of different approximations to the reference null distribution when $\alpha = 5\%$ is shown in Figure A.3.

A.4 Results for Dependent Data

This is a supplement for Section 3.1. If the data are not independent, then (3.1) is no longer true. In other words, $\overline{L}(\psi) \neq \overline{L}^S(\psi)$, where $\overline{L}(\psi) = \sum_{\ell=1}^m L^{(\ell)}(\psi)/m$ is defined in (2.1), and

$$\overline{L}^S(\psi) = \frac{1}{m} \log f_{mn}(X^{(1:m)} | \psi). \quad (\text{A.3})$$

In principle, $\overline{L}(\psi)$ should be used instead of the “stacked version” $\overline{L}^S(\psi)$, however, the stacked one is much easier to compute. Because of this reason, it is of interest to see whether the stacked version can be used generally.

To begin with, we define the stacked version of all MI statistics when $\overline{L}^S(\psi)$ is used instead of $\overline{L}(\psi)$. Let

$$\widehat{\psi}_0^S = \arg \max_{\psi \in \Psi : \theta(\psi) = \theta_0} \overline{L}^S(\psi), \quad \widehat{\psi}^S = \arg \max_{\psi \in \Psi} \overline{L}^S(\psi); \quad (\text{A.4})$$

$$\widehat{\delta}_{0,S} = 2\overline{L}^S(\widehat{\psi}_0^S), \quad \widehat{\delta}_S = 2\overline{L}^S(\widehat{\psi}^S). \quad (\text{A.5})$$

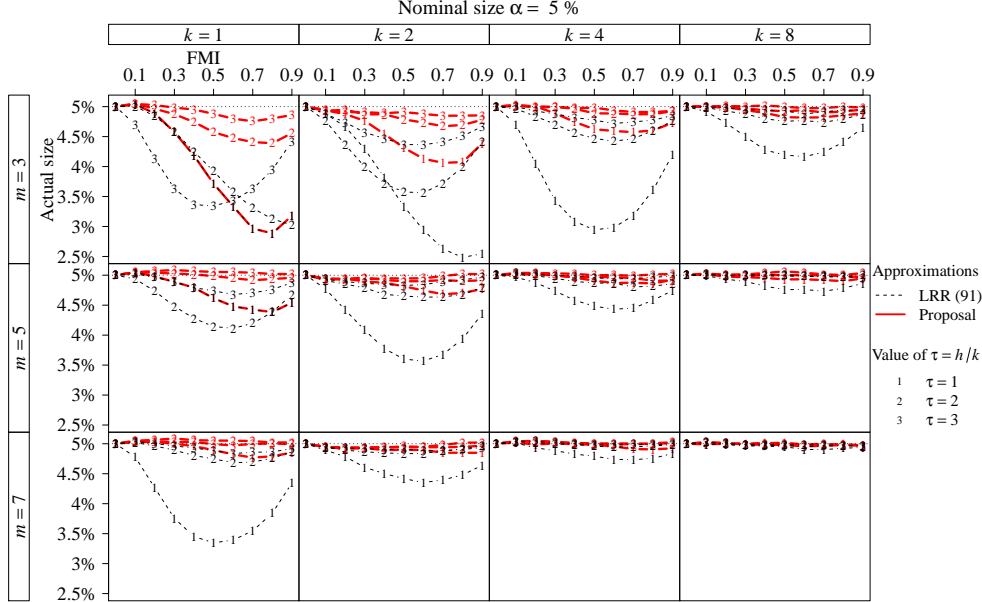


Figure A.3: The performance of two approximate null distributions when the nominal size is $\alpha = 5\%$. The vertical axis denotes $\hat{\alpha}$ or $\tilde{\alpha}$, and the horizontal axis denotes the value of f_m . The number attached to each line denotes the value of $\tau = h/k$. The proposed 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.

and

$$\hat{D}_S(\hat{r}_m) = \frac{\hat{d}_S}{k(1 + \hat{r}_m)}, \quad \text{with } \hat{d}_S = \hat{\delta}_S - \hat{\delta}_{0,S} \text{ of (A.5);} \quad (\text{A.6})$$

$$\hat{r}_S = \frac{m+1}{k(m-1)}(\bar{d}_S - \hat{d}_S), \quad \text{with } \bar{d}_S = \bar{d}_L \text{ of (1.7);} \quad (\text{A.7})$$

$$\hat{r}_S^\diamond = \frac{m+1}{h(m-1)}(\bar{\delta}_S - \hat{\delta}_S), \quad \text{with } \bar{\delta}_S = \bar{\delta}_L \text{ of (2.10);} \quad (\text{A.8})$$

and $\hat{r}_S^+ = \max(0, \hat{r}_S)$. The stacked counterparts of \hat{D}_L^\diamond and its existing counterparts \hat{D}_L and \hat{D}_L^+ (see (2.11)) then are given by

$$\hat{D}_S^\diamond = \hat{D}_S(\hat{r}_S^\diamond), \quad \hat{D}_S = \hat{D}_S(\hat{r}_S), \quad \hat{D}_S^+ = \hat{D}_S(\hat{r}_S^+). \quad (\text{A.9})$$

The approximation $\hat{d}_L \simeq \hat{d}_S$ is still true under the following conditions.

Assumption A.1. (a) Define $R(\psi) = \overline{L}^S(\psi) - \overline{L}(\psi)$, where

$$\overline{L}(\psi) = (mn)^{-1} \sum_{\ell=1}^m \log f(X^{(\ell)} | \psi) \quad \text{and} \quad \overline{L}^S(\psi) = (mn)^{-1} \log f(X^S | \psi).$$

For each m , as $n \rightarrow \infty$,

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

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

$$\sup_{\psi \in \Psi} |\overline{L}(\psi) - \underline{\mathcal{L}}(\psi)| = o_p(1).$$

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

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

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

Theorem A.2. Suppose RC_θ and Assumption A.1. Under both H_0 and H_1 , we have (i) $\hat{d}_S, \hat{r}_S \geq 0$ for all m, n ; (ii) \hat{d}_S, \hat{r}_S are invariant to the parametrization of ψ for all m, n ; and (iii) $\hat{d}_L \simeq \hat{d}_S$ and $\hat{r}_L \simeq \hat{r}_S$ as $n \rightarrow \infty$ 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 A.1 holds.

Example A.1. Consider a stationary autoregressive model of order one. Suppose the complete data $X = (X_1, \dots, X_n)^\top$ 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 $\psi = (\phi, \sigma^2)^\top$, and

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

Then, it is easy to see that condition (a) of Assumption A.1 is satisfied.

A.5 Other existing MI tests

First, we list some existing estimators of \boldsymbol{r}_m . Let $s_{W,a}^2$ be the sample variances of $\{(d_W^{(\ell)})^a\}_{\ell=1}^m$ for $a > 0$. Rubin (2004) and Li *et al.* (1991) proposed

$$\tilde{r}_{W,1} = \frac{(1 + 1/m)s_{W,1}^2}{2\bar{d}_W + \sqrt{\max\{0, 4\bar{d}_W^2 - 2ks_{W,1}^2\}}}, \quad (\text{A.10})$$

$$\tilde{r}_{W,1/2} = (1 + 1/m)s_{W,1/2}^2, \quad (\text{A.11})$$

respectively. When k is large and m is small, using (A.10) or (A.11) may lead to power loss. A trivial modification of \tilde{r}_L of (1.8), i.e., $\tilde{r}_L^+ = \max(0, \tilde{r}_L)$, is a better alternative.

Second, we list some alternative MI combining rules. Having the above estimators of \boldsymbol{r}_m , we can insert them into the following combining rules:

$$\tilde{D}'_W(\boldsymbol{r}_m) = \frac{\tilde{d}'_W}{k(1 + \boldsymbol{r}_m)}, \quad \tilde{D}_L(\boldsymbol{r}_m) = \frac{\tilde{d}_L}{k(1 + \boldsymbol{r}_m)}, \quad \tilde{D}_L^+(\boldsymbol{r}_m) = \left\{ \tilde{D}_L(\boldsymbol{r}_m) \right\}^+. \quad (\text{A.12})$$

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

$$\bar{D}'_W(\boldsymbol{r}_m) = \frac{\bar{d}'_W - \frac{k(m-1)}{m+1}\boldsymbol{r}_m}{k(1 + \boldsymbol{r}_m)}, \quad \bar{D}_L(\boldsymbol{r}_m) = \frac{\bar{d}_L - \frac{k(m-1)}{m+1}\boldsymbol{r}_m}{k(1 + \boldsymbol{r}_m)}; \quad (\text{A.13})$$

see, e.g., Li *et al.* (1991). The combining rule $\overline{D}'_{\text{W}}(\boldsymbol{\nu}_m)$ is useful when computing \overline{d}'_{W} and estimating $\boldsymbol{\nu}_m$ are simple, but the resulting power may deteriorate. If $\tilde{r}_{\text{W},1}$ or $\tilde{r}_{\text{W},1/2}$ is used for estimating $\boldsymbol{\nu}_m$, the null distribution of (A.12) and (A.13) can be approximated by $F_{k, \tilde{\text{df}}'(\boldsymbol{\nu}_m, k)}$, where $\tilde{\text{df}}'(\boldsymbol{\nu}_m, k) = (m-1)(1 + \boldsymbol{\nu}_m^{-1})^2 k^{-3/m}$; see Li *et al.* (1991).

Next, we introduce and recall some notation: (a) standard complete-data moments estimation (\mathcal{M}_{W} , \mathcal{M}_{L}) and testing procedures (\mathcal{D}_{W} , \mathcal{D}_{L}), and (b) non-standard complete-data procedures ($\tilde{\mathcal{D}}_{\text{L}}$, $\overline{\mathcal{D}}_{\text{L}}$, $\mathcal{D}_{\text{L},1}$, $\overline{\mathcal{D}}_{\text{L},1}$), where

$$\begin{aligned} \mathcal{M}_{\text{W}}(X) &= \left\{ \hat{\theta}(X), U(X) \right\}, & \mathcal{M}_{\text{L}}(X) &= \left\{ \hat{\psi}(X), \hat{\psi}_0(X) \right\}, \\ \mathcal{D}_{\text{W}}(X) &= d_{\text{W}}(\hat{\theta}(X), U(X)), & \overline{\mathcal{D}}_{\text{L},1}(\mathbb{X}) &= \frac{2}{m} \sum_{\ell=1}^m \log f(X^{(\ell)} \mid \hat{\psi}^*(\mathbb{X})), \\ \overline{\mathcal{L}}(\psi) &= \frac{1}{m} \log f_{mn}(X^{(1:m)} \mid \psi). \end{aligned}$$

Table A.2 is the full version of Table 1 in the main text. It summarizes the statistical and computational properties of different MI tests; see Section 3.2 for details.

Table A.2: 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 Section 3.2 for heading descriptions. The reference papers/book are abbreviated as follows: Rubin (2004) (R04), Li *et al.* (1991) (LMRR91) and Meng and Rubin (1992) (MR92).

Test	No.	Combining Rule		Estimator of r_m		Approx. null distribution ^a		Reference	Properties						
		Formula	Routine	Formula	Routine	Original	Proposed		Inv	Con	≥ 0	Pow	Def	Sca	EFMI
WT	WT-1	$D_W(T)^b$	\mathcal{M}_W	\tilde{r}'_W	\mathcal{D}_W	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	R04	-	+	+	-	-	-	θ^c
	WT-2	$\tilde{D}'_W(r_m)$	\mathcal{M}_W	\tilde{r}'_W	\mathcal{M}_W	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	R04	-	+ ^e	+	-	-	-	θ
	WT-3	$\tilde{D}'_W(r_m)^f$	\mathcal{M}_W	$\tilde{r}'_{W,1}$	\mathcal{D}_W	$F_{k,\tilde{d}\tilde{f}'(r_m,k)}$	NA	R04	-	-	+	-	-	-	θ
	WT-4	$\tilde{D}'_W(r_m)$	\mathcal{M}_W	$\tilde{r}'_{W,1/2}$	\mathcal{D}_W	$F_{k,\tilde{d}\tilde{f}'(r_m,k)}$	NA	LMRR91	-	-	+	-	-	-	θ
	WT-5	$\tilde{D}'_W(r_m)$	\mathcal{D}_W	$\tilde{r}'_{W,1}$	\mathcal{D}_W	$F_{k,\tilde{d}\tilde{f}'(r_m,k)}$	NA	R04	-	-	-	-	-	+	θ
	WT-6	$\tilde{D}'_W(r_m)$	\mathcal{D}_W	$\tilde{r}'_{W,1/2}$	\mathcal{D}_W	$F_{k,\tilde{d}\tilde{f}'(r_m,k)}$	NA	LMRR91	-	-	-	-	-	+	θ
LRT	LRT-1	$\tilde{D}_L(r_m)$	$\mathcal{M}_L, \tilde{\mathcal{D}}_L$	\tilde{r}_L	$\mathcal{M}_L, \tilde{\mathcal{D}}_L$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	MR92	-	-	-	-	+	- ^g	θ
	LRT-2	$\hat{D}_L(r_m)$	\mathcal{D}_L	\hat{r}_L^+	\mathcal{D}_L	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	Proposal	+	-	+	-	+	+	θ
	LRT-3	$\hat{D}_L(r_m)$	\mathcal{D}_L	\hat{r}_L^\diamond	$\mathcal{D}_{L,1}$	$F_{k,\tilde{d}\tilde{f}(r_m,h)}$	$F_{k,\tilde{d}\tilde{f}(r_m,h)}$	Proposal	+	+	+	+	+	+	ψ
	LRT-4	$\tilde{D}_L^+(r_m)$	$\mathcal{M}_L, \tilde{\mathcal{D}}_L$	\tilde{r}_L^+	$\mathcal{M}_L, \tilde{\mathcal{D}}_L$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	MR92 ^h	-	-	+	-	+	-	θ
	LRT-5	$\hat{D}_L(r_m)$	\mathcal{D}_L	\hat{r}_L	\mathcal{D}_L	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	$F_{k,\tilde{d}\tilde{f}(r_m,k)}$	Proposal	+	-	-	-	+	+	θ

^aIn actual computation, the r_m in the denominator degree of freedom of F is replaced by its corresponding estimator.

^bComputing the test statistic $D_W(T) = d_W(\hat{\theta}, T)/k$ does not require estimating r_m .

^cEFMI is not required for the test statistic $D_W(T)$, but it is required for its approximate null distribution.

^dThe approximate null distribution documented in Rubin (2004) was modified by Li *et al.* (1991). This also applies to WT-2,4,5.

^eThe estimator \tilde{r}'_W does not depend on θ_0 , but its MSE may be inflated under H_1 if a bad parametrization of θ is used.

^fThe originally proposed combining rule is $\tilde{D}'_W(r_m)$; see (A.13). Although $\tilde{D}'_W(r_m)$ is more computationally feasible, the power loss is more significant than $\tilde{D}'_W(r_m)$ after inserting an inefficient estimator $\tilde{r}'_{W,1}$ for r_m . This footnote also applies to WT-3.

^gAveraging and processing vector estimators of ψ , but not their covariance matrixes, is needed. This footnote also applies to LRT-2.

^hIt is a trivial modification of the original proposal in MR92 by replacing \tilde{r}_L with $\tilde{r}_L^+ = \max\{0, \tilde{r}_L\}$.

Table A.3: The values of parameters used in the simulation experiment in Section 4.1.

Experiment		Fixed Parameters			Variable Parameter				
No.	Variable Parameter	ρ	p	ℓ	Case 1	Case 2	Case 3	Case 4	Case 5
I	Correlation ρ	–	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 ℓ	0.4	2	–	0.1	0.3	0.5	0.7	0.9

A.6 Supplement for Section 4.1

Let \bar{X}_{obs} and S_{obs} be the sample mean and sample covariance matrix based on X_{obs} . Then, the ℓ th imputed missing data set can be produced by the following procedure, for $\ell = 1, \dots, m$.

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

Also, denote $X_i^{(\ell)} = X_i$ for $i = 1, \dots, n_{\text{obs}}$. With the ℓ th completed data set, the unconstrained MLEs for μ and Σ are

$$\hat{\mu}^{(\ell)} = \frac{1}{n} \sum_{i=1}^n X_i^{(\ell)}, \quad \hat{\Sigma}^{(\ell)} = \frac{1}{n} \sum_{i=1}^n \left(X_i^{(\ell)} - \hat{\mu}^{(\ell)} \right) \left(X_i^{(\ell)} - \hat{\mu}^{(\ell)} \right)^\top.$$

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

$$\hat{\mu}_0^{(\ell)} = \left\{ \frac{\mathbf{1}_p^\top (\hat{\Sigma}^{(\ell)})^{-1} \hat{\mu}^{(\ell)}}{\mathbf{1}_p^\top (\hat{\Sigma}^{(\ell)})^{-1} \mathbf{1}_p} \right\} \mathbf{1}_p, \quad \hat{\Sigma}_0^{(\ell)} = \hat{\Sigma}^{(\ell)} + \left(\hat{\mu}^{(\ell)} - \hat{\mu}_0^{(\ell)} \right) \left(\hat{\mu}^{(\ell)} - \hat{\mu}_0^{(\ell)} \right)^\top.$$

We first study the distribution of p -values of each test under H_0 . We use $n = 100$, $m = 3$, $\sigma^2 = 5$ and $\mu = \mathbf{1}_p$, with various values of ρ , p and ℓ specified in Table A.3. The results under parametrizations (i), (ii) and (iii) are shown

in Figures A.4, A.5 and A.6, respectively. Note that, for Wald tests under parametrization (ii), the matrix $U^{(\ell)}$ is singular in 0.25% of the replications, and those cases are removed from the analysis (which should favor the Wald tests).

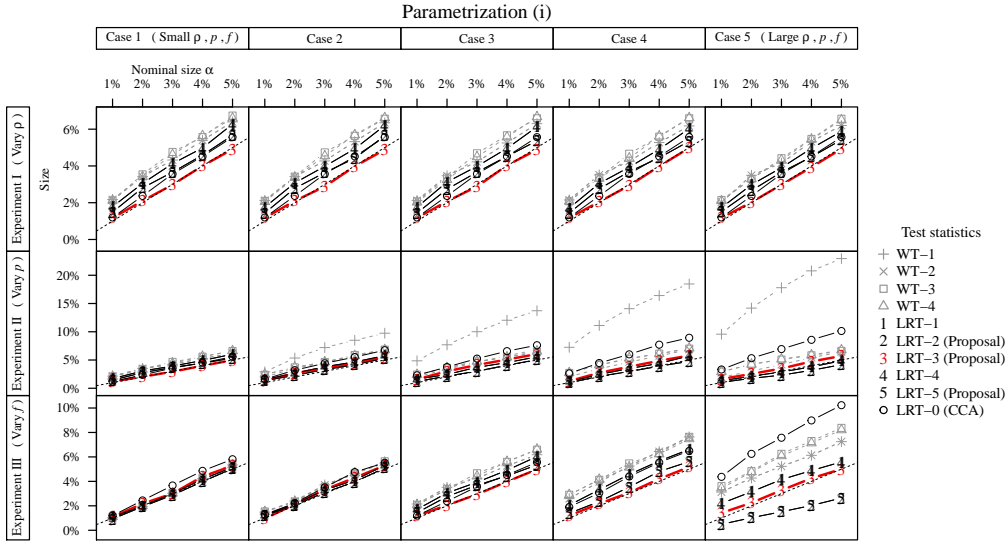


Figure A.4: The comparison between empirical size and nominal size α under parametrization (ii) for $\alpha \in (0, 5\%]$. Our most recommended proposal is LRT-3, which is highlighted red.

The empirical sizes (i.e., type-I errors) of the MI Wald tests generally deviate from the nominal size α under parametrization (ii). In contrast, the sizes of all LRTs are closer to α . However, the original L-1 and its trivial modification L-2 do not have accurate sizes when $|\rho|$ or ℓ 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 ℓ are large. 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

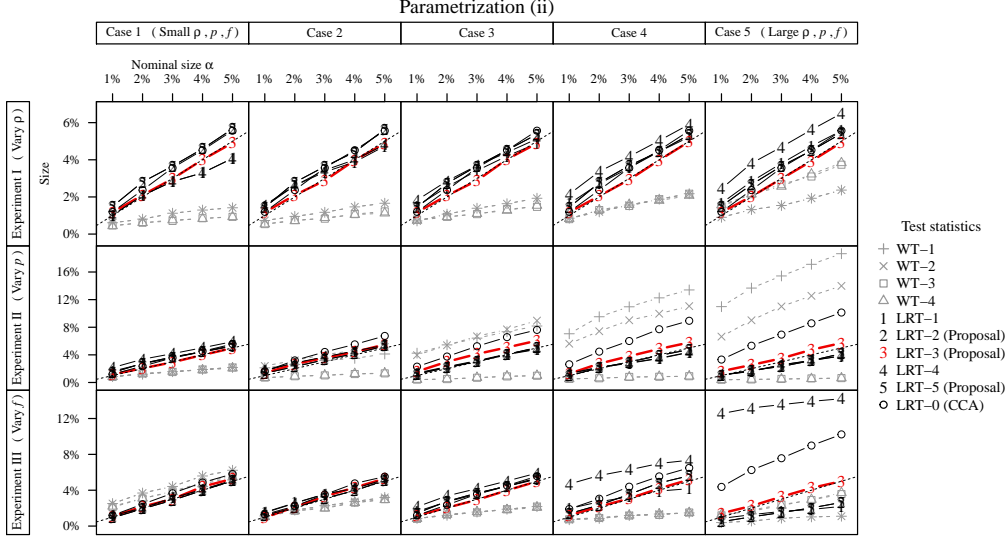


Figure A.5: The comparison between empirical size and nominal size α under parametrization (i) for $\alpha \in (0, 5\%]$. The legend in Figure A.4 also applies here.

$m = 3$.

Interestingly, as seen clearly in Figure A.5, the benchmark L-0 performs very badly for large p and ℓ . This is because the sample size per parameter, n/h , is small; for $p \geq 4$, $n/h \leq 100/14 < 8$. The asymptotic null distribution χ_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 (see Barnard and Rubin, 1999). But this is exactly what is being used for MI tests, albeit with different choices of the denominator degrees of freedom. Note also that, in some cases, nearly half of the simulated values of \tilde{r}_L and \tilde{D}_L are negative; see Table A.4. In contrast, \hat{r}_S is always non-negative in our simulation, despite the fact that it can be negative in theory.

The power curves under nominal size 0.5% and 5% are shown in Figure 2 of the main text and Figure A.7, respectively. Note that the trivial modifications LRT-2 of LRT-1 cannot retrieve all the power it should have. Tables A.5 and A.6 show the minimum and maximum of the empirical sizes over the three parametrizations considered in each test — and only one value is needed for

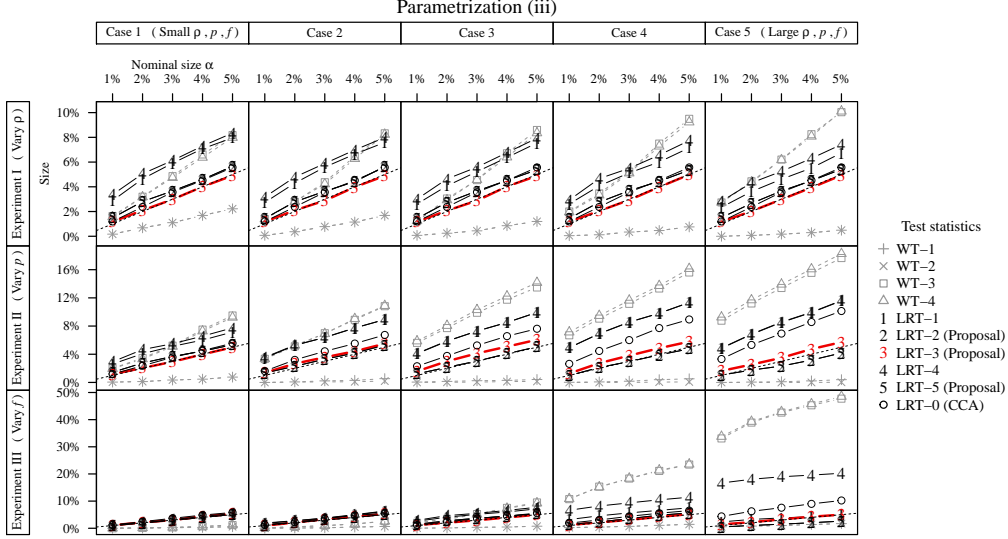


Figure A.6: The comparison between empirical size and nominal size α under parametrization (iii) for $\alpha \in (0, 5\%]$. The legend in Figure A.5 also applies here.

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

Experiment	Parametrization	Case									
		1 2 3 4 5					1 2 3 4 5				
		% of $\tilde{r}_L < 0$					% of $\tilde{D}_L < 0$				
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

those tests that are invariant to parametrization — when the nominal size is 0.5% and 5%, respectively. We see the deviations from the nominal α 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.*,

Table A.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 parametrization-invariant tests. The nominal size is $\alpha = 0.5\%$. The results under nominal size $\alpha = 5\%$ are shown in Figure A.6.

(n, m)	Range of empirical size: $[\min \hat{\alpha}, \max \hat{\alpha}]/\%$				
	(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

2016). Figures A.8 and A.9 plot the result for nominal size 0.5% and 5%, respectively. 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 $\delta \rightarrow \infty$. Nevertheless, the performances of all size 0.5% tests are less satisfactory than those for size 5% tests because larger sample sizes n are required to approximate the tail behavior well.

We also compare the performance of estimators of ν_m for different δ and parametrizations. In our experiment, we have $\nu_m = 1 + 1/m$ because we have set $\nu = 1$. The MSEs of estimators $\hat{f} = \hat{r}/(1 + \hat{r})$ of $f_m = \nu_m/(1 + \nu_m)$ are shown in Figure A.10, in log scale. Clearly, the only estimator that is consistent, invariant to parametrization and robust against δ is our proposal $\hat{f}_L^\diamond = \hat{r}_L^\diamond/(1 + \hat{r}_L^\diamond)$. It concentrates at the true value f_m quite closely even for small m and n . It verifies why L-5 has the greatest power. On the other hand, the estimator $\tilde{f}_L = \tilde{r}_L/(1 + \tilde{r}_L)$ has a large MSE when $\delta \neq 0$. It explains why L-1 is not powerful.

Table A.6: 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 parametrization-invariant tests. The nominal size is $\alpha = 5\%$.

(n, m)	Range of empirical size: $[\min \hat{\alpha}, \max \hat{\alpha}]/\%$				
	(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

A.7 Supplements for Section 4.2

Let $n_j = \sum_{i=1}^n R_{ij}$ be the number of observed j th component. Without loss of generality, assume X_{obs} is arranged in such a way that $R_{ij} \geq R_{i'j}$ for all $i < i'$ and j . To impute the missing data, it is useful to represent X_i by

$$[X_{i1} \mid \beta_1, \tau_1^2] \sim \mathcal{N}(\beta_1, \tau_1^2) \quad \text{and} \quad [X_{ij} \mid X_{i,1:(j-1)}, \beta_j, \tau_j^2] \sim \mathcal{N}(\beta_j^\top Z_{ij}, \tau_j^2),$$

for $j = 2, \dots, p$, where $\tau_1^2, \dots, \tau_p^2 \in \mathbb{R}^+$, $\beta_j \in \mathbb{R}^j$, $X_{i,1:(j-1)} = (X_{i1}, \dots, X_{i,j-1})^\top$ and $Z_{ij} = (1, X_{i,1:(j-1)}^\top)^\top$ for $j \geq 2$. Denote the (complete-case) least squares estimators of β_j and τ_j^2 respectively by

$$\hat{\beta}_j = (Z_j^\top Z_j)^{-1} Z_j^\top W_j \quad \text{and} \quad \hat{\tau}_j^2 = \frac{(W_j - Z_j \hat{\beta}_j)^\top (W_j - Z_j \hat{\beta}_j)}{n_j - j},$$

where $Z_j = (Z_{1j}, \dots, Z_{n_j j})^\top$ and $W_j = (X_{1j}, \dots, X_{n_j j})^\top$.

We assume a Bayesian imputation model with the non-informative prior $f(\beta_1, \dots, \beta_p, \tau_1^2, \dots, \tau_p^2) \propto 1/(\tau_1^2 \cdots \tau_p^2)$. For $\ell = 1, \dots, m$, denote the ℓ th imputed data set by $X^{(\ell)}$, whose (i, j) th element is $X_{ij}^{(\ell)}$. If $1 \leq j \leq p$ and $i \leq n_j$, then $X_{ij}^{(\ell)} = X_{ij}$, otherwise $X_{ij}^{(\ell)}$ is filled in by recursing the following steps for $j = 2, \dots, p$.

1. Draw a sample $(\tau_j^{(\ell)})^2$ from $\hat{\tau}_j^2(n_j - j)/\chi_{n_j - j}^2$.

2. Draw a sample $\beta_j^{(\ell)}$ from $\mathcal{N}_j(\hat{\beta}_j, (\tau_j^{(\ell)})^2(Z_j^\top Z_j)^{-1})$.
3. Draw a sample $X_{ij}^{(\ell)}$ from $\mathcal{N}((\beta_j^{(\ell)})^\top Z_{ij}^{(\ell)}, (\tau_j^{(\ell)})^2)$ for $i = n_j + 1, \dots, n$, where $Z_{ij}^{(\ell)} = (1, (X_{i,1:(j-1)}^{(\ell)})^\top)^\top$.

With the ℓ th imputed data set, the H_0 -constrained MLEs of μ and Σ are $\hat{\mu}_0^{(\ell)} = \mathbf{0}_p$ and $\hat{\Sigma}_0^{(\ell)} = (X^{(\ell)})^\top(X^{(\ell)})/n$; whereas the unconstrained counterparts are $\hat{\mu}^{(\ell)} = \mathbf{1}_n^\top X^{(\ell)}/n$ and $\hat{\Sigma}^{(\ell)} = (X^{(\ell)} - \hat{\mu}^{(\ell)})^\top(X^{(\ell)} - \hat{\mu}^{(\ell)})/n$.

The partial result is shown in Figure 3 of the main text, whereas the full version is shown in Figure A.11.

A.8 Applications to a Care-Survival Data

Meng and Rubin (1992) considered the data given in Table A.7, where i , j and k index, respectively, amount of parental care (less or more, corresponding to $i = 1, 2$), and survival status (died or survived, corresponding to $j = 1, 2$), and clinic (A or B, corresponding to $k = 1, 2$). The label k is missing for some observations. The missing mechanism was assumed to be ignorable. We consider two null hypotheses: (H_0) the clinic and parental care are conditionally independent given the survival status, and (H'_0) all three variables are independent. It is remarked that testing the conditional independence model (i.e., H_0) is useful from a modeling perspective. If H_0 cannot be rejected, then one may be tempted to adopt the more parsimonious null model (for the cell probabilities). The same model is also suggested in Little and Rubin (2002) and Meng and Rubin (1992).

Our aim is to investigate the impact on $\{\tilde{D}_S, \hat{D}_S^+, \hat{D}_S^\diamond\}$ by the parametrization of the cell probabilities

$$\pi_{ijk} = \text{P}(\text{parental care} = i, \text{survival status} = j, \text{clinic label} = k)$$

for $i, j, k \in \{1, 2\}$; and the impact on $\{\tilde{r}_L, \hat{r}_S^+, \hat{r}_S^\diamond\}$ under different null hypotheses. Here the full model parameter vector can be expressed as $\psi = (\pi_{111}, \pi_{112}, \pi_{121}, \pi_{122}, \pi_{211}, \pi_{212}, \pi_{221})^\top$. Since the restrictions imposed by H_0 are $\pi_{ijk} = (\pi_{1jk} + \pi_{2jk})(\pi_{ij1} + \pi_{ij2})$ for $j = 1, 2$, one may express the parameter of interest as $\theta = (\theta_1, \theta_2)^\top$, where $\theta_j = \pi_{ijk} - (\pi_{1jk} + \pi_{2jk})(\pi_{ij1} + \pi_{ij2})$ for

$j = 1, 2$. Then H_0 can be equivalently stated as $\theta = \theta_0$, where $\theta_0 = (0, 0)^\top$. Similarly, the parameter of interest under H'_0 can be defined.

Table A.7: Data from Meng and Rubin (1992). The notation “?” indicates missing label.

Parental care (i)		Less		More	
		Died	Survived	Died	Survived
Survival Status (j)					
Clinic Label (k)	A	3	176	4	293
	B	17	197	2	23
	?	10	150	5	90

The computation of the stacked MI estimators of $\{\pi_{ijk}\}$ is presented in A.8 of the Appendix. We consider three parametrizations: (i) $\psi_{ijk} = \pi_{ijk}$; (ii) $\psi_{ijk} = \log\{\pi_{ijk}/(1 - \pi_{ijk})\}$; and (iii) $\psi_{ij1} = \pi_{ij1}$ and $\psi_{ij2} = \pi_{ij2}/\pi_{ij1}$. Denote the p -values of tests $\{\tilde{D}_L, \hat{D}_S^+, \hat{D}_S^\diamond\}$ by $\{\tilde{p}_L, \hat{p}_S^+, \hat{p}_S^\diamond\}$, respectively. The results are summarized in Table A.8. Clearly, only $\hat{r}_S, \hat{r}_\diamond, \hat{D}_S^+, \hat{D}_S^\diamond$ are always non-negative and parametrization-invariant. Some of the values of \tilde{r}_L and \tilde{D}_L are negative, leading to the meaningless $\tilde{p}_L = 1$. For testing H_0 , we have $\hat{D}_S^+ \approx \hat{D}_S^\diamond$. For testing H'_0 , \hat{D}_S^+ and \hat{D}_S^\diamond are not very close to each other, but they both lead to essentially zero p -value. These results reconfirm the conclusions in Meng and Rubin (1992). Moreover, only \hat{r}_S^\diamond does not change under different null hypotheses.

The MI data sets are generated from a Bayesian model in Section 4.2 of Meng and Rubin (1992). The ℓ th imputed log-likelihood function is $\log f(X^{(\ell)} | \pi) = \sum_c n_c^{(\ell)} \log \pi_c$, where $X^{(\ell)}$ are the cell counts $n_c^{(\ell)}$ in the ℓ th imputed data set. Hence the unconstrained MLE of π_c is $\hat{\pi}_c^{(\ell)} = n_c^{(\ell)}/n_+^{(\ell)}$, where $n_+^{(\ell)} = \sum_c n_c^{(\ell)}$. Let $n_c^+ = \sum_{\ell=1}^m n_c^{(\ell)}$. Consequently, the joint log-likelihood based on the stacked data is

$$\log f(X^S | \pi) = \sum_{\ell=1}^m \sum_c n_c^{(\ell)} \log \pi_c = \sum_c n_c^+ \log \pi_c, \quad (\text{A.14})$$

Thus the unconstrained MLE with respect to (A.14) is $\hat{\pi}_c^S = n_c^+/n_+^+$, where $n_+^+ = \sum_c n_c^+$. Similarly, we can find the constrained MLEs under a given null.

Table A.8: The LRTs using \tilde{D}_L , \hat{D}_S^+ and \hat{D}_S^\diamond under different parametrizations in Section A.8.

Parametrization (i): identity map						
H_0 : Conditional independence			H_0 : Full independence			
m	$\tilde{r}_L, \hat{r}_S^+, \hat{r}_S^\diamond$	$\tilde{D}_L, \hat{D}_S^+, \hat{D}_S^\diamond$	$\tilde{p}_L, \hat{p}_S^+, \hat{p}_S^\diamond$	$\tilde{r}_L, \hat{r}_S^+, \hat{r}_S^\diamond$	$\tilde{D}_L, \hat{D}_S^+, \hat{D}_S^\diamond$	$\tilde{p}_L, \hat{p}_S^+, \hat{p}_S^\diamond$
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): logit transformation						
H_0 : Conditional independence			H_0 : Full independence			
m	$\tilde{r}_L, \hat{r}_S^+, \hat{r}_S^\diamond$	$\tilde{D}_L, \hat{D}_S^+, \hat{D}_S^\diamond$	$\tilde{p}_L, \hat{p}_S^+, \hat{p}_S^\diamond$	$\tilde{r}_L, \hat{r}_S^+, \hat{r}_S^\diamond$	$\tilde{D}_L, \hat{D}_S^+, \hat{D}_S^\diamond$	$\tilde{p}_L, \hat{p}_S^+, \hat{p}_S^\diamond$
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
Parametrization (iii): ratios of probabilities						
H_0 : Conditional independence			H_0 : Full independence			
m	$\tilde{r}_L, \hat{r}_S^+, \hat{r}_S^\diamond$	$\tilde{D}_L, \hat{D}_S^+, \hat{D}_S^\diamond$	$\tilde{p}_L, \hat{p}_S^+, \hat{p}_S^\diamond$	$\tilde{r}_L, \hat{r}_S^+, \hat{r}_S^\diamond$	$\tilde{D}_L, \hat{D}_S^+, \hat{D}_S^\diamond$	$\tilde{p}_L, \hat{p}_S^+, \hat{p}_S^\diamond$
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

B Proofs

Proof of Theorem 1. (i, ii) From (2.3), we know $\hat{d}_L \geq 0$ is invariant to parametrization ψ . (iii) Since \hat{d}_L is invariant to transformation of ψ , we assume, without loss of generality, that ψ admits a parameterization such that $\text{Cov}(\hat{\theta}^{(\ell)}, \hat{\eta}^{(\ell)}) \simeq \mathbf{0}$ by taking suitable linear transformation of ψ . Also write $U_\eta^{(\ell)}$ as an efficient

estimator of $\text{Var}(\hat{\eta})$ based on $X^{(\ell)}$; and recall that $U_\theta^{(\ell)} = U^{(\ell)}$ is an efficient estimator of $\text{Var}(\hat{\theta})$ based on $X^{(\ell)}$.

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

$$\bar{L}(\psi) \simeq \bar{L}(\hat{\psi}^*) - \frac{1}{2} (\psi - \hat{\psi}^*)^\top \bar{I}(\hat{\psi}^*) (\psi - \hat{\psi}^*), \quad (\text{B.1})$$

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

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

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

$$\begin{aligned} \hat{d}_L &\simeq (\hat{\psi}_0^* - \hat{\psi}^*)^\top \bar{I}(\hat{\psi}^*) (\hat{\psi}_0^* - \hat{\psi}^*), \\ &\simeq \begin{pmatrix} \theta_0 - \hat{\theta}^* \\ \hat{\eta}(\theta_0) - \hat{\eta}(\hat{\theta}^*) \end{pmatrix}^\top \begin{pmatrix} \bar{U}_\theta^{-1} & \mathbf{0} \\ \mathbf{0} & \bar{U}_\eta^{-1} \end{pmatrix} \begin{pmatrix} \theta_0 - \hat{\theta}^* \\ \hat{\eta}(\theta_0) - \hat{\eta}(\hat{\theta}^*) \end{pmatrix} \\ &\simeq (\bar{\theta}^\top - \theta_0) \bar{U}_\theta^{-1} (\bar{\theta}^\top - \theta_0) = \tilde{d}_W, \end{aligned} \quad (\text{B.3})$$

where we have used (a) $\hat{\theta}^* \simeq \bar{\theta}$; 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 $\tilde{d}_W \simeq \hat{d}_L$ (Meng and Rubin, 1992), we have $\hat{d}_L \simeq \tilde{d}_L$. \square

Proof of Proposition 1. The given condition implies that

$$\begin{aligned} \hat{\psi}^{(\ell)} &= ((\hat{\theta}^{(\ell)})^\top, (\hat{\eta}^{(\ell)})^\top)^\top, & \hat{\psi}_0^{(\ell)} &= (\theta_0^\top, (\hat{\eta}^{(\ell)})^\top)^\top, \\ \hat{\psi}^* &= ((\hat{\theta}^*)^\top, (\hat{\eta}^*)^\top)^\top, & \hat{\psi}_0^* &= (\theta_0^\top, (\hat{\eta}^*)^\top)^\top. \end{aligned}$$

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

$$\begin{aligned} \bar{d}_L - \hat{d}_L &= \frac{2}{m} \sum_{\ell=1}^m \left\{ L^{(\ell)}(\hat{\psi}^{(\ell)}) - L^{(\ell)}(\hat{\psi}_0^{(\ell)}) - L^{(\ell)}(\hat{\psi}^*) + L^{(\ell)}(\hat{\psi}_0^*) \right\} \\ &= \frac{2}{m} \sum_{\ell=1}^m \left\{ L_\dagger^{(\ell)}(\hat{\theta}^{(\ell)}) - L_\dagger^{(\ell)}(\hat{\theta}^*) \right\} \geq 0 \end{aligned}$$

since $L_\dagger^{(\ell)}(\hat{\theta}^{(\ell)}) \geq L_\dagger^{(\ell)}(\hat{\theta}^*)$ for all ℓ . \square

Proof of Corollary 1. 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 $\hat{\psi}_0^{(\ell)}$ such that

$$L^{(\ell)}(\hat{\psi}_0^{(\ell)}) = L^{(\ell)}(\hat{\psi}^{(\ell)}) - \frac{1}{2} \left(\hat{\psi}_0^{(\ell)} - \hat{\psi}^{(\ell)} \right)^\top I^{(\ell)}(\check{\psi}^{(\ell)}) \left(\hat{\psi}_0^{(\ell)} - \hat{\psi}^{(\ell)} \right),$$

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

$$\begin{aligned} L^{(\ell)}(\hat{\psi}^{(\ell)}) - L^{(\ell)}(\hat{\psi}_0^{(\ell)}) &\simeq \frac{1}{2} \left(\hat{\psi}_0^{(\ell)} - \hat{\psi}^{(\ell)} \right)^\top I^{(\ell)}(\check{\psi}^*) \left(\hat{\psi}_0^{(\ell)} - \hat{\psi}^{(\ell)} \right) \\ &\simeq \frac{1}{2} \left(\theta_0 - \hat{\theta}^{(\ell)} \right)^\top \check{U}^{-1} \left(\theta_0 - \hat{\theta}^{(\ell)} \right) \end{aligned} \quad (\text{B.4})$$

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

$$L^{(\ell)}(\hat{\psi}^*) - L^{(\ell)}(\hat{\psi}_0^*) \simeq \frac{1}{2} \left(\theta_0 - \hat{\theta}^* \right)^\top \check{U}^{-1} \left(\theta_0 - \hat{\theta}^* \right). \quad (\text{B.5})$$

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

$$\begin{aligned} \bar{d}_L - \hat{d}_L &\simeq \frac{1}{m} \sum_{\ell=1}^m \left\{ \left(\theta_0 - \hat{\theta}^{(\ell)} \right)^\top \check{U}^{-1} \left(\theta_0 - \hat{\theta}^{(\ell)} \right) - \left(\theta_0 - \hat{\theta}^* \right)^\top \check{U}^{-1} \left(\theta_0 - \hat{\theta}^* \right) \right\} \\ &= \text{tr} \left[\check{U}^{-1} \left\{ \frac{1}{m} \sum_{\ell=1}^m \left(\theta_0 - \hat{\theta}^{(\ell)} \right)^{\otimes 2} - \left(\theta_0 - \hat{\theta}^* \right)^{\otimes 2} \right\} \right] \\ &\simeq \text{tr} \left[\check{U}^{-1} \frac{1}{m} \sum_{\ell=1}^m \left\{ \left(\hat{\theta}^{(\ell)} \right)^{\otimes 2} - \bar{\theta}^{\otimes 2} \right\} \right] \simeq \text{tr} \left(\check{U}^{-1} B \right) \simeq \text{tr} \left(\mathcal{U}_{\theta,0}^{-1} \mathcal{B}_\theta \right) \end{aligned}$$

as $m, n \rightarrow \infty$, where $\mathcal{U}_{\theta,0}$ is a deterministic matrix that depends on both θ_0 and the true value of θ , and satisfies $n(\check{U} - \mathcal{U}_{\theta,0}) \xrightarrow{\text{pr}} 0$. Note that $\text{tr}(\mathcal{U}_{\theta,0}^{-1} \mathcal{B}_\theta) = k\boldsymbol{\nu}_0$, for some finite $\boldsymbol{\nu}_0$ by Assumption 2. Then $\hat{r}_L \xrightarrow{\text{pr}} \boldsymbol{\nu}_0 = \text{tr}(\mathcal{U}_{\theta,0}^{-1} \mathcal{B}_\theta) / k$, proving (ii). (But $\mathcal{U}_{\theta,0}$ may not equal to \mathcal{U}_θ , and hence \hat{r}_L may not be consistent for $\boldsymbol{\nu}_m$.)

If H_0 is true, then $\bar{\theta} \xrightarrow{\text{pr}} \theta_0$ and $\check{U} \simeq \bar{U} \simeq \mathcal{U}_\theta = \mathcal{U}_{\theta,0}$. Then, $\hat{r}_L \xrightarrow{\text{pr}} \boldsymbol{\nu}$ as $m, n \rightarrow \infty$. So, (i) follows. \square

Proof of Theorem 2. (i, ii) It is trivial by the definition of \hat{r}_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 $\hat{\psi}^{(\ell)}$ and $\hat{\psi}^*$ such that

$$L^{(\ell)}(\hat{\psi}^*) = L^{(\ell)}(\hat{\psi}^{(\ell)}) - \frac{1}{2} \left(\hat{\psi}^* - \hat{\psi}^{(\ell)} \right)^\top I^{(\ell)}(\check{\psi}^{(\ell)}) \left(\hat{\psi}^* - \hat{\psi}^{(\ell)} \right). \quad (\text{B.6})$$

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

$$\begin{aligned} \bar{\delta}_L - \hat{\delta}_L &\simeq \frac{1}{m} \sum_{\ell=1}^m \left(\hat{\psi}^* - \hat{\psi}^{(\ell)} \right)^\top \bar{I}(\hat{\psi}^*) \left(\hat{\psi}^* - \hat{\psi}^{(\ell)} \right) \\ &= \text{tr} \left\{ \bar{I}(\hat{\psi}^*) \frac{1}{m} \sum_{\ell=1}^m \left(\hat{\psi}^* - \hat{\psi}^{(\ell)} \right)^{\otimes 2} \right\} \\ &\simeq \text{tr} \left\{ \bar{I}(\hat{\psi}^*) \frac{1}{m} \sum_{\ell=1}^m \left(\hat{\psi}^{(\ell)} - \bar{\psi} \right)^{\otimes 2} \right\} \simeq \text{tr} \left(\mathcal{U}_\psi^{-1} \mathcal{B}_\psi \right) \end{aligned} \quad (\text{B.7})$$

as $m, n \rightarrow \infty$. By the assumption of EFMI of ψ , we have $\hat{r}_L^\diamond \xrightarrow{\text{PF}} \nu$. \square

Proof of Lemma 1. First, recall that, as $n \rightarrow \infty$, the observed data MLE $\hat{\theta}_{\text{obs}}$ of θ satisfies (2.4), which can be written as $[\hat{\theta}_{\text{obs}} \mid \theta] \stackrel{\mathcal{D}}{\approx} \mathcal{N}_k(\theta, \mathcal{T}_\theta)$, where $A_{1,n} \stackrel{\mathcal{D}}{\approx} A_{2,n}$ means that $A_{1,n}$ and $A_{2,n}$ have the same asymptotic distribution, i.e., there exist deterministic sequences μ_n and Σ_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 $n \rightarrow \infty$,

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

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

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

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

where $Z_1, \dots, Z_m, W \stackrel{\text{iid}}{\sim} \mathcal{N}_k(0, I_k)$. Also write $Z_\ell = (Z_{1\ell}, \dots, Z_{k\ell})^\top$, for $\ell = 1, 2, \dots, m$, and $W = (W_1, \dots, W_k)^\top$. Averaging (B.10) over ℓ , we have $\bar{\theta} \stackrel{\mathcal{D}}{\approx}$

$\hat{\theta}_{\text{obs}} + \mathcal{B}_\theta^{1/2} \bar{Z}_\bullet$, where $\bar{Z}_\bullet = m^{-1} \sum_{\ell=1}^m Z_\ell$. Since $\mathcal{B}_\theta = \nu \mathcal{U}_\theta$, we have

$$\begin{aligned} \mathcal{U}_\theta^{-1/2}(\hat{\theta}^{(\ell)} - \theta) &\stackrel{\mathcal{D}}{\approx} (1 + \nu)^{1/2} W + \nu^{1/2} Z_\ell, \\ \mathcal{U}_\theta^{-1/2}(\bar{\theta} - \theta) &\stackrel{\mathcal{D}}{\approx} (1 + \nu)^{1/2} W + \nu^{1/2} \bar{Z}_\bullet. \end{aligned}$$

Note that (2.6) implies $\mathcal{U}_\theta \simeq \bar{U}$. Under H_0 , we have $\theta = \theta_0$ and

$$\begin{aligned} \bar{d}_L &\simeq \bar{d}_W \stackrel{\mathcal{D}}{\approx} \sum_{i=1}^k \{(1 + \nu)^{1/2} W_i + \nu^{1/2} Z_{i\ell}\}^2, \\ \hat{d}_L &\simeq \tilde{d}_L \simeq \tilde{d}_W \stackrel{\mathcal{D}}{\approx} \sum_{i=1}^k \{(1 + \nu)^{1/2} W_i + \nu^{1/2} \bar{Z}_i\}^2. \end{aligned}$$

After some simple algebra, we obtain

$$\hat{r}_L^+ \stackrel{\mathcal{D}}{\approx} \frac{(m+1)\nu}{mk} \sum_{i=1}^k s_{Z_i}^2 \quad \text{and} \quad \hat{D}_L^+ \stackrel{\mathcal{D}}{\approx} \frac{m \sum_{i=1}^k \{(1 + \nu)^{1/2} W_i + \nu^{1/2} \bar{Z}_i\}^2}{mk + (m+1)\nu \sum_{i=1}^k s_{Z_i}^2},$$

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

$$\hat{r}_L^+ \stackrel{\mathcal{D}}{\approx} \frac{(m+1)\nu}{m(m-1)k} \sum_{i=1}^k H_i^2 \quad \text{and} \quad \hat{D}_L^+ \stackrel{\mathcal{D}}{\approx} \frac{(m-1)\{m + (m+1)\nu\} \sum_{i=1}^k G_i^2}{m(m-1)k + (m+1)\nu \sum_{i=1}^k H_i^2},$$

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 (2.12). \square

Proof of Theorem 3. Similar to (B.9) and (B.10), we can have a more general representation:

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

where $Z_1, \dots, Z_h, W \stackrel{\text{iid}}{\sim} \mathcal{N}_h(0, I_h)$. Also write $Z_\ell = (Z_{1\ell}, \dots, Z_{h\ell})^\top$, for $\ell =$

1, 2, \dots, m, and $W = (W_1, \dots, W_h)^\top$. Using (B.7), we have

$$\begin{aligned} \bar{\delta}_L - \hat{\delta}_L &\simeq \text{tr} \left\{ \bar{I}(\hat{\psi}^*) \frac{1}{m} \sum_{\ell=1}^m (\hat{\psi}^{(\ell)} - \bar{\psi}) (\hat{\psi}^{(\ell)} - \bar{\psi})^\top \right\} \\ &\stackrel{\mathcal{D}}{\approx} \text{tr} \left\{ \mathcal{U}_\psi^{-1} \frac{1}{m} \sum_{\ell=1}^m [(\mathcal{T}_\psi - \mathcal{U}_\psi)^{1/2} (Z_\ell - \bar{Z}_\bullet)]^{\otimes 2} \right\} \\ &= \frac{1}{m} \sum_{\ell=1}^m \text{tr} \left\{ \nu I_h (Z_\ell - \bar{Z}_\bullet)^{\otimes 2} \right\} = \frac{\nu}{m} \sum_{\ell=1}^m \sum_{i=1}^h (Z_{i\ell} - \bar{Z}_{i\bullet})^2. \end{aligned}$$

Equivalently, we can say $\bar{\delta}_L - \hat{\delta}_L \Rightarrow \nu \chi_{h(m-1)}^2 / m$ as $n \rightarrow \infty$. Hence

$$\hat{r}_L^\diamond \Rightarrow \nu \cdot \frac{m+1}{hm(m-1)} \cdot \chi_{h(m-1)}^2,$$

which is equivalent to (2.13). Note that it is true under both H_0 and H_1 . \square

Proof of Theorem 4. From the representations of \hat{d}_L^\diamond and \hat{r}_L^\diamond in Lemma 1 and Theorem 3, we know that they are asymptotically ($n \rightarrow \infty$) independent. The proof then follows the derivation for Lemma 1. \square

Proof of Theorem A.1. (i) Using the representation (A.3), we can easily see that $\hat{r}_L^\triangle \geq 0$. (ii) It suffices to show

$$m^{-1} \sum_{\ell=1}^m d_L(\hat{\psi}_0^{(\ell)} + \Delta_m, \hat{\psi}^{(\ell)} | X^{(\ell)}) \simeq \bar{d}_L - \tilde{d}_L,$$

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

$$L^{(\ell)}(\psi) \simeq L^{(\ell)}(\hat{\psi}^{(\ell)}) - \frac{1}{2} (\psi - \hat{\psi}^{(\ell)})^\top I^{(\ell)}(\hat{\psi}^{(\ell)}) (\psi - \hat{\psi}^{(\ell)}).$$

Under the parametrization of ψ in the proof of Theorem 1, we know that the upper $k \times k$ sub-matrix of $I^{(\ell)}(\hat{\psi}^{(\ell)})$ is $(U^{(\ell)})^{-1}$. Using the lower order

variability of $U^{(\ell)}$, we have $(U^{(\ell)})^{-1} \simeq \bar{U}^{-1}$ and

$$\begin{aligned} \frac{1}{m} \sum_{\ell=1}^m d_L(\hat{\psi}_0^{(\ell)} + \Delta_m, \hat{\psi}^{(\ell)} | X^{(\ell)}) &\simeq \frac{1}{m} \sum_{\ell=1}^m \left(\hat{\psi}_0^{(\ell)} + \Delta_m - \hat{\psi}^{(\ell)} \right)^\top I^{(\ell)}(\hat{\psi}^{(\ell)}) \left(\hat{\psi}_0^{(\ell)} + \Delta_m - \hat{\psi}^{(\ell)} \right) \\ &\simeq \frac{1}{m} \sum_{\ell=1}^m (\hat{\theta}^{(\ell)} - \bar{\theta})^\top \bar{U}^{-1} (\hat{\theta}^{(\ell)} - \bar{\theta}) = \bar{d}'_W - \tilde{d}'_W \simeq \bar{d}_L - \hat{d}_L. \end{aligned}$$

Therefore, the desired result follows. \square

Proof of Theorem A.2. Throughout this proof, conditions (a), (b) and (c) refer to the list given in Assumption A.1. (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

$$\begin{aligned} \bar{L}(\hat{\psi}^*) - \bar{L}(\hat{\psi}^S) &= \bar{L}(\hat{\psi}^*) - \bar{L}^S(\hat{\psi}^S) + \bar{L}^S(\hat{\psi}^S) - \bar{L}(\hat{\psi}^S) \\ &\leq \bar{L}(\hat{\psi}^*) - \bar{L}^S(\hat{\psi}^*) + \bar{L}^S(\hat{\psi}^S) - \bar{L}(\hat{\psi}^S) \\ &\leq 2 \sup_{\psi \in \Psi} \left| \bar{L}(\psi) - \bar{L}^S(\psi) \right| = O_p(1/n), \end{aligned}$$

which, together with condition (b), imply that

$$\begin{aligned} \bar{\mathcal{L}}(\psi^*) - \bar{\mathcal{L}}(\hat{\psi}^S) &= \left\{ \bar{\mathcal{L}}(\psi^*) - \bar{L}(\psi^*) \right\} + \left\{ \bar{L}(\psi^*) - \bar{L}(\hat{\psi}^S) \right\} + \left\{ \bar{L}(\hat{\psi}^S) - \bar{\mathcal{L}}(\hat{\psi}^S) \right\} \\ &\leq 2 \sup_{\psi \in \Psi} \left| \bar{L}(\psi) - \bar{\mathcal{L}}(\psi) \right| + \left\{ \bar{L}(\hat{\psi}^*) - \bar{L}(\hat{\psi}^S) \right\} = o_p(1). \quad (\text{B.11}) \end{aligned}$$

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

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

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

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

$$\bar{L}(\hat{\psi}^*) - \bar{L}(\hat{\psi}^S) = \left\{ \nabla \bar{L}(\check{\psi}) \right\}^\top \left(\hat{\psi}^* - \hat{\psi}^S \right) = o_p(1/n), \quad (\text{B.13})$$

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

$$\bar{L}(\hat{\psi}^*) - \bar{L}^S(\hat{\psi}^S) = R(\hat{\psi}^S) + o_p(1/n). \quad (\text{B.14})$$

Similar to (B.14), we have

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

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

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

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 $\left| \hat{d}_L - \hat{d}_S \right| = o_p(\hat{d}_L)$.
- (ii) Under H_1 , we have $\hat{d}_L \xrightarrow{\text{pr}} \infty$. Condition (a) and (B.11) imply that $\bar{L}(\hat{\psi}^*) - \bar{L}^S(\hat{\psi}^S) = O_p(1/n)$. Similarly, we also have $\bar{L}(\hat{\psi}_0^*) - \bar{L}^S(\hat{\psi}_0^S) = O_p(1/n)$. 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 $\hat{d}_L \simeq \hat{d}_S$ and $\bar{d}_L = \bar{d}_S$, we know $\hat{r}_L \simeq \hat{r}_S$. \square

Note that, even under the assumption of this theorem, \hat{r}_S and \hat{r}_S^\diamond are not equivalent. From (A.7) and (A.8), \hat{r}_S and \hat{r}_S^\diamond are a ‘‘difference of difference’’ estimator and a ‘‘difference’’ estimator, respectively. So, the ‘‘bias’’ of using $\bar{L}^S(\psi)$ cannot be canceled out in \hat{r}_S^\diamond .

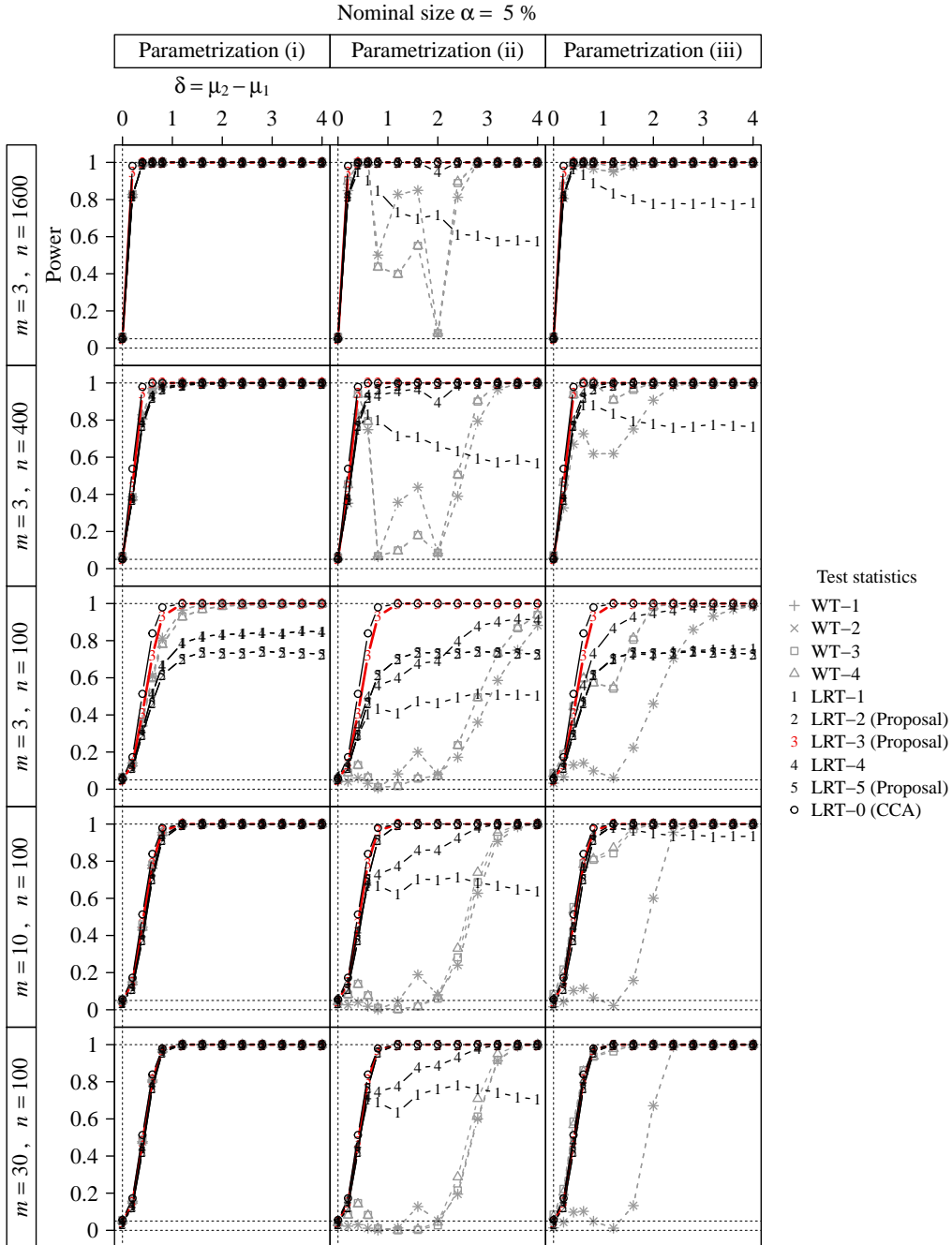


Figure A.7: The power curves under different parametrizations. The nominal size is $\alpha = 5\%$. In each plot, the vertical axis denotes the power, whereas the horizontal axis denotes the value of $\delta = \mu_2 - \mu_1$. The legend in Figure A.5 also applies here.

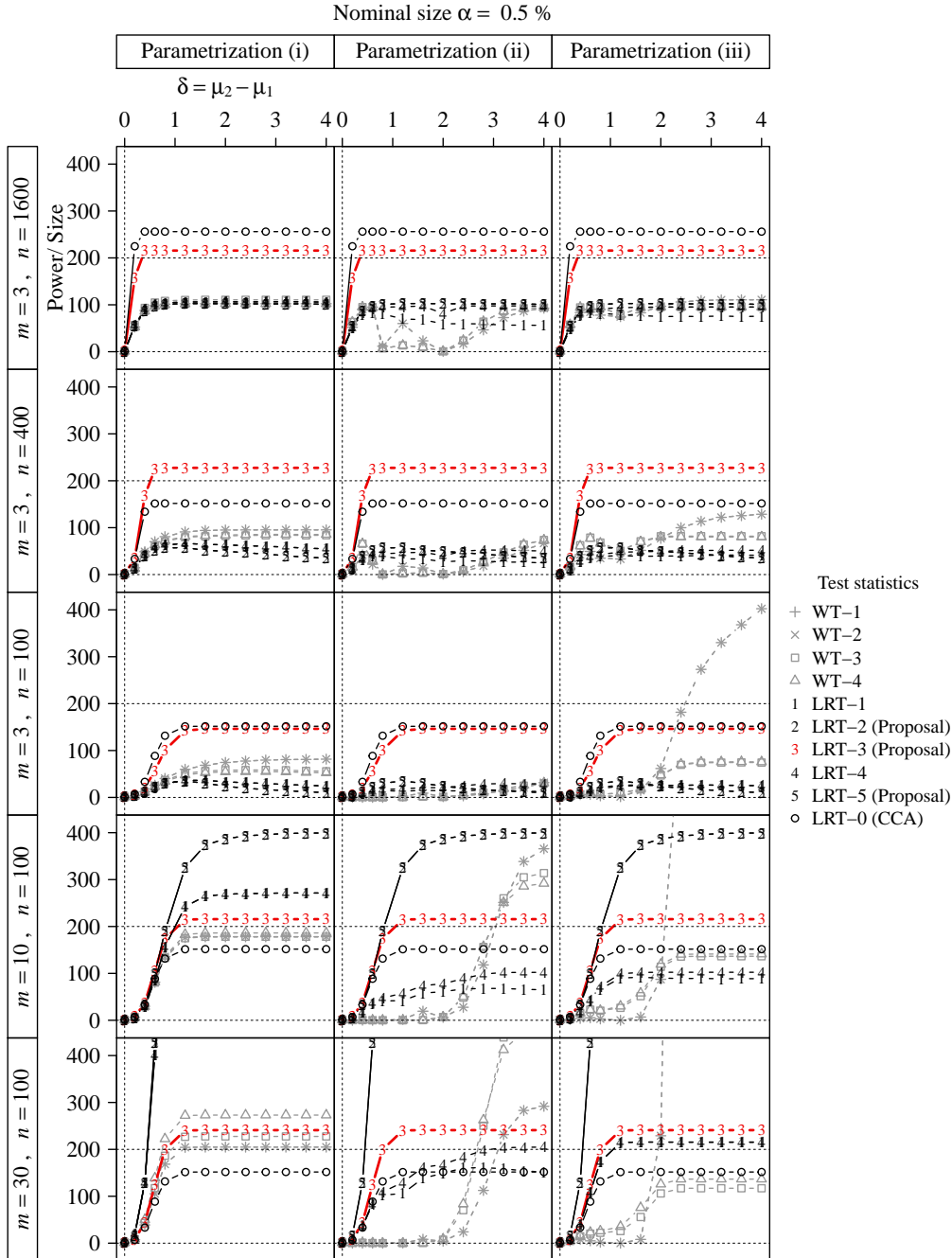


Figure A.8: The ratios of empirical power to empirical size under different parametrizations. The nominal size is $\alpha = 0.5\%$. In each plot, the vertical axis denotes the ratio, and the horizontal axis denotes $\delta = \mu_2 - \mu_1$. The legend in Figure A.5 also applies here. The results under nominal size 5% are shown in Figure A.9.

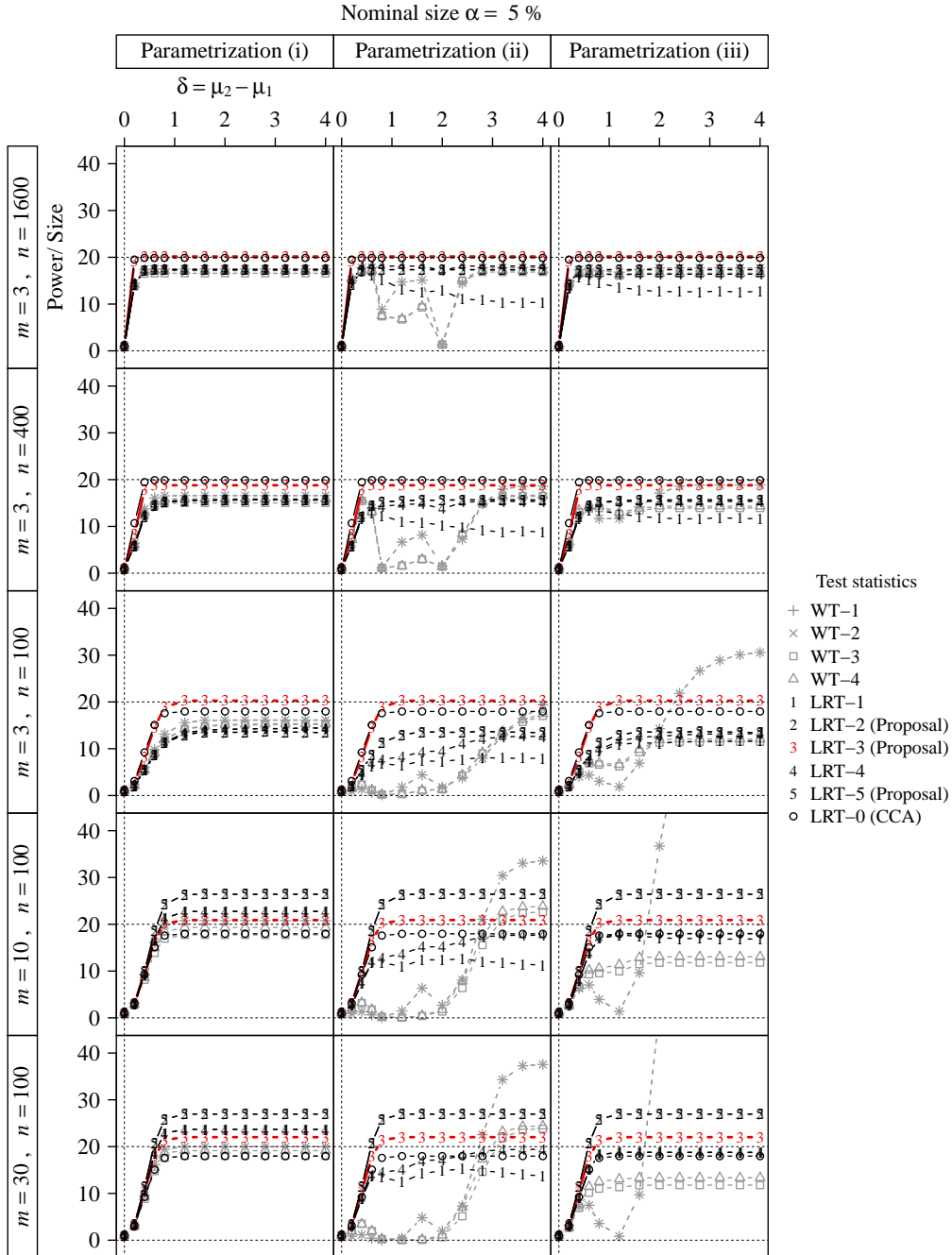


Figure A.9: The ratios of empirical power to empirical size under different parametrizations. The nominal size is $\alpha = 5\%$. In each plot, the vertical axis denotes the ratio, whereas the horizontal axis denotes $\delta = \mu_2 - \mu_1$. The legend in Figure A.5 also applies here.

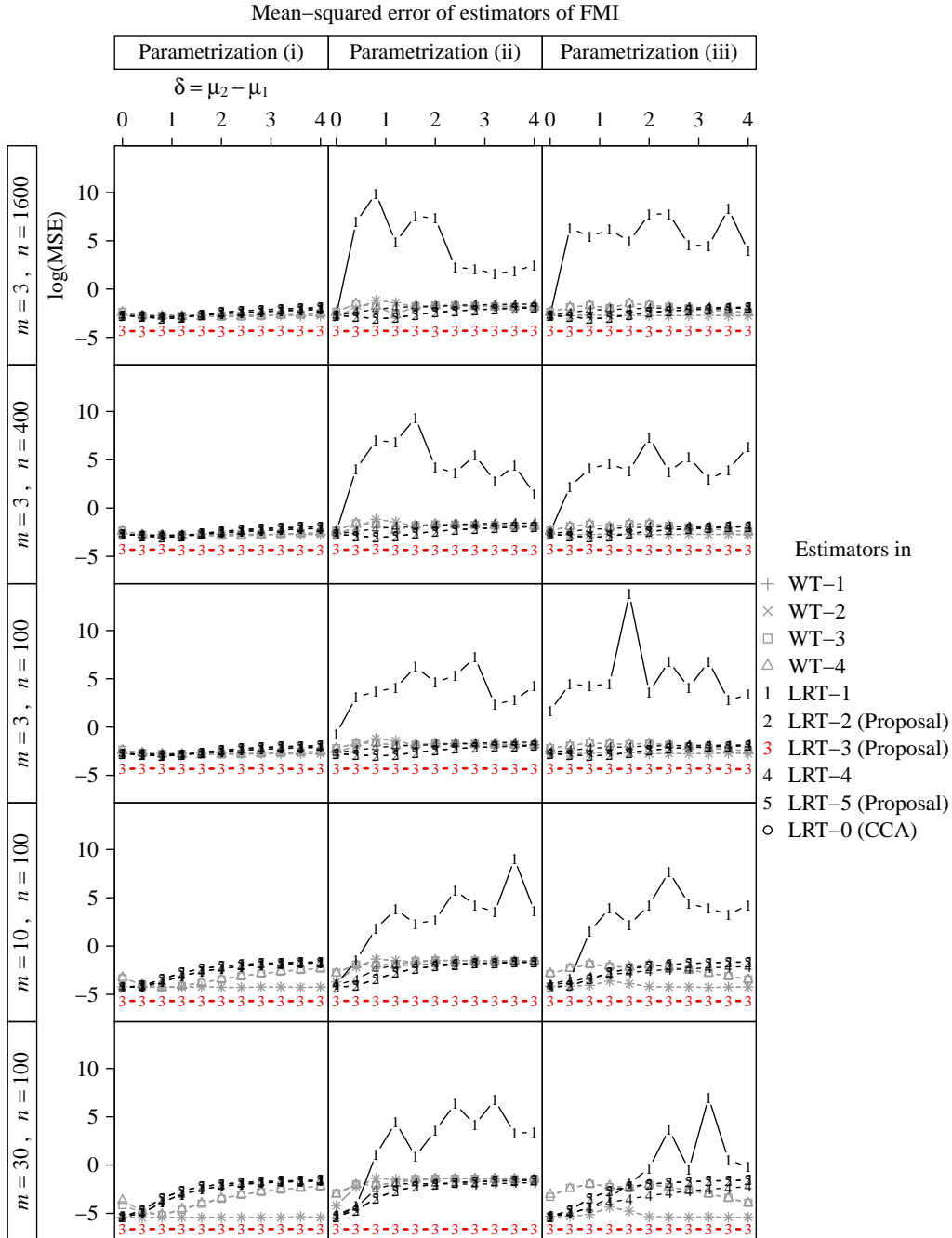


Figure A.10: The MSEs of estimators of ℓ_m used in the test statistics. The vertical axis denotes the log of MSE, whereas the horizontal axis denotes the value of $\delta = \mu_2 - \mu_1$. The legend in Figure A.5 also applies here.

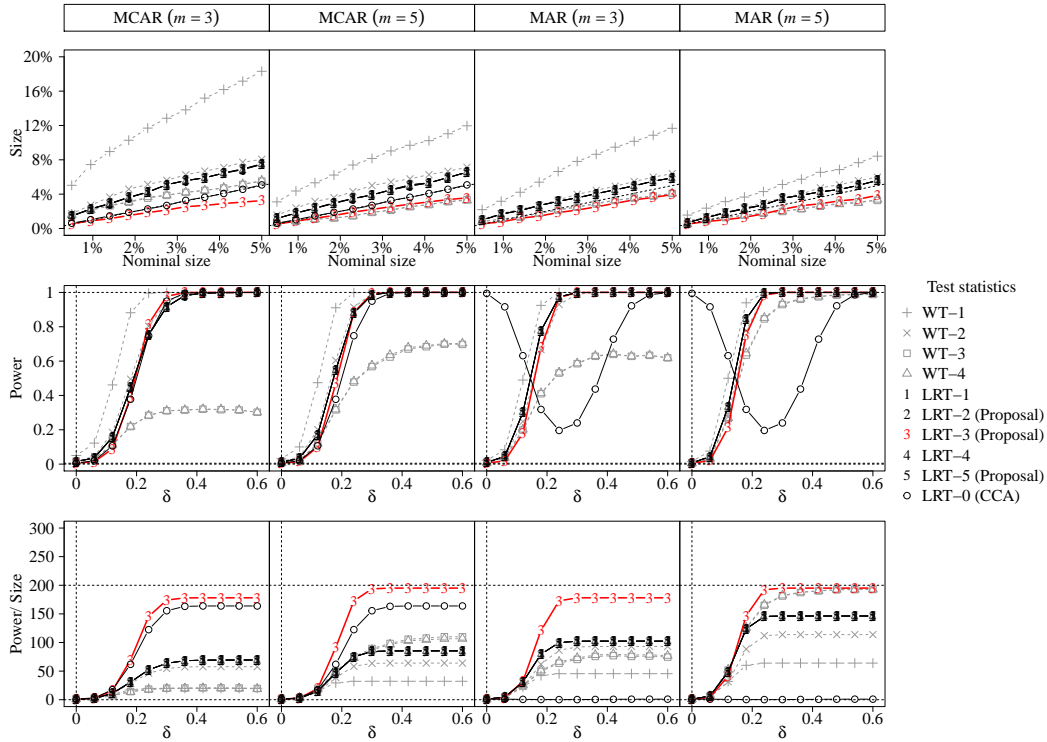


Figure A.11: The empirical size, empirical power, and their ratio. The first row of plots show the empirical sizes. The size of the complete-case test (C2) under MAR is off the chart (always equals to one) because it is invalid. The second and third rows of plots show the powers and the power-to-size ratios, respectively, where the nominal size is 0.5%.

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.
- Cox, D. R. and Reid, N. (1987) Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society B*, **49**, 1–39.
- Li, K. H., Meng, X.-L., Raghunathan, T. E. and Rubin, D. B. (1991) Significance levels from repeated p -values with multiply-imputed data. *Statistica Sinica*, **1**, 65–92.
- Little, R. J. A. and Rubin, D. B. (2002) *Statistical analysis with missing data*. Wiley, New York.
- Meng, X.-L. and Rubin, D. B. (1992) Performing likelihood ratio tests with multiply-imputed data sets. *Biometrika*, **79**, 103–111.
- Rubin, D. B. (2004) *Multiple Imputation for Nonresponse in Surveys*. Wiley, New York.
- van der Vaart, A. W. (2000) *Asymptotic Statistics*. Cambridge University Press, Cambridge.
- Wang, N. and Robins, J. M. (1998) Large-sample theory for parametric multiple imputation procedures. *Biometrika*, **85**, 935–948.