

A SUMMARY OF NOTATIONS

Table 5: Summary of notations

$g(x)$	True density
$a(x)$	Auxiliary density
$b(x)$	Noise distribution
n	Sample size
$\tilde{p}(x; \theta)$	Unnormalized model
$p(x; \theta)$	Normalized model
$q(x; \tau)$	One-parameter extended model
$x_{\text{obs}}, x_{\text{mis}}$	Observed data and missing data
$r(x)$	$q(x; \tau)/a(x)$ or $\tilde{p}(x; \theta)/a(x)$
$\Pr(\delta x, \phi)$	Selection probability
$\pi(\delta x, \phi)$	Propensity score model
η	(θ, ϕ)
ζ	(τ, ϕ)
M_{sc}	Loss function of score matching
Z_{sc}	Estimating equation of score matching
M_{nc}, M_{nc1}, M_{nc2}	Loss function of NCE
Z_{nc}, Z_{nc1}, Z_{nc2}	Estimating equation of NCE
$p(x; \theta)$	Normalized model of $\tilde{p}(x; \theta)$
$t(x_{\text{mis}}; \eta)$	Posterior $p(x; \theta)\pi(\delta x; \phi)$
θ_0	True θ
$x_i^{(*k)}$	Imputed data
$t(x)^{\otimes 2}$	$t(x)t(x)^\top$
$\hat{\eta}_p$	Initial estimator
$\hat{\eta}_{sc}$	Estimator by FISCORE and MISCORE
$\hat{\eta}_{nc,f}$	Estimator by FINCE and MINCE
μ	Baseline measure
$c_s(x; \theta)$	$\nabla_{x^s} \log \tilde{p}(x; \theta)$

B PROOF

To keep the clarity of the main points of this section, we will not specify regularity conditions. For details, see Chapter 5 in van der Vaart (1998).

Proof of Theorem 2 and 1. Direct calculation based on the original theory of M-estimator. ■

Proof of Theorem 3. First, we discuss the general derivation without using a specific form of $z_{sc}(\theta)$ so that it can be applied to NCE case. Then, we derived the specific formula for FISCORE.

we have

$$\bar{Z}_{sc}(\theta|\hat{\theta}_p) = Z_{sc,\text{obs}}(\theta) + \mathbb{E}[Z_{sc,\text{mis}}|\mathbf{x}_{\text{obs}}; \hat{\theta}_p],$$

where $Z_{sc,\text{mis}} = Z_{sc}(\theta) - Z_{sc,\text{obs}}(\theta)$. By Taylor expansion, we have

$$\mathbb{E}[Z_{sc,\text{mis}}|\mathbf{x}_{\text{obs}}; \hat{\theta}_p] = \mathbb{E}[Z_{sc,\text{mis}}(\theta_0)|\mathbf{x}_{\text{obs}}; \theta_0] + \mathbb{E}[Z_{sc,\text{mis}}(\theta)\nabla_{\theta^\top} \log p(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}}; \theta)|\mathbf{x}_{\text{obs}}; \theta_0]_{\theta_0}(\hat{\theta}_p - \theta_0) + o_p(n^{-1/2}).$$

Therefore,

$$\begin{aligned} \bar{Z}_{sc}(\theta_0|\hat{\theta}_p) &= Z_{sc,\text{obs}}(\theta_0) - \mathcal{I}_{2,sc}(\hat{\theta}_p - \theta_0) + o_p(n^{-1/2}) \\ &= -\mathcal{I}_{1,sc}(\hat{\theta}_{sc} - \theta_0) - \mathcal{I}_{2,sc}(\hat{\theta}_p - \theta_0) + o_p(n^{-1/2}) \\ &= (-\mathcal{I}_{1,sc} - \mathcal{I}_{2,sc})(\hat{\theta}_{sc} - \theta_0) - \mathcal{I}_{2,sc}(\hat{\theta}_p - \hat{\theta}_{sc}) + o_p(n^{-1/2}), \end{aligned} \tag{12}$$

where

$$\begin{aligned}
 \mathcal{I}_{1,sc} &= \mathbb{E}[\nabla_{\theta^\top} Z_{sc,\text{obs}}(\theta_0)], \\
 \mathcal{I}_{2,sc} &= -\mathbb{E}[\mathbb{E}[Z_{sc,\text{mis}}(\theta_0) \nabla_{\theta^\top} \log p(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}}; \theta_0)]] \\
 &= -\mathbb{E}[\mathbb{E}[z_{sc,\text{mis}}(\theta_0) \nabla_{\theta^\top} \log p(x_{\text{mis}}|x_{\text{obs}}; \theta_0)]] \\
 &= -\mathbb{E}[z_{sc,\text{mis}}(\theta_0) \nabla_{\theta^\top} \log p(x_{\text{mis}}|x_{\text{obs}})] \\
 &= -\mathbb{E}[\text{cov}[z_{sc,\text{mis}}(\theta_0), \nabla_{\theta^\top} \log \tilde{p}(x_{\text{mis}}|x_{\text{obs}}; \theta_0)]|x_{\text{obs}}; \theta_0].
 \end{aligned}$$

From the first line to the second line (12), we used $\mathbb{E}[Z_{sc,\text{mis}}(\theta_0)|\mathbf{x}_{\text{obs}}; \theta_0] = 0$ and Theorem 2.

In addition, since $\hat{\theta}_{sc,\infty}$ is the solution to $\bar{Z}_{sc}(\theta|\hat{\theta}_p)$. Then,

$$\begin{aligned}
 0 &= \bar{Z}_{sc}(\hat{\theta}_{sc,\infty}|\hat{\theta}_p) \\
 &= \bar{Z}_{sc}(\theta_0|\hat{\theta}_p) + \mathbb{E}[\nabla_{\theta^\top} Z_{sc}(\theta_0)](\hat{\theta}_{sc,\infty} - \theta_0) + o_p(n^{-1/2}) \\
 &= \bar{Z}_{sc}(\theta_0|\hat{\theta}_p) + \mathcal{I}_{3,sc}(\hat{\theta}_{sc,\infty} - \theta_0) + o_p(n^{-1/2}),
 \end{aligned}$$

where

$$\mathcal{I}_{3,sc} = \mathbb{E}[\nabla_{\theta^\top} Z_{sc}(\theta_0)].$$

Therefore, we get

$$\begin{aligned}
 (\hat{\theta}_{sc,\infty} - \theta_0) &= -\mathcal{I}_{3,sc}^{-1} \{(-\mathcal{I}_{1,sc} - \mathcal{I}_{2,sc})(\hat{\theta}_{sc} - \theta_0) - \mathcal{I}_{2,sc}(\hat{\theta}_p - \hat{\theta}_{sc})\} + o_p(n^{-1/2}), \\
 &= (\hat{\theta}_{sc} - \theta_0) + \mathcal{I}_{3,sc}^{-1} \mathcal{I}_{2,sc}(\hat{\theta}_p - \hat{\theta}_{sc}) + o_p(n^{-1/2}).
 \end{aligned}$$

From the first line to the second line of the last equation, we used the relation $\mathcal{I}_{3,sc} = \mathcal{I}_{1,sc} + \mathcal{I}_{2,sc}$. This is proved by

$$\begin{aligned}
 \mathcal{I}_{1,sc} + \mathcal{I}_{2,sc} &= \mathbb{E}[\nabla_{\theta^\top} (\mathbb{E}[Z_{sc}(\theta)|\mathbf{x}_{\text{obs}}; \theta])] - \mathbb{E}[\mathbb{E}[Z_{sc,\text{mis}}(\theta_0) \nabla_{\theta^\top} \log p(\mathbf{x}_{\text{mis}}|\mathbf{x}_{\text{obs}}; \theta_0)|\mathbf{x}_{\text{obs}}; \theta_0]] \\
 &= \mathbb{E}[\nabla_{\theta^\top} Z_{sc}(\theta_0)] = \mathcal{I}_{3,sc}.
 \end{aligned}$$

We go back to the special case of FISCORE.

Noting $m_{sc}(\theta) = \sum_{s=1}^{d_x} 0.5c_s^2(x) + \nabla_{x^s}(c_s(x))$, the term $z_{sc}(\theta)$ is

$$z_{sc}(\theta) = \sum_{s=1}^d \{c_s(x) \nabla_{\theta}(c_s(x)) + \nabla_{x^s}(\nabla_{\theta} c_s(x))\}.$$

Then, we have

$$\begin{aligned}
 \mathbb{E}[\nabla_{\theta^\top} z_{sc,\text{obs}}(\theta)]|_{\theta_0} &= \mathbb{E}[\nabla_{\theta^\top} \{ \mathbb{E}[z_{sc}(\theta)|x_{\text{obs}}; \theta] \}]|_{\theta_0} \\
 &= \mathbb{E}[\nabla_{\theta^\top} z_{sc}(\theta)]|_{\theta_0} + \mathbb{E}[\mathbb{E}[z_{sc}(\theta) \{ \nabla_{\theta^\top} \log p(x_{\text{mis}}|x_{\text{obs}}; \theta) \} | x_{\text{obs}}; \theta_0]]|_{\theta_0} \\
 &= \mathbb{E}[\nabla_{\theta^\top} z_{sc}(\theta)]|_{\theta_0} + \mathbb{E}[z_{sc}(\theta) \{ \nabla_{\theta^\top} \log p(x_{\text{mis}}|x_{\text{obs}}; \theta) \}]|_{\theta_0},
 \end{aligned}$$

where $\nabla_{\theta} \log p(x_{\text{mis}}|x_{\text{obs}}; \theta)$ is

$$\nabla_{\theta} \log \tilde{p}(x; \theta) - \mathbb{E}[\nabla_{\theta} \log \tilde{p}(x; \theta)|x_{\text{obs}}; \theta].$$

So, the above is equal to

$$\mathbb{E}[\nabla_{\theta^\top} z_{sc}(\theta)]|_{\theta_0} + \mathbb{E}[\text{cov}[z_{sc}(\theta), \nabla_{\theta} \log \tilde{p}(x; \theta)|x_{\text{obs}}; \theta]]|_{\theta_0}.$$

In addition,

$$\begin{aligned}
 \mathbb{E}[\nabla_{\theta^\top} z_{sc}(\theta)]|_{\theta_0} &= \mathbb{E} \left[\sum_{s=1}^d \{ \nabla_{\theta} c_s(x) \nabla_{\theta^\top} c_s(x) + c_s(x) \nabla_{\theta\theta^\top} c_s(x) + \nabla_{x^s}(\nabla_{\theta\theta^\top} c_s(x)) \} \right] |_{\theta_0} \\
 &= \mathbb{E} \left[\sum_{s=1}^d \{ \nabla_{\theta} c_s(x) \}^{\otimes 2} \right] |_{\theta_0}.
 \end{aligned}$$

From the second line to the third line, we used a partial integration trick, which is a core concept of score matching. ■

Proof of Corollary 1.

Clear from Theorem 4.

■

Proof of Corollary 2.

First, we calculate $\mathcal{I}_{1,nc}$. By noting the sampling mechanism of full data is a stratified sampling, this is calculated as follows:

$$n^{-1}(\text{var}_q[\mathbb{E}[z_{nc1}(x; \tau_0)|x_{\text{obs}}]] + \text{var}_a[z_{nc2}(y; \tau_0)]).$$

Next, we calculate $\mathcal{I}_{1,nc}$:

$$\begin{aligned} \mathcal{I}_{1,nc} &= \mathbb{E}[\nabla_{\tau^\top} Z_{nc,\text{obs}}(\tau)]|_{\tau_0} = \mathbb{E}[\nabla_{\tau^\top} z_{nc,\text{obs}}(\tau)]|_{\tau_0} = \mathbb{E}[\nabla_{\tau^\top} \{\mathbb{E}[z_{nc}(x, y; \tau)|x_{\text{obs}}; \tau]\}]|_{\tau_0} \\ &= \mathbb{E}[\nabla_{\tau^\top} z_{nc}(x, y; \tau)]|_{\tau_0} + \mathbb{E}[z_{nc}(x, y; \tau) \{\nabla_{\tau^\top} \log \bar{q}(x_{\text{mis}}|x_{\text{obs}}; \tau)\}]|_{\tau_0} \\ &= \mathcal{I}_{3,nc} - \mathcal{I}_{2,nc}. \end{aligned} \tag{13}$$

$$\tag{14}$$

where

$$\bar{q}(x_{\text{mis}}|x_{\text{obs}}; \tau) = q(x_{\text{mis}}, x_{\text{obs}}; \tau) / \int q(x_{\text{mis}}, x_{\text{obs}}; \tau) \mu(\mathrm{d}x_{\text{mis}}).$$

By some algebra, the first term in (14) is

$$\mathcal{I}_{3,nc} = \mathbb{E}[\nabla_{\tau^\top} z_{nc}(x, y; \tau)]|_{\tau_0} = \mathbb{E} \left[\frac{\nabla_{\tau} \log q(x; \tau_0)^{\otimes 2}}{1+r} \right] |_{\tau_0}.$$

In addition, the second term in (14) is

$$\begin{aligned} \mathcal{I}_{2,nc} &= -\mathbb{E}[z_{nc}(x, y; \tau) \{\nabla_{\tau^\top} \log \bar{q}(x_{\text{mis}}|x_{\text{obs}}; \tau)\}]|_{\tau_0} \\ &= -\mathbb{E}[\mathbb{E}[z_{nc1}(x; \tau)|x_{\text{obs}}] \{\nabla_{\tau^\top} \log \bar{q}(x_{\text{mis}}|x_{\text{obs}}; \tau)\}]|_{\tau_0} \\ &= -\mathbb{E}[\text{cov}[z_{nc1}(x; \tau), \nabla_{\tau} \log q(x; \tau)|x_{\text{obs}}]] \\ &= \mathcal{I}_{3,nc} - \mathbb{E} \left[\mathbb{E} \left[\frac{\nabla_{\tau} \log q(x; \tau_0)}{1+r} | x_{\text{obs}} \right] \mathbb{E} [\nabla_{\tau^\top} \log q(x; \tau_0) | x_{\text{obs}}] \right]. \end{aligned}$$

Therefore, adding the first and the second term in (14), we get

$$\mathcal{I}_{1,nc} = \mathbb{E} \left[\mathbb{E} \left[\frac{\nabla_{\tau} \log q(x; \tau_0)}{1+r} | x_{\text{obs}} \right] \mathbb{E} [\nabla_{\tau^\top} \log q(x; \tau_0) | x_{\text{obs}}] \right].$$

■

Proof of Corollary 3. By some algebra, as in the proof of Corollary 2, we obtain

$$\begin{aligned} \mathcal{I}_{1,nc} &= \mathbb{E} \left[\mathbb{E} [\nabla_{\tau} \log q(x; \tau_0) | x_{\text{obs}}]^{\otimes 2} \right], \\ \mathcal{I}_{3,nc} &= \mathbb{E} [\nabla_{\tau} \log q(x; \tau_0)^{\otimes 2}]. \end{aligned}$$

So, noting that $\mathcal{I}_{3,nc}$ is a positive definite matrix, and $\mathcal{I}_{3,nc}$ and $\mathcal{I}_{1,nc}$ are symmetric matrices, we can express $\mathcal{I}_{3,nc} = RR^\top$ and $\mathcal{I}_{1,nc} = R\Lambda R^\top$ using a nonsingular matrix R (Rao, 2008). Because $\mathcal{I}_{3,nc} - \mathcal{I}_{1,nc}$ is a positive matrix from Jensen's inequality, each element in Λ is less than 1. Then, we get

$$\mathcal{I}_{3,nc}^{-1} \mathcal{I}_{2,nc} = \mathcal{I}_{3,nc}^{-1} (\mathcal{I}_{3,nc} - \mathcal{I}_{1,nc}) = R^{-1} (I - \Lambda) R.$$

Finally,

$$(\mathcal{I}_{3,nc}^{-1} \mathcal{I}_{2,nc})^j = R^{-1} (I - \Lambda)^j R.$$

Therefore, $\{\mathcal{I}_{3,nc}^{-1} \mathcal{I}_{2,nc}\}^j$ converges to zero as j tends to infinity. ■

C Multiple missing patterns

We explain how to handle the case when the missing pattern is multiple. In the most general case, we have to introduce a missing pattern indicator δ that takes values in $0, 1, \dots, 2^K - 1$, where K is the dimension of x , for each sample. This is based on the fact that there are 2^K possible missing patterns for x . In the main manuscript, we have defined δ to take a value of 0 or 1 because there are only two missing patterns, i.e., one value is missing or not. Then, we have to introduce an importance distribution $b(x_{mis})$ separately for each missing pattern in the general case. In practice, we can simply select the importance distribution for each coordinate and take their products. Thus, the proposed methods can be applied to the general missing case. For example, in Algorithm 2, we impute missing values of each sample using the importance distribution of corresponding missing pattern. Then, W-step and M-step are essentially the same. Here is a concrete example.

Example C.1 Consider the case $x_{1,obs} = [0, NA, 1]$, $x_{2,obs} = [NA, NA, 4]$, $x_{3,obs} = [1, 2, 3]$.

In this case, there are 3 missing patterns: $(*, *, *)$, $(*, NA, *)$, $(NA, NA, *)$, where $*$ means an observed value and NA means a missing value. Then, we introduce a missing indicator δ to take values in $0, 1, 2$. Namely, $\delta = 0$ corresponds to $(*, *, *)$, $\delta = 1$ corresponds to $(*, NA, *)$, and $\delta = 2$ corresponds to $(NA, NA, *)$. Then, for $x_{1,obs} = [0, NA, 1]$, we have $\delta_1 = 1$. For $x_{2,obs} = [NA, NA, 4]$, we have $\delta_2 = 2$. Since x_3 is observed without missing, then $\delta_3 = 0$. (6) is rewritten:

$$\frac{1}{n} \sum_{i=1}^n \sum_{k=0}^2 \mathbb{I}(\delta_i = k) \mathbb{E}[z_{nc}(x_i; \tau) \mid x_{i,obs}; \theta].$$

D COMPARISON BETWEEN FINCE and VNCE

Here, we compare FINCE and variational NCE (VNCE) (Rhodes and Gutmann, 2019). From (7), the difference between the estimator proposed in this paper and VNCE (Rhodes and Gutmann, 2019) is clearly shown. Mainly, there are two differences: (1) VNCE attempts to maximize the observed likelihood directly, whereas FINCE attempts to solve the observed estimating equation, (2) VNCE assumes that the dimension of $a(x)$ is the same as the dimension of x_{obs} , whereas FINCE assumes that the dimension of $a(x)$ is the same as the dimension of x .

More specifically, an ideal loss function in VNCE is

$$\begin{aligned} \arg \max_s J_{VNCE}(\tau, s(x_{mis})) &= J_{VNCE}(\tau, q(x_{mis} \mid x_{obs})) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\log \left\{ \frac{1}{1 + \frac{q(x_{mis} \mid x_{i,obs}) a(x_{i,obs})}{q(x_{mis}, x_{i,obs})}} \right\} \mid x_{i,obs} \right] + \frac{1}{n} \sum_{j=1}^n \log \left\{ \frac{a(y_j)}{a(y_j) + \mathbb{E}[q(y_j, y_{j,mis}) \mid y_j]} \right\} \\ &= \frac{1}{n} \sum_{i=1}^n \log \left\{ \frac{q(x_{i,obs})}{q(x_{i,obs}) + a(x_{i,obs})} \right\} + \frac{1}{n} \sum_{j=1}^n \log \left\{ \frac{a(y_j)}{a(y_j) + q(y_j)} \right\}, \end{aligned} \quad (15)$$

where $q(x_{obs}) = \int q(x_{mis}, x_{obs}) \mu(dx_{mis})$. On the other hand, the loss function of our proposed estimator is (7). In general, the efficiencies of the two loss function are not directly comparable. The following points highlight the comparison between two methods.

- In terms of inferences, our proposed methods (FINCE, FISCORE) are superior to VNCE because it is difficult to achieve the upper bound in (15) in VNCE.
- Unless the family of variational distribution includes the true posterior, VNCE does not have consistency. On the other hand, FINCE has consistency and also asymptotic normality without requiring such conditions.
- In terms of the scalability, VNCE is superior to the proposed methods because VNCE does not require any sampling methods.
- FINCE can be applied even if the missing data mechanism is MAR or MNAR. However, VNCE cannot be directly applied when the missing mechanism is MAR or NMAR.

E INFERENCE OF FISCORE WHEN m IS FIXED

We consider an asymptotic result of FISCORE when m is fixed. Actually, the estimating equation $Z_{sc,m}$ is not unbiased estimator for \bar{Z}_{sc} because a self normalizing importance sampling is used rather than importance sampling (Owen, 2013). This means that the derived estimator is theoretically not consistent; however, practically, a self normalized importance sampling is preferable to importance sampling because of its robustness. Here, we consider the case when the weight is defined as $w(x|x_{\text{obs}}) = p(x_{\text{mis}}|x_{\text{obs}}; \theta_0)/b(x)$.

As in the proof of Theorem 3, we have

$$\hat{\theta}_{sc,m} - \theta_0 = -\mathcal{I}_{3,sc}^{-1} Z_{sc,m}(\theta_0|\hat{\theta}_p) + o_p(n^{-1/2}) \quad (16)$$

This term is decomposed into two terms: $-\mathcal{I}_{3,sc}^{-1} \bar{Z}_{sc}(\theta_0|\hat{\theta}_p)$ and $-\mathcal{I}_{3,sc}^{-1} \{Z_{sc,m}(\theta_0|\hat{\theta}_p) - \bar{Z}_{sc}(\theta_0|\hat{\theta}_p)\}$. These two terms in (16) are independent. The first term is equal to $\hat{\theta}_{sc,\infty} - \theta_0$, of which the asymptotic property is shown in Theorem 3. The second term converges to the normal distribution with mean 0 and variance $\mathcal{I}_{3,sc}^{-1} \mathbb{E}[\text{var}_b\{Z_{sc,m}(\theta_0|\hat{\theta}_p)\}]\mathcal{I}_{3,sc}^{-1}$.

Theorem 4 When $\hat{\theta}_p = \hat{\theta}_{sc}$, the asymptotic variance of $\hat{\theta}_{sc,m}$ is equal to

$$\mathcal{I}_{1,sc}^{-1} \mathcal{J}_{1,sc} \mathcal{I}_{1,sc}^{\top-1} + m^{-1} \mathcal{I}_{3,sc}^{-1} \mathcal{J}_{2,sc} \mathcal{I}_{3,sc}^{\top-1},$$

where $w(x|x_{\text{obs}}) = p(x_{\text{mis}}|x_{\text{obs}}; \theta_0)/b(x)$ and

$$\mathcal{J}_{2,sc} = n^{-1} \mathbb{E}[\mathbb{E}_{b(x_{\text{mis}})}[w^2(x)\{z_{sc}(\theta_0)^{\otimes 2}\}|x_{\text{obs}}]] - n^{-1} \mathbb{E}[\mathbb{E}[z_{sc}(\theta_0)|x_{\text{obs}}]^{\otimes 2}].$$

F EXTENSION TO MULTIPLE IMPUTATION: MISCORE AND MINCE

MI was originally developed with Bayesian flavor (Rubin, 1987; Meng, 1994). In this paper, we consider frequentist MI rather than Bayesian MI (Tsiatis, 2006) to avoid the additional computation. In addition, it is shown that frequentist MI is asymptotically more efficient than Bayesian MI (Wang and Robins, 1998; Robins and Wang, 2000).

In MI, the crucial assumption is that the sample can be obtained from $p(x_{\text{mis}}|x_{\text{obs}}; \theta)$. When the missing data mechanism is MAR, it is easy to sample from $p(x_{\text{mis}}|x_{\text{obs}}; \theta)$ using the MCMC based on (5). The algorithm is described as in Algorithm 4. In this paper, this approach is referred to as MISCORE. MINCE is also defined similarly. Nevertheless, we do not recommend Algorithm 4 for the practical reason of its instability and computational burden.

Algorithm 4: MISCORE

1 **repeat**

2 W-step: Take a set of m samples from $x_{\text{mis}}^{*k} \sim p(x_{\text{mis}}|x_{\text{obs}}; \hat{\theta}_t)$ using MCMC for each i

3 M-step: Update the solution to the following function with respect to θ as $\hat{\theta}_{t+1}$:

$$\frac{1}{nm} \sum_{i=1}^n \sum_{k=1}^m m_{sc}(x_i^{*k}; \theta).$$

4 **until** $\hat{\tau}_t$ converges;

Dues to the challenges associated with Algorithm 4, we recommend the following algorithm. This algorithm is similar to the one in Levine and Casella (2001). In the original MISCORE, a set of samples is generated at every step. This requires tremendous computational cost and causes instability. In Algorithm 5, by constructing a \sqrt{n} -consistent estimator based on FISCORE at each step and updating by MISCORE one time, this limitation is overcome.

Algorithm 5: One step MISCORE

1 **repeat**

2 Do W-step and M-step in Algorithm 2 (FISCORE)

3 **until** $\hat{\tau}_t$ converges;

4 Do W-step and M-step in Algorithm 4 (MISCORE)

Table 6: Monte Carlo median square error and bias

Setting 1			
n		MINCE	MISCORE
500	(bias)	0.15	0.10
	(mse)	0.037	0.024
1000	(bias)	0.13	0.12
	(mse)	0.030	0.020

Table 6 illustrates the experimental result. We generated a set of 50 samples for each i using MCMC in the last step. Compared with FINCE and FISCORE, the performance of one step MISCORE is worse. Perhaps, more step is needed.

The asymptotic property is obtained as follows.

Corollary 4 When $\hat{\theta}_p = \hat{\theta}_{sc}$ and m is fixed, the asymptotic variance of $\hat{\theta}_{sc,\infty}$ is equal to

$$\mathcal{I}_{1,sc}^{-1} \mathcal{J}_{1,sc} \mathcal{I}_{1,sc}^{\top-1} + m^{-1} \mathcal{I}_{3,sc}^{-1} \mathcal{J}_{2,sc} \mathcal{I}_{3,sc}^{\top-1},$$

where

$$\mathcal{J}_{2,sc} = n^{-1} \{E[z_{sc}(\theta_0)^{\otimes 2}] - E[E[z_{sc}(\theta_0)|x_{\text{obs}}]^{\otimes 2}]\},$$

and other terms are the same as in Theorem 3.

Proof of Corollary 4. We just replace $b(x_{\text{mis}})$ with $p(x_{\text{mis}}|x_{\text{obs}}; \theta_0)$ in Theorem 4. ■

Finally, there are two things to note about MISCORE and MINCE. When the missing data mechanism is MNAR, we have to sample from $\tilde{p}(x_{\text{mis}}|x_{\text{obs}}, \delta; \eta) \propto \tilde{p}(x_{\text{mis}}|x_{\text{obs}}; \theta) \pi(\delta|x_{\text{mis}}, x_{\text{obs}}; \phi)$. In this case, the distribution becomes a doubly-intractable distribution (Miller et al., 2006; Murray et al., 2006), and it is generally difficult to sample. Secondly, when we use a Bayesian multiple imputation assuming the prior distribution $\rho(\theta)$, even if the missing mechanism is MAR, we have to sample from $\tilde{p}(x_{\text{mis}}, \theta|x_{\text{obs}}) \propto \tilde{p}(x_{\text{mis}}, x_{\text{obs}}; \theta) \rho(\theta)$. Often, data augmentation is utilized for this purpose (Tanner and Wong, 1987). However, even if the data augmentation is applied, we still have to deal with doubly-intractable distributions to calculate $\Pr(\theta|x) \propto \rho(\theta)p(x; \theta)$.

G EXTENSION TO CONTRASTIVE DIVERGENCE METHODS

Although there are several variations of contrastive divergence methods (Younes, 1989; Tieleman, 2008), the basic idea is that θ is updated by adding the gradient of log-likelihood $\log p(\mathbf{x}; \theta)$ with respect to θ :

$$\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \log \tilde{p}(x_i; \theta) - E_{p(x; \theta)}[\nabla_{\theta} \log \tilde{p}(x; \theta)],$$

multiplying some learning rate. When some data is not observed, the expected gradient becomes

$$\frac{1}{n} \sum_{i=1}^n E[\nabla_{\theta} \log \tilde{p}(x_i; \theta)|x_{i,\text{obs}}; \theta] - E[\nabla_{\theta} \log \tilde{p}(x; \theta)].$$

The expectation of the first term is taken under $p(x_{\text{mis}}|x_{\text{obs}}; \theta)$. It is possible to sample from MCMC like (5) without involving doubly-intractable distributions (Miller et al., 2006). Therefore, the gradient is approximated as

$$\frac{1}{nm} \sum_{i=1}^n \sum_{k=1}^m \nabla_{\theta} \log \tilde{p}(x_i^{*k}; \theta) - \frac{1}{n} \sum_{j=1}^n \nabla_{\theta} \log \tilde{p}(y_j; \theta),$$

where $x_i^{*k} \sim p(x_{\text{mis}}|x_{i,\text{obs}}; \theta)$ and $y_j \sim p(y; \theta)$. We refer the updating method using the above gradient as MICD.

We can still use a FI approach for the approximation. By introducing an auxiliary distribution with a density $b(x)$, the gradient is approximated as

$$\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^m w_{ik} \nabla_{\theta} \log \tilde{p}(x_i^{*k}; \theta) - \frac{1}{n} \sum_{j=1}^n \nabla_{\theta} \log \tilde{p}(y_j; \theta).$$

where $x_i^{*k} \sim b(x)$, $w_{ik} \propto \tilde{p}(x_i^{*k}; \theta)/b(x_i^{*k})$, $y_j \sim p(y; \theta)$. We refer this approach to FICD.

Furthermore, by introducing a noise distribution with a density $a(y)$ to prevent using MCMC totally, the gradient is approximated as

$$\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^m w_{ik} \nabla_{\theta} \log \tilde{p}(x_i^{*k}; \theta) - \frac{1}{n} \sum_{j=1}^n r_j \nabla_{\theta} \log \tilde{p}(y_j; \theta),$$

where $x_i^{*k} \sim b(x)$, $w_{ik} \propto \tilde{p}(x_i^{*k}; \theta)/b(x_i^{*k})$, $y_j \sim a(y)$, and $r_j \propto \tilde{p}(y_j; \theta)/a(y_j)$. In this case, the gradient is essentially equivalent to the loss function of FINCE when $f(x) = x \log x$ by profiling-out c .

H DETAILED ALGORITHM OF FINCE WITH MNAR DATA

The algorithm is described as in Algorithm 6.

Algorithm 6: FINCE with MNAR data

- 1 Initialize $t = 0$, $\hat{\zeta}_0 = (\hat{c}_0, \hat{\theta}_0, \hat{\phi}_0)$
- 2 Take n samples $\{y_j\}_{j=1}^n$ from $a(y)$.
- 3 For i with $\delta_i = 0$, take m samples $\{x_{i,\text{mis}}^{*k}\}_{k=1}^m$ from $b(x)$.
- 4 For i with $\delta_i = 1$, set m samples $\{x_i^{*k}\}_{k=1}^m$ to $x_i^{*k} = x_i$
- 5 **repeat**
- 6 **W-Step:**
- 7 For i with $\delta_i = 0$; $w_{ik} \propto q(x_i^{*k}; \hat{\tau}_t) \pi(\delta_i | x_i^{*k}; \hat{\phi}_t) / b(x_{i,\text{mis}}^{*k})$,
- 8 For i with $\delta_i = 1$; $w_{ik} = 1/m$.
- 9 **M-step:** Solve the following equation for $\hat{\zeta}_{t+1}$ w.r.t ζ :

$$\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^m w_{ik} z_{nc1}(x_i^{*k}; \tau) + Z_{nc2}(\mathbf{y}; \zeta) = 0,$$

$$\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^m \nabla_{\phi} w_{ik} \log \pi(\delta_i | x_i^{*k}; \phi) = 0.$$

- 10 $t = t + 1$
 - 10 **until** $\hat{\tau}_t$ converges;
-

I VARIANCE ESTIMATORS OF FISCORE AND FINCE

I.1 FISCORE

The variance estimator of FISCORE in the case of Corollary is defined as follows: $\hat{\mathcal{L}}_{1,sc}^{-1} \hat{\mathcal{J}}_{1,sc} \hat{\mathcal{L}}_{1,sc}^{\top -1} |_{\hat{\theta}}$, where

$$\begin{aligned}\hat{\mathcal{L}}_{1,sc} &= \frac{1}{n} \sum_{i=1}^n \{ \hat{\mathcal{L}}_{1,sc1}(x_{i,\text{obs}}) + \hat{\mathcal{L}}_{1,sc2}(x_{i,\text{obs}}) - \hat{\mathcal{L}}_{1,sc3}(x_{i,\text{obs}}) \}, \\ \hat{\mathcal{L}}_{1,sc1}(x_{i,\text{obs}}) &= \sum_{k=1}^m w(x_i^{*k}; \theta) \left(\sum_{s=1}^d \nabla_{\theta} c_s(x_i^{*k}) \right)^{\otimes 2}, \\ \hat{\mathcal{L}}_{1,sc2}(x_{i,\text{obs}}) &= \sum_{k=1}^m w(x_i^{*k}; \theta) z_{sc}(x_i^{*k}) \nabla_{\theta^{\top}} \log \tilde{p}(x_i^{*k}; \theta), \\ \hat{\mathcal{L}}_{1,sc3}(x_{i,\text{obs}}) &= \left(\sum_{k=1}^m w(x_i^{*k}; \theta) z_{sc}(x_i^{*k}) \right) \left(\sum_{k=1}^m w(x_i^{*k}; \theta) \nabla_{\theta^{\top}} \log \tilde{p}(x_i^{*k}; \theta) \right), \\ \hat{\mathcal{J}}_{1,sc} &= \frac{1}{n(n-1)} \sum_{i=1}^n \left(\sum_{k=1}^m w(x_i^{*k}; \theta) z_{sc}(x_i^{*k}; \theta) - \bar{z} \right)^{\otimes 2}, \\ \bar{z} &= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^m w(x_i^{*k}; \theta) z_{sc}(x_i^{*k}; \theta), \\ z_{sc}(\theta) &= \sum_{s=1}^d \{ c_s(x) \nabla_{\theta} (c_s(x)) + \nabla_{x^s} (\nabla_{\theta} c_s(x)) \}, \quad c_s(x; \theta) = \nabla_{x^s} \log \tilde{p}(x; \theta).\end{aligned}$$

Next, consider an loss function and a variance estimator in truncated exponential family cases (Hyvärinen, 2007). Assume that $\tilde{p}(x; \theta)$ is given by

$$\log \tilde{p}(x; \theta) = \sum_{k=1}^d \theta_k F_k(x).$$

Let us denote two matrices: $d_{\theta} \times d$ matrix $K_1(x)$ with elements $\nabla_{x^b} F_a$ ($1 \leq a \leq d_{\theta}, 1 \leq b \leq d$) and $d_{\theta} \times 1$ matrix, $K_{i,2}(x)$ with elements $\nabla \nabla_{x^i} F_a$ ($1 \leq a \leq d_x$). The loss function is written as $n^{-1} \sum_{i=1}^n z_{sc,t}(x_i; \theta)$.

$$z_{sc,t}(x; \theta) = 0.5 \theta^{\top} K_1(x) K_1(x)^{\top} \theta + \theta^{\top} \sum_{i=1}^{d_x} K_{i,2}(x).$$

The variance estimator is obtained almost in the same by replacing $z_{sc}(x)$ with $z_{sc,t}(x)$. The only modification is

$$\hat{\mathcal{L}}_{1,sc1}(x_{i,\text{obs}}) = \sum_{k=1}^m w(x_i^{*k}; \theta) K_1(x_i^{*k}) K_1(x_i^{*k})^{\top}.$$

I.2 FINCE

The variance estimator of FINCE in the case of Corollary 2 is defined as follows: $\hat{\mathcal{I}}_{1,nc}^{-1} \hat{\mathcal{J}}_{1,nc} \hat{\mathcal{I}}_{1,nc}^{-1} |_{\hat{\tau}}$, where

$$\begin{aligned} \hat{\mathcal{I}}_{1,nc} &= \frac{1}{n} \sum_{i=1}^n \left(\sum_{k=1}^m w(x_i^{*k}; \tau) \frac{\nabla_{\tau} \log q(x_i^{*k}; \tau)}{1 + q(x_i^{*k})/a(x_i^{*k})} \right) \left(\sum_{k=1}^m w(x_i^{*k}; \tau) \nabla_{\tau^{\top}} \log q(x_i^{*k}; \tau) \right), \\ \hat{\mathcal{J}}_{1,nc} &= \hat{\mathcal{J}}_{1,nc1} + \hat{\mathcal{J}}_{1,nc2}, \\ \hat{\mathcal{J}}_{1,nc1} &= \frac{1}{n(n-1)} \sum_{i=1}^n \left(\sum_{k=1}^m w(x_i^{*k}; \tau) z_{nc1}(x_i^{*k}; \tau) - \bar{z}_{nc1} \right)^{\otimes 2}, \\ \hat{\mathcal{J}}_{1,nc2} &= \frac{1}{n(n-1)} \sum_{i=1}^n (z_{nc2}(y_i; \tau) - \bar{z}_{nc2})^{\otimes 2}, \\ \bar{z}_{nc1} &= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^m w(x_i^{*k}; \tau) z_{nc1}(x_i^{*k}; \tau), \quad \bar{z}_{nc2} = \frac{1}{n} \sum_{j=1}^n z_{nc2}(y_j; \tau), \\ z_{nc1}(\tau) &= -\frac{\nabla_{\tau} \log q(x; \tau)}{1 + r}, \quad z_{nc2}(\tau) = \frac{r \nabla_{\tau} \log q(x; \tau)}{1 + r}. \end{aligned}$$