

Statistica Sinica Preprint No: SS-2017-0069

Title	Rejoinder: Please Visit the Wild Arboretum of Multi-Phase Inference
Manuscript ID	SS-2017-0069
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202017.0069
Complete List of Authors	Xiao-Li Meng
Corresponding Author	Xiao-Li Meng
E-mail	meng@stat.harvard.edu

Rejoinder: Please Visit the Wild Arboretum of Multi-Phase Inference¹

Xianchao Xie and Xiao-Li Meng

1 This was not an easy article to write or to publish. As with most statis-
2 tical theory, the difficulty was not in proving theorems, but in formulating
3 the relevant ones that can convey statistical insights and provide practical
4 guidelines. A further challenge for multi-phase inference lies in finding the
5 most intuitive and simplest ways to illustrate and explain the intricate re-
6 lationships among different phases and their consequences, especially those
7 that are counter-intuitive. It therefore took us a while to pave an entry path
8 into the multi-phase forest, and it took even longer for us to convince enough
9 visitors that it is not a dangerous jungle but rather a wild arboretum with
10 many flowers and fruits, some of which are rather low-hanging.

11 We are therefore very grateful to the editors of *Statistica Sinica* for or-
12 ganizing a general tour of this relatively new landscape of statistical founda-
13 tion, and to our eight brave VIPs (Very Insightful Participants) of the tour.
14 Judging from their comments, we see that we have had a mixed success (or
15 failure) in our attempt to provide an informative and enticing tour guide.
16 Some shared our desire to greatly explore this landscape because the current
17 single-phase theory does not address the increasingly common multi-phase
18 reality. We particularly thank Banks-Peña, Draper and Reiter for their
19 endorsements with additional examples going beyond the multiple imputa-
20 tion setting. Others indicated that we need to do a better job to spell out
21 the practical relevance of our findings (e.g., Little-Zhou) and to demystify

¹We thank NSF and JTF for partial financial support, and Steven Finch for proofreading.

22 the complex world with multiple Gods and parties (Desmond and Raghu-
23 nathan). Below, by addressing some major points raised by the VIPs, we
24 hope to improve and enhance our tour guide, although we are mindful that
25 multi-phase reality will always be more complex than any single brochure
26 can possibly capture.

27 **1. How valid is our concept of validity?**

28 Several discussants (e.g., Desmond, Draper, Yang-Kim) raised the ques-
29 tion of the usefulness of the concept of *confidence validity*, which permits a
30 confidence procedure to cover more than its nominal coverage. We partic-
31 ularly thank Desmond for a very helpful investigation of the origin of, and
32 possible motivations for, Neyman's definition of this concept; Draper's per-
33 sonal touch, being one of Neyman's academic grandsons, is also appreciated.
34 We also agree with Desmond that historically the allowance for over-coverage
35 was mostly motivated by its mathematical convenience to deal with discrete-
36 ness. Nevertheless, it reflected an implicit preference of Neyman and many of
37 our founding generations to rather err on over-covering than under-covering.
38 For practical purposes, it is trivial to come up with many examples of harm-
39 ful consequences of being either overly confident or inadequately confident.
40 However, statistical inference is *not* a symmetric game. It is a game of
41 exclusion and contradiction, not inclusion or confirmation.

42 Regardless of whether our information comes from a (reliable) prior or
43 data or both, we use inferential tools to *sharpen our inference*. That is, we
44 reduce the region of plausible states of our inferential target by excluding
45 those pre-inferential states that are now deemed to be implausible because
46 they have reached a critical level of conflict with available information, as
47 determined by a criterion specified by our inferential procedure. From this
48 perspective, over-covering is simply a necessary step to ensure that the actual
49 exclusion criterion used is in itself not in conflict with what is called for by
50 our procedure. That is, we exclude a target state only when we are *sure*
51 that it has satisfied the exclusion criteria we adopted; otherwise we have to
52 give it the benefit of doubt. Over-covering is therefore not as much an issue

53 of being conservative, but rather a means to ensure rigorousness and hence
54 replicability.

55 Indeed, the consideration of replicability of research is a compelling rea-
56 son to prefer overestimating the uncertainties in our inference, which typi-
57 cally implies over-coverage, than underestimating them, when the exact as-
58 sessment (and hence exact coverage) cannot be achieved. Exact assessment,
59 such as under perfect normality, is never achieved in practice – just con-
60 sidering all kinds of errors and approximations we make, from data defects
61 to modeling frailties to computational corner-cutting. Much of the current
62 crisis of non-replicable research in sciences, especially in the medical, life
63 and social sciences, is due to our asymmetric incentive system, which effec-
64 tively encourage researches to rush into “discoveries” based on quantitative
65 evidence that does not stand up to scrutiny. Ignoring or under-assessing un-
66 certainties, due to a whole host of mishandling, e.g., selection bias, multiple
67 comparisons, over-fitting, etc., is a common cause.

68 Handling model mis-specifications, for which uncongeniality can be viewed
69 as a special (though unavoidable) case, via variance doubling is not a univer-
70 sal recipe. But it is the simplest and most applicable way of combating the
71 common tendency of underestimating actual uncertainty, leading to many
72 falsely significant results. Surely there will be cases where variance dou-
73 bling can result in missed opportunities, due to loss of power, for example,
74 as in Yang-Kim’s simulation. This could have rather serious consequences,
75 such as a delayed release of life-saving medication, and therefore we must
76 be particularly cautious of applying it in those cases where false negative
77 has graver consequences than false positive. Nevertheless, if a more sophis-
78 ticated *and* justifiable approach is unavailable or non-implementable, then
79 variance doubling is likely the lesser of the evils, the other being to ignore
80 the issue (say) un-congeniality all together. This is because variance dou-
81 bling has the added benefit of at least partially “covering” the omissions of
82 other kinds, such as failing to take into account model uncertainty.

83 Raghunathan raised a deeper question about validity for multi-phase in-
84 ferences, especially in the context of multiple imputation for public data files,

85 where there are potentially many analysts. Even assuming every analyst is
86 perfectly trained to do absolutely the best job based on the information
87 s/he has, we still have many model classes to contemplate and each one can
88 lead to its own version of validity, as Raghunathan's " x -analyst" example
89 illustrates, where x can take on many values. Which validity were/are we
90 talking about then?

91 As argued in Liu and Meng [1], to define validity meaningfully we first
92 need to determine the relevant replication setting, over which we can then
93 determine whether some properties are replicable. In a multi-phase set-
94 ting especially with multiple analysts, there are multiple ways of defining
95 meaningful replications, including the marginal, conditional, and joint ones
96 articulated by Raghunathan. Furthermore, shall we treat (some of) the pre-
97 analysis phases, such as an imputation phase, fixed, or should it be a part
98 of our replications? As we argued in the paper, whereas it is natural to
99 consider all kinds of replications, currently we are able only to obtain useful
100 theory under the "grand replications", that is, with respect to God's model
101 that generates the variations for all phases. Theories under more restrictive
102 replications, especially permitting mis-specifications, are challenging. But
103 we hope the more challenging a problem might be, the more enticing it is
104 for adventurous minds.

105 **2. How efficient is our formulation of self-efficiency?**

106 Yang-Kim is correct that self-efficiency can easily be violated by very
107 common procedures, such as ordinary least squares (applied to heteroscedas-
108 tic models), as we demonstrated in the on-line supplementary appendix, bor-
109 rowing an example from Meng and Xie [2]. Yang-Kim is also correct that
110 when self-efficiency is violated, it is possible to recast the problem as an un-
111 congeniality issue, because the latter is formulated via model embedding. A
112 *self-inefficient* procedure with respect to one model can be self-efficient with
113 respect to another; the *model* here includes both the process that generates
114 the original complete data and the missing-data mechanism.

115 The examples in Yang-Kim also provide a good demonstration of the
116 need to be explicit about the procedure being evaluated and with respect to
117 what models—or more generally *replications* (see Liu and Meng [1])—the
118 evaluation is made. If we understand the notation in Yang-Kim correctly,
119 we surmise that the commonality of their three examples is as follows. We
120 have i.i.d. triplets $\{(Y_i, X_i, R_i)\}_{i=1}^n$, where Y_i is the outcome subject to
121 missingness, X_i is the covariate, which is always observed, and R_i is the
122 missing-data indicator, taking value one when Y_i is (fully) observed and
123 zero otherwise. Our estimand θ is the marginal mean of $g(Y)$ for some
124 pre-specified g , and our estimator is the simple average over the observed
125 sample:

$$\hat{\theta}_{\text{obs}} = \frac{\sum_{i=1}^n R_i g(Y_i)}{\sum_{i=1}^n R_i}. \quad (1)$$

126 We emphasize that the concept of self-efficiency is defined for the *observed-*
127 *data procedure*, not the *complete-data procedure*, as stated in Yang-Kim,

$$\hat{\theta}_{\text{com}} = \frac{1}{n} \sum_{i=1}^n g(Y_i), \quad (2)$$

128 because $\hat{\theta}_{\text{obs}}$ trivially specifies $\hat{\theta}_{\text{com}}$ as a special case when all $R_i = 1$, but
129 clearly not vice versa.

130 The usefulness of $\hat{\theta}_{\text{obs}}$ as defined in (1) is well-known to depend on the
131 missing data mechanism (MDM). Yang-Kim invoked the safe assumption of
132 MAR, but upon checking the cited article by Yang and Kim [3], it seems
133 Yang-Kim’s assumption is a more restrictive (but common) one, that is, Y_i
134 and R_i are conditionally independent given the covariate X_i , for all $i =$
135 $1, \dots, n$. Under such an assumption, it is easy to show that $\hat{\theta}_{\text{obs}}$ is unbiased
136 for θ , and we can rely on the asymptotic result given by Theorem 4 of
137 our paper to determine the self-efficiency of $\hat{\theta}_{\text{obs}}$. However, for the linear
138 form (1), we can derive exact results for any sample sizes, which can render
139 statistical insights without any distraction of approximation.

140 Specifically, by the definition of self-efficiency as given in Section 6 of
141 our paper, $\hat{\theta}_{\text{obs}}$ is self-efficient with respect to a given MSE norm, which is

142 the same as Var when $\hat{\theta}_{\text{obs}}$ is unbiased, if and only if $\hat{\theta}_{\text{com}}$ is orthogonal to
143 $\hat{\theta}_{\text{obs}} - \hat{\theta}_{\text{com}}$, that is,

$$\text{Cov}(\hat{\theta}_{\text{com}}, \hat{\theta}_{\text{obs}} - \hat{\theta}_{\text{com}}) = 0. \quad (3)$$

144 But the linearity of (1) renders the linear decomposition

$$\hat{\theta}_{\text{com}} = r\hat{\theta}_{\text{obs}} + (1-r)\hat{\theta}_{\text{mis}}, \quad \text{where} \quad \hat{\theta}_{\text{mis}} = \frac{\sum_{i=1}^n (1-R_i)g(Y_i)}{\sum_{i=1}^n (1-R_i)} \quad (4)$$

145 and $r = \sum_{i=1}^n R_i/n$ is the proportion of the observed data size. Suppose
146 our MSE calculation is conditioning on the *missing-data pattern*, that is, the
147 values of $\{R_i\}_{i=1}^n$. Then under the conditional independence assumption of
148 Y_i and R_i given X_i , $\text{Cov}(\hat{\theta}_{\text{obs}}, \hat{\theta}_{\text{mis}}) = 0$. Consequently, (3) is equivalent to

$$r\text{Var}(\hat{\theta}_{\text{obs}}) = \text{Var}(\hat{\theta}_{\text{com}}) \iff \text{Var}(\hat{\theta}_{\text{obs}}) \propto \frac{1}{n_{\text{obs}}}. \quad (5)$$

149 That is, for the sample average (1) as an estimation procedure, it is (exactly)
150 self-efficient, with respect to the MDM as previously specified, if and only
151 if the variance of the procedure follows (exactly) the well-known inverse-
152 sample-size rule (for all samples sizes or a sample size sufficiently large).
153 But this is trivially true when Y_i 's are i.i.d.

154 We were therefore puzzled initially when we read Yang-Kim's statement
155 that (1) is self-efficient only in the first case of their first example. Since (5)
156 is a sufficient and necessary condition (assuming $\text{Cov}(\hat{\theta}_{\text{obs}}, \hat{\theta}_{\text{mis}}) = 0$), we
157 know that in order for this statement to hold, we must consider a different
158 variance operation for which (5) will hold only for the first case of Yang-
159 Kim's first example. Given Yang-Kim's regression-like setting, the obvious
160 alternative choice would be the conditional variance $\text{Var}(\hat{\theta}_{\text{obs}}|\vec{X})$, where $\vec{X} =$
161 (X_1, \dots, X_n) . Indeed, for this choice of replications (i.e., with \vec{X} fixed),
162 $\text{Var}(\hat{\theta}_{\text{obs}}|\vec{X})$ is free of \vec{X} for the first case in Yang-Kim's Example 1, where
163 $g(Y) = Y$ and only its conditional mean depends on X , not its conditional
164 variance. For other cases in Yang-Kim, the conditional variance of $g(Y)$
165 given X is not free of X either because $g(Y)$ is not linear in Y (e.g., $g(Y) =$
166 $I(Y < C)$ as in their Example 1) or $g(Y)$ is linear in Y , but $E(Y|X)$ itself

167 is not linear in X (e.g., the log-linear example in their Example 3, where
168 $E(Y|X) = \exp\{X^\top\beta + \sigma^2/2\}$).

169 However, even if we adopt this conditional evaluation when the estimand
170 is defined unconditionally, we still cannot conclude that the procedure in
171 Yang-Kim’s Example 2 is not self-efficient because this example is a special
172 case of the first case of their Example 1, by setting the regression intersection
173 to be zero. We therefore wonder if Yang-Kim used some other variance
174 operation for determining the procedure (1) is self-efficient in the first case
175 of Example 1, but not for a special case of it as in Example 2.

176 Our puzzle notwithstanding, Yang-Kim’s general message is the one that
177 we share, that is, one should not take self-efficiency for granted. Fortunately,
178 there are other ways to ensure the consistency of Rubin’s variance combining
179 rules, as Chen reported. Moreover, as we demonstrated in Section 8 of our
180 paper, it is possible for uncongeniality to effectively cancel self-inefficiency
181 to produce a consistent variance estimator by Rubin’s combining rule, high-
182 lighting the intricate nature of multi-phase inference.

183 **3. EM, MI, and FI – are they cousins?**

184 Yang-Kim also raised the issue of the links between MI to EM and to
185 Fractional Imputation (FI). As we stated in Section 4.2 of our paper, “per-
186 forming MI with an infinite number of imputations (and with the plug-in
187 predictive imputation) is the same as carrying out the final EM iteration.”
188 This is because the E-step of the EM algorithm evaluates the conditional
189 expectation of the complete-data score function $S(\theta; Z_{\text{com}})$ with respect to
190 $p(Z_{\text{mis}}|Z_{\text{obs}}, \theta = \theta^{(t)})$, where $Z_{\text{com}} = \{Z_{\text{obs}}, Z_{\text{mis}}\}$ with Z_{obs} and Z_{mis} denot-
191 ing respectively the observed data and missing data. That is, at the $(t+1)th$
192 iteration of EM, we utilize the so-called Q-function in the EM literature (see
193 van Dyk and Meng [4] for an overview):

$$Q(\theta|\theta^{(t)}) = E \left[S(\theta; Z_{\text{com}}) | Z_{\text{obs}}, \theta = \theta^{(t)} \right]. \quad (6)$$

194 Therefore, at the last iteration of EM, we compute $Q(\theta|\theta^*)$, where $\theta^* =$
195 $\lim_{t \rightarrow \infty} \theta^{(t)}$. This is equivalent to using an infinite number of draws from

196 $p(Z_{\text{mis}}|Z_{\text{obs}}, \theta = \theta^*)$, that is, an infinite number of imputations from the
197 “plug-in” predictive posterior to perform multiple imputation inference.

198 Multiple imputation, although is closely related to EM as a method to
199 deal with missing data problems, is designed to handle more general sit-
200 uations where subsequent analysis with complete data can use any valid
201 estimating method in addition to MLE. For example, if the subsequent
202 complete-data analysis uses an estimating equation

$$U(\theta, ; Z_{\text{com}}) = 0, \quad (7)$$

203 then, as shown in our paper, the point estimator from MI is asymptotically
204 equivalent to solving the following observed-data estimating equation:

$$E(U(\theta; Z_{\text{com}})|Z_{\text{obs}}) = 0, \quad (8)$$

205 where the conditional distribution $p(Z_{\text{mis}}|Z_{\text{obs}})$ is the predictive distribution
206 of Z_{mis} from a Bayesian model, which is also asymptotically equivalent to
207 its frequentist’s counterpart

$$E(U(\theta; Z_{\text{com}})|Z_{\text{obs}}; \theta = \theta^*) = 0, \quad (9)$$

208 where θ^* is the observed-data MLE.

209 The fractional imputation, as described in Yang and Kim [3], seems to
210 accomplish the same task as MI but via importance sampling. Specifically,
211 it seeks to approximate (9) by

$$E(U(\theta; Z_{\text{com}})|Z_{\text{obs}}; \theta = \theta^*) \approx \sum_j w_j \cdot U(\theta; Z_{\text{com}}^{(j)}|Z_{\text{obs}}; \theta = \theta^*) = 0, \quad (10)$$

212 where $w_j \propto p(Z_{\text{com}}^{(j)}|Z_{\text{obs}}; \theta = \theta^*)/h(Z_{\text{com}}^{(j)}|Z_{\text{obs}})$ is the (standardized) weight
213 of the importance sampling with $h(Z_{\text{com}}^{(j)}|Z_{\text{obs}})$ as its (pre-chosen) proposal
214 distribution. The accuracy of this weighting approach, as is well-understood,
215 depends on the choice of the proposal.

216 If our understanding of FI is correct, then there is a link between FI
217 to another cousin in the big family of missing-data approaches, that is,
218 Stochastic EM (SEM; see Celeux et al. [5]; not to be confused with the SEM

219 algorithm of Meng and Rubin [6] for computing variance estimators). SEM
220 uses Monte Carlo draws from $p(Z_{\text{com}}^{(j)}|Z_{\text{obs}}; \theta = \theta^{(t)})$ to form a Monte Carlo
221 estimator of (6), and then it iterates just as the standard EM. Because the
222 resulting iterative sequence now depends on noise introduced by the Monte
223 Carlo draws, it is stochastic. Clearly we can introduce importance sampling
224 in approximating (6) as well, where the proposal density can vary with
225 iteration—preferred for statistical efficiency, or fixed at some $h(Z_{\text{com}}^{(j)}|Z_{\text{obs}})$ —
226 preferred for computational efficiency, or a hybrid of them to achieve a
227 sensible compromise. In that sense, SEM to FI is like EM to MI, as FI and
228 MI can be viewed as the final iteration of SEM and EM, respectively.

229 Another closely related cousin is the ES algorithm investigated by Elashoff
230 and Ryan [7], which replaces (9) by

$$E(U(\theta; Z_{\text{com}})|Z_{\text{obs}}; \theta = \theta^{(t)}) = 0, \quad (11)$$

231 and then *solves* (hence the “S” in “ES”) it to obtain $\hat{\theta}^{(t+1)}$. This generalizes
232 the EM algorithm for maximizing likelihood estimation to solving a more
233 general estimating equation with incomplete data. A special case of ES is
234 the iterative Projection-Solution algorithm for quasi-likelihood in Heyde and
235 Morton [8], as Desmond cited; we also fully agree with Desmond that pro-
236 jection of the estimating equation, as in (11), is more powerful and fruitful
237 than projection of estimators, at least for finite-sample properties. MI then
238 can be viewed as the final iteration of ES, but with the Expectation step
239 carried out via Monte Carlo.

240 4. A clean theory of the messy world of pre-processing?

241 A common theme of the multi-phase examples provided by the VIPs
242 is that they are all *messy*. Some are necessarily so, such as protecting
243 confidentiality, as outlined by Reiter, because it would forever be a struggle
244 between protecting privacy and preserving information. We simply cannot
245 have both: complete protection and full information. Others are avoidable,
246 such as those unsettling zeros produced by the team that did not share

247 the same VP as Draper. But the messiest of all are those cases where the
248 analysts have little idea about what was done to their data, which is rather
249 the rule than the exception, as in many cases of pre-processing. Could then
250 there be any “clean” theory to deal with such messiness?

251 Draper outlined the idea of a Bayesian composition model, borrowing the
252 notion of function composition, $f_2(f_1(D))$, where D denotes data, asking
253 how f_i 's should be constructed to preserve as much information as possible.
254 A similar question was asked in Blocker and Meng [9], in the context of dis-
255 tributed pre-processing, that is, what the analyst received is in the form of
256 $\{g_j(D_j)\}_{j=1}^J$ from a system with J pre-processors (e.g., one for each experi-
257 ment). The question then is what are the computationally economical and
258 yet information-preserving choices of $g_j, j = 1, \dots, J$? We can see clearly
259 the competing nature of our goals: computationally, the most economical
260 choice would be (say) to set all $g_j \equiv 0$, which is ridiculous as it preserves
261 no information. On the other hand, choosing $g_j(D_j) = D_j, j = 1, \dots, J$ will
262 preserve whatever information is contained in the data, but it achieves no
263 computational saving or any other kind of desirable pre-processing (such as
264 privacy protection). Furthermore, preserving information is not a meaning-
265 ful requirement without specifying the meaning of information or for what
266 purposes (e.g., estimation? testing? prediction?).

267 But even in the classic context of sufficiency with respect to a well-
268 specified parametric family, it is not easy at all to obtain a “clean” theory
269 for the most economical lossless data compression. Blocker and Meng [9] ob-
270 tained sufficient conditions, as well as necessary conditions, but not sufficient
271 *and* necessary conditions for such g'_j 's, $j = 1, \dots, I$. A simple example suffices
272 to illustrate the difficulty. Suppose $I = 1$ and the data $D_1 = \{Y_j\}_{j=1}^J$ are
273 i.i.d. Poisson observations with mean θ . The preprocessor however chooses
274 the convenient (and very wrong) model $N(\mu, 1)$, and hence he preserves its
275 sufficient statistic $\bar{Y}_n = \sum_{i=1}^n Y_i/n$. However, since \bar{Y}_n is also the sufficient
276 statistic for θ under the Poisson model, there is no information lost even if
277 the pre-processor used an entirely wrong model, which does not even share
278 the same support with the correct model. This indicates the difficulties with

279 establishing if-and-only-if conditions for pre-processing, since we can obtain
280 the same results with very different models.

281 The problem becomes even harder when sufficient statistics are difficult
282 to come by, as Banks-Peña questioned, and when information in the data
283 is hard to quantify; and most challengingly, when the pre-processor is not
284 well-informed of, or just unable to model, the purposes of analyses by down-
285 stream users. But we hope these challenges will help to entice those with
286 strong adventurous spirits to join us in our search for a “clean theory” about
287 pre-processing. By clean theories we mean those that can either shed lights
288 on the treacherous paths, or those that can lead to practical and effective
289 (though not necessarily optimal) procedures, such as our variance doubling
290 rule.

291 **5. Is bias-variance trade-off also critical for multi-phase inference?**

292 Yes, very much so. Yang-Kim’s question on robustness of modeling, by
293 analysts and by imputer, lies at the heart of statistical inference, and to
294 answer it sensibly one must have full grasp of one of a very few fundamental
295 principles of statistics, namely, the ubiquity of robustness-efficiency trade-
296 off, a.k.a, bias-variance trade-off. “Some questions” raised by Banks-Peña,
297 especially the last one, emphasized the very trade-off. Chen’s emphasis on
298 paying attention to (analysis) model selection touches on the same issue, be-
299 cause the most critical balancing act of any model selection procedure is to
300 ensure capturing replicable signals but not to overfit the idiosyncratic indi-
301 vidualities. Not incidently, this need for balancing presents a grand challenge
302 for building a framework toward accumulating statistical evidence underly-
303 ing *individualized inference/prediction*, but that is the subject for another
304 hard-to-write paper. An initial attempt was made in Meng [10] for estab-
305 lishing a multi-resolution framework supporting individualized inference, as
306 one of the framework trio. (The other two cover multi-phase inference, for
307 which our current paper is a sequel, and multi-source inference, as in Meng
308 [11].)

309 Both Banks-Peña and Desmond raised the possibility of an all-encompassing
310 Bayesian modelling strategy for multi-phase inference. Indeed any (serious)
311 Bayesian can, and probably is compelled to, model the entire multi-phase
312 as a whole, which has the added benefit of being coherent. But then there is
313 a bias-variance trade-off. Given the *uncongenial* nature of the multi-phase
314 paradigm, literally and technically, such modeling would necessarily need
315 critical assumptions that are known to be false or minimally cannot be con-
316 firmed by reality, because otherwise there would not be any uncongeniality
317 in the first place. And even seemingly “good” pre-processing models can
318 (and often) lead to provably undesirable results, as Banks-Peña’s Carlo-Bob
319 example further illustrated. This is what makes the multi-phase inference
320 paradigm interesting, intriguing, and inspiring. The many examples from
321 government statistics agencies, especially under the mandate of disclosure
322 protection, as succinctly summarized by Reiter, and from industrial and
323 business sectors, as vividly illustrated by Desmond, highlighted the urgent
324 need of developing this paradigm.

325 Indeed, as Desmond correctly recognized, ultimately the multi-phase
326 paradigm needs to handle an unholy trinity: missingness, misspecification,
327 and uncongnialty. In comparison to this grand goal, what we presented in
328 the current paper is only one of many needed building blocks. We are there-
329 fore humbled by the kind encouragements from the VIPs, especially the
330 extremely flattering endorsement from Banks-Peña, Chen, Desmond and
331 Reiter. We also particularly thank Draper for his *RSS* style vote for thanks,
332 regardless of whether he would propose or second, especially because initially
333 we did plan to seek such a vote. Ultimately, our long journey of dealing with
334 uncongenality led us to the welcoming arms of *Statistica Sinica*, to which
335 we are deeply grateful.

336 References

- 337 [1] K. Liu, X.-L. Meng, There is individualized treatment. Why not indi-
338 vidualized inference?, *Annual Review of Statistics and Its Applications*
339 3 (2016) 79–111.

- 340 [2] X.-L. Meng, X. Xie, I got more data, my model is more refined, but
341 my estimator is getting worse! Am I just dumb?, *Econometric Reviews*
342 33 (2014) 218–250.
- 343 [3] S. Yang, J. K. Kim, Fractional imputation in survey sampling: A
344 comparative review, *Statistical Science* 31 (2016) 415–432.
- 345 [4] D. A. van Dyk, X.-L. Meng, Cross-fertilizing strategies for better EM
346 mountain climbing and DA field exploration: A graphical guide book,
347 *Statistical Science* 25 (2010) 429–449.
- 348 [5] G. Celeux, D. Chauveau, J. Diebolt, Stochastic versions of the EM
349 algorithm: an experimental study in the mixture case, *Journal of Sta-*
350 *tistical Computation and Simulation* 55 (1996) 287–314.
- 351 [6] X.-L. Meng, D. B. Rubin, Using EM to obtain asymptotic variance-
352 covariance matrices: The SEM algorithm, *Journal of the American*
353 *Statistical Association* 86 (1991) 899–909.
- 354 [7] M. Elashoff, L. Ryan, An EM algorithm for estimating equations, *Jour-*
355 *nal of Computational and Graphical Statistics* 13 (2004) 48–65.
- 356 [8] C. Heyde, R. Morton, Quasi-likelihood and generalizing the EM algo-
357 rithm, *Journal of the Royal Statistical Society. Series B (Methodologi-*
358 *cal)* (1996) 317–327.
- 359 [9] A. W. Blocker, X.-L. Meng, The potential and perils of preprocessing:
360 Building new foundations, *Bernoulli* 19 (2013) 1176–1211.
- 361 [10] X.-L. Meng, A trio of inference problems that could win you a Nobel
362 prize in statistics (if you help fund it), in: *Past, Present, and Future of*
363 *Statistical Science* (Eds: Lin et. al.), CRC Press, 2014, pp. 537–562.
- 364 [11] X.-L. Meng, Statistical paradises and paradoxes in big data (I): The
365 bigger the data, the surer we miss our target?, Technical Report, De-
366 partment of Statistics, Harvard University (2017).