# Imputation in Clinical Research

**Hansheng Wang**
*Peking University, Beijing, People's Republic of China*

## INTRODUCTION

Missing values or incomplete data are commonly encountered in clinical research and are studied by many authors.[1–3] Basically, the causes of missing values in a study can be classified into two categories. The first category includes the reasons that are not directly related to the study. For example, a patient may be lost to follow-up because he/she moves out of the area. This category of missing values can be considered as *missing completely at random*. The second category includes the reasons that are related to the study. For example, a patient may withdraw from the study due to treatment-emergent adverse events. In practice, it is not uncommon to have multiple assessments from each subject. Subjects with all observations missing are called unit nonrespondents. Because unit nonrespondents do not provide any useful information, these subjects are usually excluded from the analysis. On the other hand, the subjects with some, but not all, observations missing are referred to as item nonrespondents. In practice, excluding item nonrespondents from the analysis is considered against the *intent-to-treat* (ITT) principle and, hence, not acceptable. In clinical research, the primary analysis is usually conducted based on ITT population, which includes all randomized subjects with at least posttreatment evaluation. As a result, most item nonrespondents may be included in the ITT population. In practice, excluding item nonrespondents may seriously decrease power/efficiency of the study.

To account for item nonrespondents, two methods are commonly considered. The first method is the so-called likelihood-based method. Under a parametric model, the marginal likelihood function for the observed responses is obtained by integrating out the missing responses. The parameter of interest can then be estimated by the maximum likelihood estimator (MLE). Consequently, a corresponding test (e.g., likelihood ratio test) can be constructed. The merit of this method is that the resulting statistical procedures are usually efficient. The drawback is that the calculation of the marginal likelihood could be difficult. As a result, some special statistical or numerical algorithms are commonly applied for obtaining the MLE. For example, the expectation–maximization (EM) algorithm is one of the most popular methods for obtaining the MLE when there are missing data. The other method for item nonrespondents is imputation. Compared with the likelihood-based method, the method of imputation is relatively simple and easy to apply. The idea of imputation is to treat the imputed values as the observed values and then apply the standard statistical software for obtaining consistent estimators. However, it should be noted that the variability of the estimator obtained by imputation is usually different from the estimator obtained from the complete data. In this case, the formulas designed to estimate the variance of the complete data set cannot be used to estimate the variance of estimator produced by the imputed data. As an alternative, two methods are considered for estimation of its variability. One is based on Taylor's expansion. This method is referred to as the "linearization method." The merit of the linearization method is that it requires less computation. However, the drawback is that its formula could be very complicated and/or nontrackable. The other approach is based on resampling method (e.g., bootstrap and jackknife). The drawback of the resampling method is that it requires an intensive computation. The merit is that it is very easy to apply. With the help of a fast-speed computer, the resampling method has become much more attractive in practice.

Note that imputation is not only popular in clinical research, it is also very popular in many other statistical fields such as sample survey. However, the imputation methods in clinical research are more diversified due to the complexity of the study design relative to sample survey. As a result, the statistical properties of many commonly used imputation methods in clinical research are still unknown, while most imputation methods used in sample survey are well studied. Hence, the imputation methods in clinical research provide a unique challenge and also an opportunity for the statisticians in the area of clinical research. In what follows, we will summarize the most commonly used imputation methods and investigate their statistical properties. Recent development will also be discussed.

## LAST OBSERVATION CARRY FORWARD/ ENDPOINT ANALYSIS

Last observation carry forward (LOCF) and endpoint analysis (EPA) are probably the most commonly used

imputation methods in clinical research. For illustration purposes, one example is described below.

Consider a randomized, parallel-group clinical trial comparing $r$ treatments. Each patient is randomly assigned to one of the treatments. According to the protocol, each patient should undergo $s$ consecutive visits. Let $y_{ijk}$ be the observation from the $k$th subject in the $i$th treatment group at visit $j$. The following statistical model is usually considered.

$$y_{ijk} = \mu_{ij} + \epsilon_{ijk}, \text{ where } \epsilon_{ijk} \sim N(0, \sigma^2) \tag{1}$$

where $\mu_{ij}$ represents the fixed effect of the $i^{\text{th}}$ treatment at visit $j$. If there are no missing values, the primary comparison between treatments will be based on the observations from the last visit ($j = s$) because this reflects the treatment difference at the end of the treatment period. However, it is not necessary that every subject completes the study. Suppose that the last evaluable visit is $j^* < m$ for the $k$th subject in the $i$th treatment group. Then the value of $y_{ij^*k}$ can be used to impute $t_{isk}$. After imputation, the data at endpoint are analyzed by the usual analysis of variance (ANOVA) model. We will refer to the procedure described above as LOCF.

Note that the method of LOCF is usually applied according to the ITT principle. The ITT population include all randomized subjects. In clinical research, although the LOCF is commonly employed, the statistical motivation and the consequence are not clear. In what follows, we attempt to provide two different approaches to understand the statistical properties of LOCF.

## Bias-Variance Trade-Off

The objective of a clinical study is usually to assess the safety and efficacy of a test treatment under investigation. Statistical inferences on the efficacy parameters are usually obtained. In practice, a sufficiently large number of sample size is required to obtain a reliable estimate and to achieve a desired power for establishment of the efficacy of the treatment. The reliability of an estimator can be evaluated by bias and by variability. A reliable estimator should have a small or zero bias with small variability. Hence, the estimator based on LOCF and the estimator based on completers are compared in terms of their bias and variability.

For illustration purposes, we focus on only one treatment group with two visits. Assume that there are a total of $n = n_1 + n_2$ randomized subjects, where $n_1$ subjects complete the trial, while the remaining $n_2$ subjects only have observations at visit 1. Let $y_{ik}$ be the response from the $k$th subject at the $i$th visit and $\mu_i = E(y_{ik})$. The parameter of interest is $\mu_2$. The estimator based on completers is given by

$$\bar{y}_c = \frac{1}{n_1} \sum_{k=1}^{n_1} y_{i2k}$$

On the other hand, the estimator based on LOCF can be obtained as

$$\bar{y}_{\text{LOCF}} = \frac{1}{n} \left( \sum_{i=1}^{n_1} y_{i2k} + \sum_{i=n_1+1}^{n} y_{i1k} \right)$$

It can be verified that the bias of $\bar{y}_c$ is 0 with variance $\sigma^2 = n_1$, while of the bias of $\bar{y}_{\text{LOCF}}$ is $n_2(\mu_1 - \mu_2)/n$ with variance $\sigma^2/(n_1 + n_2)$.

As noted, although LOCF may introduce some bias, it decreases the variability. In a clinical trial with multiple visits, usually, $\mu_j \approx \mu_s$ if $j \approx s$. This implies that the LOCF is recommended if the patients withdraw from the study at the end of the study. However, if a patient drops out of the study at the very beginning, the bias of the LOCF could be substantial. As a result, it is recommended that the results from the analysis based on LOCF be interpreted with caution.

## Hypothesis Testing

In practice, the LOCF is viewed as a pure imputation method for testing the null hypothesis of

$$H_0 : \mu_{1s} = \cdots = \mu_{rs}$$

where $\mu_{ij}$ are defined in Eq. 1.

Zhong and Shao[4] provided another look of statistical properties of the LOCF under the above null hypothesis. More specifically, they partitioned the total patient population into $s$ subpopulations according to the time when the number of patients drop out from the study. Note that in their definition, the patients who complete the study are considered a special case of ''drop out'' at the end of the study. Then $\mu_{ij}$ represents the population mean of the $j$th subpopulation under treatment $i$. Assume that the $j$th subpopulation under the $i$th treatment accounts for $p_i \times 100\%$ of the overall population under the $i$th treatment. They argue that the objective of the intend-to-treat analysis is to test the following hypothesis test

$$H_0 : \mu_1 = \cdots = \mu_r \tag{2}$$

where $\mu_i = \sum_{j=1}^{s} p_{ij} \mu_{ij}$. Based on the above hypothesis, Zhong and Shao[4] indicated that the LOCF bears the following properties:

1.  In the special case of $r = 2$, the asymptotic ($n_i \to \infty$) size of the LOCF under $H_0$ is $\leq \alpha$ if and only if

$$\lim \left( \frac{n_2 \tau_1^2}{n} + \frac{n_1 \tau_2^2}{n} \right) \leq \lim \left( \frac{n_1 \tau_1^2}{n} + \frac{n_2 \tau_2^2}{n} \right)$$

where

$$\tau_i^2 = \sum_{j=1}^{s} p_{ij}(\mu_{ij} - \mu_i)^2$$

The LOCF is robust in the sense that its asymptotic size is $\alpha$ if $\lim (n_1/n) = (n_2/n)$ or $\tau_1^2 = \tau_2^2$. Note that, in reality, $\tau_1^2 = \tau_2^2$ is impractical unless $\mu_{ij} = \mu_i$ for all $j$. However, $n_1 = n_2$ (as a result $\lim (n_1/n) = \lim (n_2/n)$) is very typical, in practice. The above observation indicates in such a situation ($n_1 = n_2$) that LOCF is still valid.

2. When $r = 2$, $\tau_1^2 \neq \tau_2^2$, and $n_1 \neq n_2$, the LOCF has an asymptotic size smaller than $\alpha$ if

$$(n_2 - n_1)\tau_1^2 < (n_2 - n_1)\tau_2^2 \tag{3}$$

or larger than $\alpha$ if " $<$ " in Eq. 3 is replaced by " $>$ ."

3. When $r \geq 3$, the asymptotic size of the LOCF is generally not $\alpha$ except for some special case (e.g., $\tau_1^2 = \tau_2^2 = \cdots = \tau_r^2 = 0$).

Because the LOCF usually does not produce a test with asymptotic significance level $\alpha$ when $r \geq 3$, Zhong and Shao[4] proposed the following testing procedure based on the idea of poststratification. The null hypothesis $H_0$ should be rejected if $T > \chi_{1-\alpha, r-1}^2$, where $\chi_{1-\alpha, r-1}^2$ is the $(1-\alpha)$th quantile of a chi-square random variable with $r-1$ degrees of freedom and

$$T = \sum_{i=1}^{r} \frac{1}{\hat{V}_i} \left( \bar{y}_{i\cdot\cdot} - \frac{\sum_{i=1}^{r} \bar{y}_{i\cdot\cdot}/\hat{V}_i}{\sum_{i=1}^{r} 1/\hat{V}_i} \right)^2$$

$$\hat{V}_i = \frac{1}{n_i(n_i - 1)} \sum_{j=1}^{s} \sum_{k=1}^{n_{ij}} (y_{ijk} - \bar{y}_{i\cdot\cdot})^2$$

Under the Eq. model 1 and the null hypothesis of 2, this procedure has the exact type I error $\alpha$.

## MEAN/MEDIAN IMPUTATION

Missing ordinal responses are also commonly encountered in clinical research. For those types of missing data, mean or median imputation is commonly considered. Let $x_i$ be the ordinal response from the $i$th subject, where $i = 1, \cdots, n$. The parameter of interest is $\mu = E(x_i)$. Assume that $x_i$ for $i = 1, \cdots, n_1 < n$ are observed and the rest are missing.

Median imputation will impute the missing response by the median of the observed response (i.e., $x_i$, $i = 1, \cdots, n_1$). The merit of median imputation is that it can keep the imputed response within the sample space as the original response by appropriately defining the median. The sample mean of the imputed data set will be used as an estimator for the population mean. However, as the parameter of interest is population mean, the median imputation may lead to biased estimates.

As an alternative, mean imputation will impute the missing value by the sample mean of the observed units (i.e., $1/n_1 \sum_i^{n_i} x_i$. The disadvantage of the mean imputation is that the imputed value may be out of the original response sample space. However, it can be shown that the sample mean of the imputed data set is a consistent estimator of population mean. Its variability can be assessed by jackknife method proposed by Rao and Shao.[5]

In practice, usually, each subject will provide more than one ordinal response. The summation of those ordinal responses (total score) is usually considered as the primary efficacy parameter. The parameter of interest is the population mean of the total score. In such a situation, mean/median imputation can be carried out for each ordinal response within each treatment group.

## REGRESSION IMPUTATION

The method of regression imputation is usually considered when covariates are available. Regression imputation assumes a linear model between the response and the covariates. The method of regression imputation has been studied by various authors.[6,7]

Let $y_{ijk}$ be the response from the $k$th subject in the $i$th treatment group at the $j$th visit. The following regression model is considered:

$$y_{ijk} = \mu_i + \beta_i x_{ij} + \epsilon_{ijk} \tag{4}$$

where $x_{ij}$ is the covariate of the $k$th subject in the $i$th treatment group. In practice, the covariates $x_{ij}$ could be demographic variables (e.g., age, sex, and race) or patient's baseline characteristics (e.g., medical history or disease severity). Model 4 suggests a regression imputation method. Let $\hat{\mu}_i$ and $\hat{\beta}_i$ denote the estimators of $\mu_i$ and $\beta_i$ based on complete data set, respectively. If $y_{ijk}$ is missing, its predicted mean value $y_{ijk}^* = \hat{\mu}_i + \hat{\beta}_i x_{ij}$ is used to impute. The imputed values are treated as true responses and the usual ANOVA is used to perform the analysis.

## MARGINAL/CONDITIONAL IMPUTATION FOR CONTINGENCY TABLES

In an observational study, two-way contingency tables can be used to summarize two-dimensional categorical data. Each cell (category) in a two-way contingency table is defined by a two-dimensional categorical variable

(A, B), where A and B take values in $\{1, \cdots, a\}$ and $\{1, \cdots, b\}$, respectively. Sample cell frequencies can be computed based on the observed responses of (A, B) from a sample of units (subjects). Statistical interest includes the estimation of cell probabilities and testing hypotheses of goodness of fit or the independence of the two components A and B. In an observational study, there can be more than one strata. It is assumed that within a strata, sampled units independently have the same probability $\pi_A$ to have missing B and observed A, $\pi_B$ to have missing A and observed B, $\pi_C$ to have observed A and B. (The probabilities $\pi_A$, $\pi_B$; and $\pi_C$ may be different in different imputation classes.) As units with both A and B missing are considered as unit nonrespondent, they are excluded in the analysis. As a result, without loss of generality, it is assumed that $\pi_A + \pi_B + \pi_C = 1$.

For a two-way contingency table, it is very important for an appropriate imputation method to keep imputed values in the appropriate sample space. Whether in calculating the cell probability or in testing hypotheses (e.g., testing independence or goodness of fit), the corresponding statistical procedures are all based on the frequency counts of a contingency table. If the imputed value is out of the sample space, additional categories will be produced which is of no practical meaning. As a result, two hot deck imputation methods are thoroughly studied by Wang[8] and Wang and Shao.[9]

## Simple Random Sampling

Consider a sampled unit with observed A = i and missing B. Two imputation methods were studied by Wang and Shao.[9] The marginal (or unconditional) random hot deck imputation method imputes B by the value of B of a unit randomly selected from all units with observed B. The conditional hot deck imputation method imputes B by the value of B of a unit randomly selected from all units with observed B and A = i. All nonrespondents are imputed independently.

### Point estimation

After imputation, the cell probabilities $p_{ij}$ can be estimated using the standard formulas in the analysis of data from a two-way contingency table by treating imputed values as observed data. Denote these estimators by $\hat{p}^{\mathrm{I}}_{ij}$, where $i = 1, \cdots, a$, and $j = 1, \cdots, b$. Let

$$\hat{p}^{\mathrm{I}} = (\hat{p}^{\mathrm{I}}_{11}, \ldots, \hat{p}^{\mathrm{I}}_{1b}, \ldots, \hat{p}^{\mathrm{I}}_{a1}, \ldots, \hat{p}^{\mathrm{I}}_{ab})'$$

and

$$p = (p_{11}, \ldots, p_{1b}, \ldots, p_{a1}, \ldots, p_{ab})'$$

where $p_{ij} = P(A = i, B = j)$. Intuitively, marginal random hot deck imputation leads to consistent estimators

of $p_{i\cdot} = P(A = i)$ and $p_{\cdot j} = P(B = j)$, but not $p_{ij}$. Wang and Shao[9] showed that $\hat{p}^{\mathrm{I}}$ under conditional hot deck imputation are consistent, asymptotically unbiased, and asymptotically normal.

**Theorem 1.** Assume that $\pi_C > 0$. Under conditional hot deck imputation,

$$\sqrt{n}(\hat{p}^{\mathrm{I}} - p) \to_d N(0, MPM' + (1 - \pi_C)P)$$

where $P = \mathrm{diag}\{p\} - pp'$ ($\mathrm{diag}\{a\}$ is a diagonal matrix whose ith diagonal element is the Ith component of the vector $a$),

$$\begin{aligned} M &= \frac{1}{\sqrt{\pi_C}} (I_{a\times b} - \pi_A \mathrm{diag}\{p_{B|A}\} I_a \otimes U_b \\ &\quad - \pi_B \mathrm{diag}\{p_{A|B}\} U_a \otimes I_b), \end{aligned}$$

$$p_{A|B} = (p_{11}/p_{\cdot 1}, \cdots, p_{1b}/p_{\cdot b}, \cdots, p_{a1}/p_{\cdot 1}, \cdots, p_{ab}/p_{\cdot b})',$$

$$p_{B|A} = (p_{11}/p_{1\cdot}, \cdots, p_{1b}/p_{1\cdot}, \cdots, p_{a1}|p_{a\cdot}, \cdots, p_{ab}/p_{a\cdot})'$$

where $I_a$ denotes an $a$-dimensional identity matrix, $U_b$ denotes a $b$-dimensional square matrix with all components being 1, and $\otimes$ is the Kronecker product.

### Goodness-of-fit test

A direct application of Theorem 1 is to obtain a Wald-type test for goodness of fit. Consider the null hypothesis of the form $H_0$: $p = p_0$, where $p_0$ is a known vector. Under $H_0$,

$$X_W^2 = n(\hat{p}^* - p_0^*)' \hat{\Sigma}^{*-1} (\hat{p}^* - p_0^*) \to_d \chi_{ab-1}^2$$

where $\chi_v^2$ denotes a random variable having the chi-square distribution with $v$ degrees of freedom, $\hat{p}^*$ ($p_0^*$) is obtained by dropping the last component of $\hat{p}^{\mathrm{I}}$ ($p_0$), and $\hat{\Sigma}^*$ is the estimated asymptotic covariance matrix of $\hat{p}^*$, which can be obtained by dropping the last row and column of $\hat{\Sigma}$, the estimated asymptotic covariance matrix of $\hat{p}^{\mathrm{I}}$.

Noting the fact that the computation of $\hat{\Sigma}^{*-1}$ is complicated, Wang and Shao proposed a simple correction of the standard Pearson chi-square statistic by matching the first-order moment, an approach developed by Rao and Scott.[9] Let

$$X_G^2 = n \sum \frac{(\hat{p}^{\mathrm{I}}_{ij} - p_{ij})^2}{p_{ij}}$$

It is noted that under conditional imputation the asymptotic expectation of $X_G^2$

$$\frac{1}{\pi_C}(ab + \pi_A^2 a + \pi_B^2 b - 2\pi_A a - 2\pi_B b + 2\pi_A \pi_B + 2\pi_A \pi_B \delta)$$

$$- \pi_C ab + (ab - 1)$$

Let (see equation at bottom of the page below). Then the asymptotic expectation of $X_G^2/\lambda$ is $ab - 1$, which is the first-order moment of a standard chi-square variable with $ab - 1$ degrees of freedom. Thus, $X_G^2/\lambda$ can be used just like a normal chi-square statistic to test the goodness of fit. However, it should be noted that this is just an approximated test procedure which is not asymptotically correct. According to a Wang and Shao's simulation study, this test performs reasonably well with moderate sample sizes.

## TESTING FOR INDEPENDENCE

Testing for the independence between $A$ and $B$ can be performed by the following chi-square statistic when there is no missing data

$$X^2 \;=\; n \sum_{i,j} \frac{(\hat{p}_{ij}^{\mathrm{I}} - \hat{p}_{i\cdot}^{\mathrm{I}} \hat{p}_{\cdot j}^{\mathrm{I}})^2}{\hat{p}_{i\cdot}^{\mathrm{I}} \hat{p}_{\cdot j}^{\mathrm{I}}} \;\to_d\; \chi^2_{(a-1)(b-1)}$$

It is of interest to know what the asymptotic behavior of the above chi-square statistic is under both marginal and conditional imputation. It is found that under the null hypothesis of $A$ and $B$ are independent and conditional hot deck imputation

$$X^2 \to_d (\pi_C^{-1} + 1 - \pi_C)\chi^2_{(a-1)(b-1)}$$

and under marginal hot deck imputation

$$X^2_{\mathrm{MI}} \to_d \chi^2_{(a-1)(b-1)}$$

## Results Under Stratified Simple Random Sampling

### When number of strata is small

Stratified samplings are also commonly used in medical study. For example, a large epidemiology study is usually conducted by several large centers. Those centers are usually considered as strata. For those types of study, the number of strata is not very large; however, the sample size within each strata is very large. As a result, imputation is usually carried out within each strata.

Within the $h$th stratum, we assume that a simple random sample of size $n_h$ is obtained and samples across strata are obtained independently. The total sample size is $n = \sum_{h=1}^{H} n_h$, where $H$ is the number of strata and $n_h$ is the sample size in stratum $h$. The parameter of interest is the overall cell probability vector $p = \sum_{h=1}^{H} w_h p_h$, where $w_h$ is the $h$th stratum weight. The estimator of $p$ based on conditional imputation is given by $\hat{p}^{\mathrm{I}} = \sum_{h=1}^{H} w_h \hat{p}_h^{\mathrm{I}}$. Assume that $n_n = n \to \rho$ as $n \to \infty 1$, $h = 1, \cdots, H$. Then a direct application of Theorem 1 leads to

$$\sqrt{n}(\hat{p}^{\mathrm{I}} - p) \to_d N(0, \Sigma)$$

where

$$\Sigma \;=\; \sum_{h=1}^{H} \frac{w_h^2}{\rho_h} \Sigma_h$$

and $\Sigma_h$ is the $\Sigma$ in Theorem 1 but restricted to the $h$th stratum.

### When number of strata is large

In a medical survey, it is also possible to have the number of strata ($H$) very large, while the sample size within each strata is small. A typical example is that if a medical survey is conducted by family, then the family can be considered as a strata and all the members within the family become the samples from this strata. In such a situation, the method of imputation within stratum is impractical because it is possible that within a stratum, there are no completers. As an alternative, Wang and Shao[9] proposed the method of imputation across strata under the assumption that $(\pi_{h,A}, \pi_{h,B}, \pi_{h,C})$, where $h = 1, \cdots, H$, is constant. More specifically, let $n_{h,ij}^{\mathrm{C}}$ denote the number of completers in the $h$th stratum such that $A = i$ and $B = j$. For a sampled unit in the $k$th imputation class with observed $B = j$ and missing $A$, the missing value is imputed by $i$ according to the conditional probability

$$p_{ij}|B, k \;=\; \frac{\sum\limits_{h} w_h n_{h,ij}^{\mathrm{C}}/n_h}{\sum\limits_{h} w_h n_{h,\cdot j}^{\mathrm{C}}/n_h}$$

Similarly, the missing value of a sampled unit in the $k$th imputation class with observed $A = i$ and missing $B$ can be imputed by $j$ according to the conditional probability

$$p_{ij}|A, k \;=\; \frac{\sum\limits_{h} w_h n_{h,ij}^{\mathrm{C}}/n_h}{\sum\limits_{h} w_h n_{h,i\cdot}^{\mathrm{C}}/n_h}$$

Again, $\hat{p}^{\mathrm{I}}$ can be computed by ignoring imputation classes and treating imputed values as observed data. The following result establishes the asymptotic normality of $\hat{p}^{\mathrm{I}}$

$$\lambda \;=\; \frac{\frac{1}{\pi_C}(ab + \pi_A^2 a + \pi_B^2 b - 2\pi_A a - 2\pi_B b + 2\pi_A \pi_B + 2\pi_A \pi_B \delta) - \pi_C ab + (ab - 1)}{ab - 1}$$

based on the method of conditional hot deck imputation across strata.

**Theorem 2.** Let $(\pi_{h,A}, \pi_{h,B}, \pi_{h,C}) = (\pi_A, \pi_B, \pi_C)$ for all $h$. Assume further that $H \to \infty$ and that there are constants $c_j$, for $j = 1, \cdots, 4$, such that $n_h \leq c_1$, $c_2 \leq H w_h \leq c_3$, and $p_{h,ij} \geq c_4$ for all $h$. Then

$$\sqrt{n}(\hat{p}^{\mathrm{I}} - p) \to_d N(0, \Sigma)$$

where $\Sigma$ is the limit of

$$n\left( \sum_h \frac{w_h^2}{n_h} \Sigma_h + \Sigma_A + \Sigma_B \right)$$

Details regarding the expression of $\Sigma_h$, $\Sigma_A$, and $\Sigma_B$ and their procedure can be found in Wang.[8]

## CONCLUSION

Missing values or incomplete data are commonly encountered in clinical research. How to handle the incomplete data is always a challenge to the statisticians in practice. Imputation as one of very popular methodology to compensate for the missing data is widely used in biopharmaceutical research. As compared to its popular-
ity, however, its theoretical properties are far away from well understood. Further research is needed.

## REFERENCES

1. Little, R.J.; Rubin, D.B. *Statistical Analysis with Missing Data*; Wiley: New York, 1987.
2. Schafer, J.L. *Analysis of Incomplete Multivariate Data*; Chapman and Hall: London, 1997.
3. Kalton, G.; Kasprzyk, D. The treatment of missing data. Surv. Methodol. **1986**, *12*, 1–16.
4. Zhong, B.; Shao, J. *Last Observation Carry-Forward and Intention-to-Treat Analysis*; 2002, Submitted.
5. Rao, J.N.K.; Shao, J. Jackknife variance estimation with survey data under hot deck imputation. Biometrika **1992**, *79*, 811–822.
6. Srivastava, M.S.; Carter, E.M. The maximum likelihood method for non-response in sample surveys. Surv. Methodol. **1986**, *12*, 61–72.
7. Shao, J.; Wang, H. Sample correlation coefficients based on survey data under regression imputation. JASA **2002**, *97*, 544–552.
8. Wang, H. Two-way contingency tables with marginally and conditionally imputed nonrespondents. Ph.D. Thesis; Department of Statistics, University of Wisconsin-Madison, 2001.
9. Rao, J.N.K.; Scott, A.J. On simple adjustments to chi-square tests with sample survey data. J. Annal. Stat. **1987**, *15*, 1–12.