

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34

# **Incidence of Missing Item Scores in Personality Measurement, and Simple Item-Score Imputation**

Joost R. van Ginkel, Leiden University  
Klaas Sijtsma, Tilburg University  
L. Andries van der Ark, Tilburg University  
Jeroen K. Vermunt, Tilburg University

January 14, 2009

First Author's address:  
Joost R. Van Ginkel  
Leiden University  
Faculty of Social and Behavioural Sciences  
Data Theory Group  
PO Box 9555  
2300 RB Leiden  
The Netherlands  
Tel: +31-(0)71-527 3620  
Email: [jginkel@fsw.leidenuniv.nl](mailto:jginkel@fsw.leidenuniv.nl)

## Abstract

The focus of this study was the incidence of different kinds of missing data problems in personality research and the handling of these problems. Missing-data problems were reported in approximately half of more than 800 articles published in three leading personality journals. In these articles, unit-nonresponse, attrition, and planned missingness were distinguished but missing item scores in trait measurement were reported most frequently. Listwise deletion was the most frequently used method for handling all missing-data problems. Listwise deletion is known to reduce the accuracy of parameter estimates and the power of statistical tests and often to produce biased statistical analysis results. This study proposes a simple alternative method for handling missing item scores, known as two-way imputation, which leaves the sample size intact and has been shown to produce almost unbiased results based on multi-item questionnaire data.

*Keywords: incidence of missing data, missing item scores, Two-way imputation, questionnaire data, multiple imputation of item scores.*

## Introduction

Multi-item questionnaires, inventories, and checklists—henceforth, generically called questionnaires—are widely used for measuring personality traits. Multiple items are used to cover all relevant aspects of a trait in an effort to measure the trait validly, and to control measurement error to a degree that the total score on the questionnaire is reliable. Examples of traits measured by means of multi-item questionnaires are obsessive-compulsive disorder, depression, and anxiety. The Obsessive-compulsive inventory (Foa, Kozak, Salkovskis, & Amir, 1998) is a well-known questionnaire for measuring obsessive-compulsive disorder, the Beck depression inventory II (e.g., Segal, Coolidge, Cahill, & O’Riley, 2008) measures depression, and the Beck anxiety inventory (e.g., Morin et al., 1999) measures anxiety.

Even when respondents have been instructed explicitly to respond to all items and not leave any responses open, data collection by means of multi-item questionnaires regularly suffers from missing item scores. Often the researcher is in the dark with respect to the reasons of this item nonresponse. In many cases, re-approaching respondents is an unrealistic option because of anonymity guarantee or financial or other restraints. Thus, the researcher often has to accept the incidence of the missing item scores and make a decision how to handle this problem in the statistical analysis of the data. One popular strategy is to leave out the cases that have at least one missing score and analyze only the complete cases. This strategy is called listwise deletion.

Our experience is that listwise deletion is an immensely popular method for handling missing item scores but it has a few serious drawbacks. By definition, it always reduces the sample size, which has the effect of reducing the accuracy of estimation and the power of statistical testing. In addition, under many circumstances listwise deletion may even cause more harm by producing biased statistical results (Little & Rubin, 2002; Schafer, 1997). For example, means and correlations may be distorted, which may affect the outcomes of methods such as the Student’s *t* test and factor analysis. Also, see Burton and Altman (2004), who corroborated the dominance of listwise deletion in the context of cancer research.

The large-scale application of listwise deletion suggests that researchers may not always realize the potentially damaging effects of listwise deletion on their research outcomes and also may not be aware of the availability of simple and

1 statistically superior methods for handling missing data that keep these damaging  
2 effects to a minimum. Thus, this study has two purposes. First, by means of a  
3 literature search we focus on the incidence of several kinds of missing-data problems  
4 that are reported in the literature on personality research. These missing-data  
5 problems also include missing item scores in multiple-item questionnaires, which  
6 constitute a large portion of the general missing-data problem. Also, we record the  
7 methods used in practice to handle missing-data problems. Second, we suggest a  
8 simple and statistically superior alternative to listwise deletion, that does not have the  
9 damaging effect of listwise deletion in multi-item trait measurement. We illustrate the  
10 method by solving the missing item-score problems in a real data set.

11

12 **Missingness Mechanisms and Real-Data Analysis**

13

14 An example using a real data set (Vorst, 1992; also, see Van der Ark, 2007) collected  
15 by means of a Dutch translation of the Adjective Checklist (ACL; Gough & Heilbrun,  
16 1980) may illustrate the problem of item nonresponse, which leads to missing item  
17 scores. The 218 items of the ACL are divided across 22 subscales; see Table 1. A  
18 sample of  $N = 433$  students from the University of Amsterdam provided ordered  
19 scores on a five-point rating scale, scored 0 (completely disagree) to 4 (completely  
20 agree). The data were completely observed; thus, there were no missing item scores.  
21 The completeness of the real data enabled us to manipulate mechanisms that created  
22 item nonresponse so as to illustrate what listwise deletion can do to the statistical  
23 results, but first we consider the complete data results.

24

25 INSERT TABLE 1 ABOUT HERE

26

27 Suppose a researcher uses the total score on the ACL Aggression subscale  
28 (items 101-110) and the ACL Dominance subscale (items 21-30) to test the  
29 hypothesis that aggressive people tend to be more dominant than non-aggressive  
30 people. To this end, (s)he uses a median split of the total scores on Aggression to  
31 divide the respondents into ‘aggressive’ respondents and ‘non-aggressive’  
32 respondents. The researcher is interested in the mean difference in the total  
33 Dominance score between aggressive and non-aggressive people. To test whether this  
34 difference is significant, (s)he performs a two-sample  $t$ -test with the dichotomized

1 aggression score as the independent variable, and the total Dominance score as  
2 dependent variable. The researcher is also interested in the range, the mean, and the  
3 reliability of the Dominance subscale in the total sample. Table 2 (first row) shows  
4 that Cronbach’s (1951) alpha equaled 0.807, and that the relationship between  
5 aggression and dominance was significant ( $p = .024$ ).

6  
7  
8

INSERT TABLE 2 ABOUT HERE

9         The statistical literature (Little & Rubin, 2002, p. 12; Schafer, 1997)  
10 distinguishes three mechanisms that may produce missing scores on variables.  
11 Listwise deletion always leads to a reduced sample size irrespective of which  
12 mechanism caused the missing item scores, but it leads to biased results under two of  
13 the mechanisms. Unfortunately, these are the mechanisms that are the most likely to  
14 cause missing-data problems in practical research. Thus, for a better understanding of  
15 the problems involved in using listwise deletion and the solutions of these problems, it  
16 is necessary to understand these three mechanisms. Each is explained next, and their  
17 effects on data analysis after the application of listwise deletion are illustrated using  
18 the ACL data.

19

20 *The Missing Completely at Random Mechanism*

21

22 The first mechanism produces missing item scores as if they constituted a simple  
23 random sample from all scores in the data. There is no relation to the value of the item  
24 score that is missing, nor to any other variable. In this case, the missing item scores  
25 are *missing completely at random* (MCAR; Little & Rubin, 2002, p. 12). This is the  
26 only situation in which listwise deletion is guaranteed not to result in biased  
27 outcomes. However, reduction of the sample size and its effects on accuracy and  
28 power are unavoidable.

29         The MCAR mechanism in the Dominance data was simulated by randomly  
30 drawing entries from the data matrix, which consisted of 433 rows (respondents) and  
31 10 columns (Dominance items), removing the item scores corresponding to these  
32 entries, and considering the resulting data matrix as suffering from item nonresponse.  
33 For this example, entries were drawn with a probability equal to .05 and without  
34 replacement; this produced a sample of 217 entries [433 (respondents) × 10 (items) ×

1 0.05 (probability) = 216.5], and the corresponding item scores were removed.  
2 Listwise deletion resulted in a 40% reduction of the sample; that is,  $N = 258$  complete  
3 cases were left for statistical analysis.

4 Because the reduced sample was a simple random sample drawn from the  
5 complete sample, we did not expect biased results. Table 2 (second row) shows that  
6 Cronbach's alpha dropped from 0.807 to 0.802, which reflects sampling error. The  
7 mean and the range of the test score were also similar to those found in the complete  
8 sample. However, a smaller sample size leads to a loss of power, which was apparent  
9 from a nonsignificant  $t$ -test compared to a significant result in the complete sample.  
10 Also, the mean difference has become smaller, which also reflects sampling error.  
11 Thus, listwise deletion may have important consequences for the outcomes of  
12 research.

13

14 *The Missing at Random Mechanism*

15

16 The second mechanism also produces missing item scores as if they constituted a  
17 random sample from the data, but the missingness is related to one or more observed  
18 variables in the data; hence, the missing item scores do not constitute a simple random  
19 sample. Missing scores are now said to be *missing at random* (MAR; Rubin, 1976;  
20 Little & Rubin, 2002, p. 12). The next example may further clarify the MAR  
21 mechanism.

22 Suppose we distinguish decent citizens from indecent citizens (e.g., due to  
23 hazardous traffic behavior, littering the street, not waiting in line at the bakery). A  
24 median split of the ACL Communality subscale total score produced groups of decent  
25 people and indecent people. Suppose that indecent people have a probability of not  
26 responding to items in the Dominance subscale that is three times as high as the  
27 corresponding probability for decent people. Thus, whether scores on dominance  
28 items are missing depends on the total score on Communality, which is an observed  
29 variable in the data. As this variable explains the missingness, it may be used to fix  
30 the missing data problem. Because listwise deletion ignores such explanatory  
31 variables, it now produces biased statistical results.

32 The MAR mechanism was simulated by randomly drawing 217 entries from  
33 the data (i.e., 5% missingness), such that respondents low on Communality had a  
34 probability of missing a Dominance-item score that was three times higher than

1 respondents high on Community. After the corresponding item scores were  
2 removed, listwise deletion resulted in a 39% reduction of the sample, leaving  $N = 265$   
3 cases for statistical analysis. Table 2 (third row) shows that Cronbach's alpha  
4 increased by 0.003, and that the  $t$ -test was not significant. The mean test score was  
5 similar to the mean test score in the complete-data example and the MCAR example.  
6 However, the maximally observed test score decreased from 40 to 38. Hence, the  
7 MAR mechanism produced results that are slightly worse than the MCAR  
8 mechanism.

### 9 10 *The Miscellaneous Category: Not Missing at Random Mechanisms*

11  
12 The third category contains all the mechanisms that produce missingness that is  
13 related to the value that is missing or to one or more variables that are not in the data  
14 of the study under consideration. These mechanisms produce missingness such that  
15 item scores are *not missing at random* (NMAR; Little & Rubin, 2002, p. 12). The  
16 problem here is that the researcher has no knowledge of the causes of the missingness,  
17 and thus is not in a position to solve the problem adequately. Because of the solution  
18 of NMAR problems requires knowledge that is inaccessible, one may resort to  
19 solutions assuming MAR in an effort to fix the problem as much as possible.

20 NMAR was simulated by removing 217 item scores (i.e., 5% missingness),  
21 such that for scores of 3 and higher, the probability of being missing was three times  
22 as high as for scores lower than 3. Table 2 (fourth row) shows that, compared to the  
23 original data, Cronbach's alpha increased by 0.011. The mean test score was  
24 underestimated. The maximum test score decreased from 40 to 38. The  $t$ -test is not  
25 significant.

### 26 27 Study 1: Incidence of Missing Data in Personality Measurement

28  
29 In Study 1, we investigated the frequency with which particular types of missing data  
30 were reported in articles discussing personality-trait measurement. Prior to discussing  
31 the results from the first study, we discuss the four types of missing data that were  
32 frequently reported: item nonresponse, unit nonresponse, attrition, and planned  
33 missingness. Because we already discussed item nonresponse, we now limit attention  
34 to *unit nonresponse*, *attrition*, and *planned missingness*.

1 Unit nonresponse occurs when a participant drawn into the sample refuses to  
2 take part in the investigation, so that for this person no observed data exist. De Leeuw  
3 and Hox (1988), Dillman (1991), and Groves and Couper (1998) have extensively  
4 studied the statistical handling of unit nonresponse.

5 Attrition occurs when participants drop out of a longitudinal study in which  
6 they are objected to repeated observation. Dropout may be due to loss of interest or  
7 motivation to proceed, having moved to another city, and in medical and health  
8 studies due to complete recovery, becoming too ill to further participate, or passing  
9 away as a result of the illness. Fleming and Harington (1991) and Andersen, Borgan,  
10 Gill, and Kleiding (1993) discuss methods for statistically dealing with attrition.

11 Planned missingness results from the researcher's intentional planning. For  
12 example, in a medical screening using multiple tests, for reasons of efficiency the  
13 researcher may not administer all tests to all participants. Eggen and Verhelst (1992)  
14 and Mislavy and Wu (1988) discuss statistical methods for handling planned  
15 missingness in the context of educational measurement.

#### 16 17 *Method*

18  
19 We used the following strategy for studying the incidence of missing-data problems in  
20 personality measurement. A total of 832 articles from six recent volumes (1995, 1997,  
21 2000, 2002, 2005, and 2007), four issues per volume, of three personality journals  
22 (*Psychological Assessment*, *Personality and Individual Differences*, and *Journal of*  
23 *Personality Assessment*) were screened for report of missing-data problems. The four  
24 issues per volume were selected as follows: *Psychological Assessment* is issued four  
25 times per year, *Personality and Individual Differences* is issued monthly (arbitrarily,  
26 the January, April, August, and December issues were selected), and *Journal of*  
27 *Personality Assessment* is issued six times per year (arbitrarily, the February, July,  
28 August, and December issues were selected). When multiple types of missingness  
29 were reported within the same article, the article was counted multiply. This yielded a  
30 total count of 927 cases within 832 articles.

#### 31 32 *Results*



1 Table 3 shows that 30% of the 927 cases pertained to item nonresponse (third  
2 column). Unit nonresponse and attrition are typical of survey studies and longitudinal  
3 studies, which are types of research that are not published as regularly in the three  
4 journals as personality measurement studies. Several articles specified the number of  
5 participants who provided incomplete score patterns but did not mention the type of  
6 missing data, and a few articles reported the removal of participants but not whether  
7 removal was due to missing scores or other reasons (e.g., random responding).  
8 Articles that mentioned nonresponse but did not mention the type of nonresponse  
9 were classified as 'Not Clear' (Table 3).

10

11

INSERT TABLE 3 ABOUT HERE

12

13 Table 4 shows descriptive statistics (mean, standard deviation, minimum, and  
14 maximum) of the proportion of incomplete score patterns computed across the 369  
15 cases where the proportion of incomplete cases was reported. The distribution of the  
16 proportion of incomplete score patterns is positively skewed, which means that most  
17 articles reported little missing data, and a small number of articles (6%) reported a  
18 large proportion of incomplete score patterns (30% or more). For item nonresponse,  
19 the percentage of incomplete item-score patterns on average equaled 9%. Thus, on  
20 average listwise deletion would result in a sample reduction of approximately 9%.  
21 Some articles reported the presence of missing item scores, but not the percentage of  
22 incomplete score patterns.

23

24

INSERT TABLE 4 ABOUT HERE

25

1 *Discussion*

2  
3 Almost half of the articles reported missing-data problems. Assuming that some  
4 articles failed to report such problems, the incidence of missing-data problems in  
5 personality measurement may even be greater. Item nonresponse was reported more  
6 often than other types of missing data. Item nonresponse occurs frequently in  
7 personality trait measurement using multi-item questionnaires. Item nonresponse is a  
8 serious problem in data analysis that calls for effective solutions that are easy to  
9 understand and implement.

10  
11 **Study 2: Handling Missing Data in Personality Measurement**

12  
13 In Study 2, we investigated the methods researchers in personality measurement  
14 typically use for handling missing-data problems.

15  
16 *Method*

17  
18 The observations were the 927 missing-data problems used in Study 1. The  
19 independent variable was missing-data type, which had six levels: unit nonresponse,  
20 attrition, item nonresponse, planned missingness, not clear, and none reported (Table  
21 3). The dependent variable was the method researchers in personality measurement  
22 use to handle missing-data problems. Seven principal methods for missing-data  
23 handling were found to be used in the 832 articles: follow-up, listwise deletion,  
24 available-case analysis, single imputation, direct maximum likelihood, variable  
25 deletion, and prorating. In addition, four variations or combinations of principal  
26 methods were identified: listwise deletion with a check for MCAR and MCAR not  
27 rejected; listwise deletion with a check for MCAR but MCAR rejected; Available  
28 case analysis with a check for MCAR and MCAR not rejected; and a combination of  
29 follow-up and listwise deletion with a check for MCAR). Also, two rest categories  
30 were identified and categorized as ‘other’ and ‘none reported’. Addition of these  
31 missing-data handling methods led to a dependent variable having  $7 + 4 + 2 = 13$   
32 levels. The seven principal methods were also used to handle item nonresponse. These  
33 methods and another method known as *multiple imputation* are discussed below.  
34 Some of the methods are illustrated using an incomplete-data example (see, Sijtsma &

1 Van der Ark, 2003), which is shown in Table 5. This data set contains the scores of 8  
2 fictitious respondents on 5 items.

3

4 INSERT TABLE 5 ABOUT HERE

5

6 *Follow-up.* Perhaps the best way to deal with missing data is re-approaching  
7 respondents with incomplete score patterns in an effort to obtain the scores that are  
8 missing. When successful, data that were initially missing become observed, and  
9 statistical analyses may be carried out without any problems, and without running the  
10 risk of obtaining biased results. For an example, see Huisman, Krol, and Van  
11 Sonderen (1998) who re-approached patients in a study with respect to the waiting list  
12 problem in orthopedic practices. Unfortunately, however, due to many different  
13 restraints, in many studies follow-up is not feasible.

14 *Listwise deletion.* Consider the data in Table 5. Suppose a researcher plans  
15 computing Cronbach's alpha for the total score on the items  $X_1$ ,  $X_2$ , and  $X_3$ , and the  
16 correlation between the items  $X_4$  and  $X_5$ . Listwise deletion uses cases 2, 4, and 7 for  
17 both computing Cronbach's alpha and the correlation. Advantages of listwise deletion  
18 are that statistical analyses can be done without any modifications on the data and that  
19 all statistical analyses are done on the same subsample. Disadvantages are that the  
20 reduction of the sample size results in a loss of estimation precision and a reduced  
21 power in hypothesis testing. Furthermore, unless the missing scores are MCAR  
22 statistics may be biased. Listwise deletion may be preceded by a check whether  
23 MCAR is a reasonable assumption. This check may entail testing whether respondents  
24 with completely observed item-score patterns and respondents with incomplete or  
25 blank item-score patterns differ significantly with respect to demographic variables  
26 such as gender or ethnicity. For example, when the background variable 'age' is  
27 observed for all respondents, a two-sample  $t$ -test may be used to test whether  
28 respondents with complete score patterns differ systematically with respect to age  
29 from respondents with incomplete score patterns. For categorical background  
30 variables, such as gender, chi-square tests may be used. See, for example, Hishinuma  
31 et al. (2000), and Cole, Hoffman, Tram, and Maxwell (2000) who used this strategy  
32 for checking the MCAR assumption.

33 *Available-case analysis.* Loss of power may be reduced when all cases are  
34 used in the statistical analysis, which have observed values on the variables that are

1 effective in the analyses. This option is called available-case analysis. When applied  
2 to the data from Table 5, available-case analysis uses cases 1, 2, 4, 6, 7, and 8 for  
3 computing Cronbach's alpha for the total score on the items  $X_1$ ,  $X_2$ , and  $X_3$ . For  
4 computing the correlation between the items  $X_4$ , and  $X_5$ , available-case analysis uses  
5 cases 2, 3, 4, and 7. Available-case analysis (Little & Rubin, 2002, pp. 53-54) is the  
6 default option for missing-data handling in SPSS (2008).

7       Compared to listwise deletion, a disadvantage of available-case analysis is that  
8 different statistical analyses that use different variables may be based on (partly)  
9 different sub-samples with different sample sizes. A disadvantage shared with listwise  
10 deletion is that statistics may be biased unless the missingness mechanism is MCAR.  
11 Kim and Curry (1977) showed that available-case analysis is superior to listwise  
12 deletion when correlations among variables are modest. Haitovsky (1968) and Azen  
13 and Van Guilder (1981) showed that listwise deletion is superior to available-case  
14 analysis when correlations among variables are large. Little and Rubin (2002, p. 55)  
15 argued that both options are generally unsatisfactory.

16       Because listwise deletion and available-case analysis result in a loss of power  
17 and possibly biased results, researchers should be cautious using these methods. It  
18 may be recommended to use these methods only when the reduced sample is large  
19 and when it has been checked whether there are systematic differences on the  
20 background variables between the completely observed cases and the incomplete  
21 cases, so that the MCAR assumption at least is plausible.

22       *Single imputation.* Single imputation replaces the missing scores by plausible  
23 scores, so that cases that have missing scores can be included in the statistical  
24 analyses. We discuss two possibilities.

25       Deterministic imputation replaces the empty cells in the data matrix by  
26 estimates of the item scores. For example, Saggino and Kline (1995) replaced each  
27 missing score on variable  $X$  by the sample mean of  $X$  based on the available scores,  
28 and Sheviin and Adamson (2005) replaced each missing score by the expected value  
29 from a regression model. Table 6 (upper left panel) shows how variable-mean  
30 imputation is done in the incomplete-data example in Table 5. The imputed scores are  
31 derived readily by computing the means for each variable (last row). For example, the  
32 imputed score on variable  $X_1$  is computed as  $(2 + 3 + 4 + 1 + 5 + 1 + 3)/7 = 2.71$ . Note  
33 that the resulting imputed scores are not necessarily integer scores. Depending on the  
34 application, imputed scores may be analyzed as real numbers (e.g., as in factor

1 analysis, which treats rating-scale scores as continuous) or they may be rounded to the  
2 nearest feasible integer (e.g., as in item analysis using item response models, which  
3 treat rating-scale scores as discrete).

4  
5 INSERT TABLE 6 ABOUT HERE  
6

7 Table 6 (upper right panel) also shows the completed data set that results from  
8 deterministic regression imputation. Imputations were done using SPSS 16.0 (Analyze,  
9 Missing Value Analysis). The imputed scores are less easily derived because the  
10 computation procedure that SPSS uses is rather complicated.

11 The advantage of deterministic imputation is that it provides the researcher  
12 with a complete data set, which may be used for further statistical analysis. A  
13 disadvantage is that variances and covariances are biased downward (Schafer, 1997,  
14 p. 2).

15 Stochastic imputation improves upon deterministic imputation by imputing a  
16 value that includes a random error; For example, in regression imputation the imputed  
17 value includes a normally distributed random error with variance equal to the error  
18 variance of the regression model. Thus, the imputed values have the same variance as  
19 the observed scores. Stochastic imputation keeps the covariance structure intact but in  
20 subsequent statistical analyses the imputed scores are treated as if they were observed  
21 without taking the uncertainty about these imputed values into account. As a result,  
22 the standard errors of the statistics are too small.

23 Table 6 (lower left panel) shows how stochastic variable mean imputation is  
24 done. Here, the imputed values are random draws from from a normal distribution  
25 rather than a mean substitution. For example, the imputed score on variable  $X_1$  is a  
26 random draw from a normal distribution with a mean of 2.71 and a standard deviation  
27 of 1.50 (last row).

28 Because the detailed explanation of how the computations for both  
29 deterministic and stochastic regression imputation are carried out would be too  
30 involved, we only show the syntax that performs the imputations in SPSS. Here, it is  
31 assumed that the incomplete data set is named example.sav and located in the directory  
32 C:\imputation\, and that the completed data files are called deterministic.sav and  
33 stochastic.sav. The resulting syntax file is shown in Figure 1. Note that the 12th line  
34 (SET SEED = 2 .) is only added to reproduce the results from the example (Table 6) for

1 stochastic regression imputation. To obtain imputed values that differ from the  
2 example, this line may be removed.

3  
4 INSERT FIGURE 1 ABOUT HERE

5  
6 *Multiple imputation.* Multiple imputation improves upon stochastic imputation  
7 by substituting multiple random values (i.e., not necessarily integer scores) for each  
8 missing score, resulting in several plausible complete versions of the data. These  
9 completed data sets are then analyzed by standard statistical procedures, and the  
10 results are combined into one overall result, using rules proposed by Rubin (1987,  
11 chap. 3). Schafer (1997, p. 106) recommends doing the statistical analyses on three,  
12 four, or five completed data sets.

13 An advantage of multiple imputation compared to single imputation is that  
14 statistical analysis takes the uncertainty about the missing data into account, so that  
15 standard errors of statistics are not biased downwards. Moreover, whereas listwise  
16 deletion and available-case analysis only lead to valid inferences when scores are  
17 MCAR, multiple imputation also leads to valid inferences when scores are MAR. A  
18 disadvantage of multiple imputation is that the method is rather involved and only  
19 available in software packages that are not frequently used among personality  
20 researchers. Examples of software are SAS 8.1, in the procedure PROC MI (Yuan,  
21 2000), S-plus 8 for Windows (2007), AMOS 6.0 (Arbuckle & Wothke, 2006), the  
22 stand-alone program NORM (Schafer, 1998), ICE in Stata 10.0 (StataCorp, 2007), the  
23 MICE library in S-plus, and the stand-alone program WinMICE V1.0 (Jacobusse,  
24 2005).

25 Table 7 shows 5 completed versions of the incomplete data set in Table 5.  
26 Multiple imputation was done using the program NORM (Schafer, 1998). Cronbach's  
27 alpha for the total score on the items  $X_1$ ,  $X_2$ , and  $X_3$  may be obtained as the mean of  
28 the five alpha values obtained from the five imputed data sets. The same goes for the  
29 correlation between the variables  $X_4$  and  $X_5$ . To test the significance of the correlation,  
30 an overall standard error has to be computed across the five imputed data sets using  
31 Rubin's (1987) rules. See Rubin (1987, Chap. 3) for an extensive discussion of these  
32 rules.

33  
34 INSERT TABLE 7 ABOUT HERE

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34

*Direct maximum likelihood estimation.* Direct maximum likelihood estimation (e.g., Allison, 2002) entails estimating the parameters from a statistical model while ignoring the unobserved scores but without deleting cases. Thus, unlike listwise deletion and available-case analysis, direct maximum likelihood estimation uses all observed item scores instead of using only the scores of respondents with complete item-score patterns. The method is used for the estimation of, for example, item response theory models, latent class models, and structural equation models. An advantage of direct maximum likelihood estimation is that all cases are used to estimate the model. A disadvantage of the method is that, like most multiple imputation methods, it is relatively complex and can only be used in nonstandard statistical procedures and nonstandard statistical software packages. The method cannot be used in popular procedures like principal components analysis and analysis of variance (ANOVA). Moreover, SPSS (2008) does not allow using the method even for procedures that are suited for it, such as factor models or loglinear models.

*Prorating test scores.* Prorating test scores entails computing a respondent's test score across his/her observed scores and then rescaling the resulting score. Together with the total scores for respondents with complete data, these resulting scores are used as dependent variable in statistical analyses. In Table 5, the test score of person 2 is computed as  $3 + 5 + 4 + 5 + 5 = 22$ , and the *prorated* test score of person 1 is computed as  $[(1 + 1 + 2) / 3] \times 5 = 6.67$ .

This method does not explicitly impute scores but is equivalent to substituting for each missing value the person mean across a respondent's available scores. This procedure is common practice, and is even recommended in manuals of many personality-trait questionnaires (e.g., Bracken & Howell, 2004; Hare, 2003). However, from a statistical point of view, prorating test scores is a suboptimal method. First, it does not take the differences between item means into account. Second, because the mean test score across the remaining items does not have an error component, the variance of the test score is biased downwards.

*Variable deletion.* Variable deletion leaves out variables with missing scores from the statistical analysis. Thus, for items it is the counterpart of listwise deletion. The missing-data literature does not explicitly mention this procedure as a useful method but researchers often use it. For example, when information on gender is missing for some respondents a researcher may decide not to use gender as an

1 independent variable in statistical tests but to use it only for describing the  
2 demographic characteristics of the sample. See, for example, Watson, et al. (2007),  
3 who reported that “The sample consisted of 376 women and 121 men (2 participants  
4 did not specify their sex)”. Another example of variable deletion may concern a  
5 particular item, which has so many missing values that the researcher may decide to  
6 leave it out of the reliability analysis and compute test scores across the remaining  
7 items. In the data example of Table 5, a researcher may decide that item  $X_4$  has too  
8 many missing values to be useful for any statistical analysis. Thus, (s)he may decide  
9 not to compute the correlation between items  $X_4$  and  $X_5$ . Because variable deletion  
10 does not result in a selective dropout of respondents, it gives valid results in statistical  
11 analyses but limits the substantive meaning of the research.

### 12 13 *Results*

14  
15 Table 8 shows that listwise deletion is by far the most frequently used missing-data  
16 method, followed by available-case analysis. Single imputation was used 19 times,  
17 and multiple imputation was not used at all. Some studies used several methods of  
18 handling nonresponse. Each method was counted separately, leading to a total of 1025  
19 cases of missing-data handling rather than 927 as in Table 3. Only few studies  
20 checked whether MCAR was plausible prior to deleting the cases from the analyses.  
21 All of these studies, regardless of the outcome of this check, conducted the statistical  
22 analyses based on the complete cases, and only in the discussion section they  
23 mentioned that the sample was probably not completely representative, thus resulting  
24 in limited generalizability.

25 Two articles reported a combination of follow-up and listwise deletion  
26 preceded by a check for MCAR (row 12). Specifically, Iversen and Rundmo (2002)  
27 reported that “A control study was conducted to find out if the group of respondents  
28 who had replied to the questionnaire differed significantly from those who did not.  
29 Fifty subjects were contacted by phone and interviewed using the same questionnaire  
30 as in the survey. Results from this study showed that the final sample was  
31 representative of the population of Norwegian drivers with regard to age, gender and  
32 education.”

33  
34 INSERT TABLE 8 ABOUT HERE



1

2 *Discussion*

3

4 Personality-trait measurement using multiple-item questionnaires predominantly uses  
5 listwise deletion for handling missing data problems. The popularity of listwise  
6 deletion probably resides in its simplicity but researchers seem to be unaware of its  
7 potential problems. We give two possible explanations. First, it may be incorrectly  
8 assumed that missing scores make a score pattern useless so that the pattern better be  
9 discarded from the data analysis. Second, it may be incorrectly assumed that deleting  
10 cases only reduces power, whereas the bias resulting from nonresponse may not be  
11 appreciated. We noted that missing data were often discussed as if they were nothing  
12 more than a nuisance in the data-collection process, which could simply be remedied  
13 by collecting enough data so that after listwise deletion enough cases were left for  
14 analysis.

15 Sometimes, listwise deletion is a good solution for missing item-score  
16 problems. For example, respondents who have almost no observed data may be  
17 discarded from the data analyses. Also, when only a few respondents out of a  
18 relatively large sample have incomplete item-score records leaving them out of the  
19 analysis has little effect on the outcomes of statistical analysis. For example, Boyd-  
20 Wilson, Walkey, McClure, and Green (2000) deleted two incomplete cases from a  
21 total sample of  $N = 205$ . However, listwise deletion was used so frequently that it  
22 seems safe to conclude that it is often used inappropriately.

23 The popularity and dominance of listwise deletion seems to have the effect of  
24 hiding simple, user-friendly and statistically superior alternatives for the handling of  
25 item nonresponse from the researchers' statistical toolbox. Given the availability of  
26 such alternatives and the established inferiority of listwise deletion in many research  
27 situations, next we discuss an attractive method for handling item nonresponse in  
28 multi-item questionnaires for personality-trait measurement.

29

30 A Simple Method to Handle Item Nonresponse in Multi-Item Questionnaire Data

31

32 For multiple-item questionnaire data, the most promising simple imputation method is  
33 *two-way multiple imputation with error* (abbreviated Method TW; Little & Su, 1987,  
34 discussed the core of Method TW in the context of incomplete longitudinal data, and

1 Bernaards & Sijtsma, 2000, proposed using the method for questionnaire data; also  
 2 see Van Ginkel et al., 2007a; 2007b, and Van Ginkel, Van der Ark, Sijtsma, &  
 3 Vermunt, 2007). In the Appendix we show how Method TW can be used by means of  
 4 SPSS (2008).

5 Method TW is based on a typical ANOVA layout. We assume that the scores  
 6 of  $N$  persons to  $J$  items measuring a single personality trait are incomplete. Let  $PM_i$   
 7 denote the mean item score of person  $i$  based on his/her available item scores, let  $IM_j$   
 8 denote the mean score of item  $j$  based on all scores available for this item, and let  $OM$   
 9 be the overall mean of all available item scores in the  $N \times J$  data matrix. A  
 10 deterministic imputation method may use  $TW_{ij} = PM_i + IM_j - OM$  to impute a score  
 11 for a missing value in cell  $(i, j)$  of the data matrix, and a probabilistic imputation  
 12 method adds an error term  $\varepsilon_{ij}$  and then imputes  $TW_{ij}^* = TW_{ij} + \varepsilon_{ij}$ . Depending on the  
 13 application, imputed  $TW_{ij}^*$  scores are analyzed as real numbers (e.g., as in factor  
 14 analysis) or rounded to the nearest feasible integer (e.g., as in item analysis using item  
 15 response models).

16 The computation of  $TW_{ij}^*$  is illustrated next using the data example in Table 5  
 17 for person 5 and variable  $X_1$ . It may be verified that  $PM_5 = (3 + 3 + 4) / 3 = 3.33$ ,  $IM_1$   
 18  $= (2 + 3 + 4 + 1 + 5 + 1 + 3) / 7 = 2.71$ , and  $OM = 95 / 33 = 2.88$ ; hence,  $TW_{51} = 3.33$   
 19  $+ 2.71 - 2.88 = 3.16$ . The other values of  $TW_{ij}$  from the example in Table 5 are shown  
 20 in Table 9.

21

22 INSERT TABLE 9 ABOUT HERE

23

24 Next, the error  $\varepsilon_{ij}$  is drawn from a normal distribution with mean 0 and  
 25 variance  $S_\varepsilon^2$ ;  $S_\varepsilon^2$  is the error variance in the observed data, which is computed as  
 26 follows. First, for each observed item score  $X_{ij}$  the corresponding  $TW_{ij}$  score is  
 27 computed. The  $TW_{ij}$  scores are considered to be the expected scores of the two-way  
 28 model, had the  $X_{ij}$  scores been missing. Second, the sum of the squared differences,  
 29  $(X_{ij} - TW_{ij})^2$ , is computed across all observed cells, and this sum is divided by the  
 30 number of observed scores minus 1 (denoted by  $M$ ; in Table 5,  $M = 33 - 1 = 32$ ).  
 31 Thus, we find that  $S_\varepsilon^2 = \sum \sum (X_{ij} - TW_{ij})^2 / M$ .

1 Multiple imputation based on five independent draws of the error is done as  
2 follows. For the data in Table 5 the error variance equals 0.901 (it may be noted that  
3 for computing a  $TW_{ij}$  score, the corresponding observed  $X_{ij}$  score is treated as  
4 missing; as a result, the person and item means vary with each cell  $(i, j)$ , and the  
5 person and item means in Table 9 cannot be used throughout the computation of the  
6 error variance. These details are ignored here). Assume that five randomly drawn  
7 error terms are:  $\varepsilon_{51}^{(1)} = -0.1601879$ ,  $\varepsilon_{51}^{(2)} = -1.0220348$ ,  $\varepsilon_{51}^{(3)} = 0.4451876$ ,  $\varepsilon_{51}^{(4)} =$   
8  $2.5191623$ , and  $\varepsilon_{51}^{(5)} = -0.6389984$ . For producing consecutive data matrices, each of  
9 these values is added to  $TW_{51} = 3.17$ , which yields five different values (rounded to  
10 two decimals):  $TW_{51}^* = 3.01, 2.15, 3.61, 5.69, \text{ and } 2.53$ , respectively. Each of these  
11 values is imputed in the data matrix in Table 5 (thus treating scores as continuous).  
12 The same procedure is followed for the other missing values (not shown here), which  
13 yields five different completed data sets. Statistical analyses are done on all five data  
14 sets separately, and the results are combined using Rubin's (1987, chap. 3) rules.

15 Simulation results (Van Ginkel et al., 2007a; 2007b, Van Ginkel, Van der Ark,  
16 Sijtsma, & Vermunt, 2007) have shown that Method TW produces statistical results  
17 with very little or no bias at all, even when missing item scores are NMAR and the  
18 percentage of missing item scores increases up to 15% (in these studies, this  
19 corresponded to only 4% completely observed cases on average). A plausible  
20 explanation why Method Two-Way works so well in case of NMAR is because  
21 multiple items are used to measure the same construct. Even if some extreme NMAR  
22 missingness results in many missing item scores for certain respondents, these  
23 respondents will usually have responded to some items measuring the same construct.  
24 The observed item scores contain enough information to predict the missing item  
25 scores reasonably well. Only in case of extremely high percentages of missingness,  
26 Method Two-Way will result in biased estimates (see, Van Buuren, 2009). This is an  
27 important finding implying that a researcher may safely use Method TW to impute  
28 item scores in multiple-item questionnaires for measuring personality traits.

29 To illustrate the usefulness of Method TW, we simulated item nonresponse in  
30 the multiple-item ACL Dominance subscale (Table 1) for item scores that were either

1 MCAR, MAR<sup>1</sup>, or NMAR, thus producing three different incomplete data sets. We  
2 used Method TW to impute scores in each of the three data sets, and computed the  
3 values of Cronbach's alpha, the mean test score, the minimum and maximum  
4 observed test scores, and the *t*-test, with Aggression as the independent variable and  
5 the Dominance test score as the dependent variable (Table 10).

6  
7 INSERT TABLE 10 ABOUT HERE  
8

9 Almost all results produced by multiple imputation using Method TW were  
10 closer to the results produced by the complete data than the results produced by  
11 listwise deletion (cf. Table 1). For the MAR data set, the maximum test score was  
12 underestimated, but less than for listwise deletion (cf. Table 1, fourth column). For the  
13 three completed data sets, the *t*-test (last three columns) was significant, as in the  
14 original data.

15 First, it may be noted that when a test contains more than one subscale,  
16 Method TW may be applied to each subscale separately. Two other versions of  
17 Method TW, not discussed here, use the multidimensionality of the data for imputing  
18 scores; see Van Ginkel et al. (2007b) for more details. Second, Method TW should be  
19 applied only if  $PM_i$  can be interpreted as an indicator of the trait level of person  $i$   
20 (method TW capitalizes on each of the  $J$  items holding information on the other  
21 items).  $PM_i$  cannot be interpreted as an indicator of the trait level of person  $i$  if items  
22 are included that do not measure the intended trait, such as gender or social economic  
23 status, or if a respondent has excessively many missing values. In the former case,  
24 other methods such as multiple imputation under the latent class model may be used  
25 (Vermunt, Van Ginkel, Van der Ark, & Sijtsma, in press), and in the latter case such  
26 exceptional cases may be removed before method TW is used.

27  
28 

### General Discussion

  
29

30 Item nonresponse occurs frequently in personality measurement. Even though  
31 multiple imputation is a highly recommended procedure in the statistical literature for

---

<sup>1</sup> It may be noted that even though the missingness only depends on the fully observed variable 'Communality group', the default application of Method Two-Way does not impute scores separately for 'Communality group = 1' and 'Communality group = 2'. Therefore, technically, Method Two-Way treats this condition as NMAR.

1 dealing with item nonresponse, this method appears to be used rarely if ever in  
2 personality measurement. Instead, the inferior listwise deletion method is by far the  
3 most popular method for handling missing item scores.

4         The screening of three leading personality journals underlined the need for  
5 simple, user-friendly and statistically correct methods to deal with item nonresponse  
6 in questionnaire data. Method TW has these properties and may be used for the  
7 imputation of item scores. SPSS macros for multiple item-score imputation are  
8 available as freeware from the Internet (<http://www.uvt.nl/mto/software2.html>; Van  
9 Ginkel & Van der Ark, 2005a; 2005b). In an empirical-data example, it was shown  
10 that method TW accurately recovered several statistics typical of the psychometric  
11 analysis of questionnaire data. Thus, method TW may be a good alternative for  
12 listwise deletion and other missing-data handling methods for handling missing item  
13 scores in personality measurement. Method TW is appropriate for multi-item  
14 questionnaire data, in which the items all measure aspects of one underlying  
15 personality trait and a total score is typically used for measuring individuals but the  
16 method may also be extended to multidimensional questionnaire data.

## APPENDIX

SPSS syntax is available to conduct the following types of statistical analyses on test data with missing item scores using Method TW.

1. *Computation of a statistic without standard error* (e.g., reliability statistics such as Cronbach's alpha and corrected item-total correlations; descriptive statistics such as the mean, standard deviation, median, maximum, and minimum; correlation coefficients, loadings from factor analysis). As an example we show how to compute Cronbach's alpha for a dominance test containing 10 items.
2. *Computation of a statistic with standard error*. Note that in several cases SPSS does not provide standard errors and they have to be computed by the researcher. As an example we show how to compute the mean score on a dominance test containing 10 items, its standard error, and 95% confidence interval.
3. *All t-tests and univariate regression analyses* can be computed in a straightforward way. As an example, we show how to compare the mean scores on a dominance test of a group of non-aggressive and a group of aggressive respondents using a two-sample *t*-test.
4. For *other analyses* (multivariate regression, multi-level analysis, ANOVA, significance tests for correlations, mixed models) the procedures are more involved and we refer to Van Ginkel (2006) for detailed information.

Statistical analyses that cannot be performed include MANOVA and structural equation models.

The necessary files for the exemplary statistical analyses can be obtained from <http://www.uvt.nl/mto/software2.html> in the zip file *imputation.zip*, which contains four files:

- *ACL.sav*: An SPSS data file containing the item scores of 433 persons to 10 dominance items (V021 to V030), 5% of the scores are missing (MCAR); and their scores on variable *Naggress* (score 1 indicates non-aggressive behavior, score 2 indicates aggressive behavior).

1 • imputation.sps: An SPSS syntax file performing statistical analyses on the  
2 incomplete data file ACL.sav, using Method TW<sup>2</sup>.

3 • tw.sps: An SPSS syntax file containing preprogrammed macro tw.

4 • mi.sps: An SPSS syntax file containing preprogrammed macro mi.

5 These four files should be unpacked and moved to the same directory. Without  
6 loss of generality we assume that this directory is called C:/imputation/. The  
7 analyses are performed by running imputation.sps, which is discussed next.

8

9 The file imputation.sps contains four steps.

10 • *Step 1: Preliminary commands* (lines 1-7). Determining the working directory  
11 (lines 4-5). If the unzipped files are not in C:/imputation/ the FILE HANDLE  
12 command (line 5) should be modified before use. Line 7 ensures that the  
13 results in the Appendix are reproduced exactly; this line should be removed if  
14 imputation.sps is modified for other data sets. Line 7 suppresses the printing of  
15 syntax commands in the output. The command prevents that the many syntax  
16 commands from mi.sps and tw.sps are printed in the output.

17 • *Step 2: Creating five completed data sets* (lines 9-16). Line 13 reads the  
18 preprogrammed macro tw.sps. Five completed versions of acl.sav are created  
19 by the command TWOWAY. Subcommand /SELECT specifies the items to  
20 which Method TW is applied and subcommand /M specifies the number of  
21 required completed data sets; here  $M = 5$ . Running TWOWAY results in a single  
22 SPSS data file containing five completed versions of ACL.sav. This file,  
23 which is automatically called ACL\_imp.sav, contains all five completed  
24 datasets appended one after another. An additional variable called  
25 imputation\_# has been added, which indicates the data set number.

26 • *Step 3: Conducting statistical analysis* (lines 18-56). First, data file  
27 ACL\_imp.sav is read and split into five separate data sets (lines 20-22). In  
28 SPSS, the split file option may be found under task bar: Data, Split File.  
29 Second, five Cronbach's alphas are computed using the command RELIABILITY  
30 (line 31). RELIABILITY is preceded by the command OMS and followed by the

---

<sup>2</sup> This file is based on the package tw.zip (Van Ginkel & Van der Ark, 2005a, 2005b; Van Ginkel, 2006). This package is more general than the syntax presented here and has an extensive manual. To allow a brief yet concise explanation of Method TW, we have modified these general files and collected them in a single syntax file.

1 command OMSEND. These commands direct SPSS output into an SPSS data  
2 file<sup>3</sup>. The resulting file reliability.sav contains the five values of Cronbach's  
3 alpha. Similarly, the mean test score and the standard deviation are computed  
4 using DESCRIPTIVES and the output is directed to descriptives.sav (lines 42-  
5 44), and the *t*-test is performed and the output is directed to ttest.sav (lines 46-  
6 56).

- 7 • *Step 4: Combining the results of the five statistical analyses* (lines 58-86).  
8 First, the five Cronbach's alphas, collected in reliability.sav, are combined (lines  
9 60-62). The Cronbach's alpha that should be reported is obtained by simply  
10 taking the mean of the Cronbach's alphas of the five data sets. The output  
11 shows that Cronbach's alpha equals .8105. Second, the mean test scores  
12 (Mean) and standard deviations (Std.Deviation), collected in descriptives.sav, are  
13 combined (lines 64-74). This is a little bit more involved. The standard error of  
14 the mean is not provided by SPSS and must be computed separately as  
15  $S.E.Mean = Std.Deviation / \sqrt{N}$  (line 66). Furthermore, the even lines in  
16 descriptives.sav contain no information and they are removed (line 65). The  
17 command RULESMI gives the correct combination of the statistic and standard  
18 error. The output shows that the mean test score equals 24.398, its standard  
19 error equals 0.292, and the 95% confidence interval is [23.825; 24.972]; the  
20 remaining statistics (*t* statistic, *df*, and *p*-value) can be ignored here. Third, in a  
21 similar way the results of the *t*-test are combined (lines 76-87). Note that  
22 ttest.sav contains the results for both 'equal variances assumed' and for 'equal  
23 variances not assumed' whereas we are only interested in *t*-tests where equal  
24 variances are assumed. The other results are deleted in line 77. For the  
25 command RULESMI the difference in mean test scores (MeanDifference; line 84)  
26 and its standard error (Std.ErrorDifference; line 85) are provided. The number of  
27 degrees of freedom in a two-sample *t*-test equals  $N-2 = 433-2 = 431$  (line 86).  
28 The output shows that the difference in mean test scores equals -1.201 with  
29 standard error 0.589. The corresponding *T* statistic equals  $T = -2.039$ ,  $df =$   
30  $390.652$ ,  $p = .042$ , indicating a significant difference between aggressive and  
31 non-aggressive respondents.  
32

---

<sup>3</sup> For other analyses, other OMS options may have to be specified, which can be found under task bar: Utilities, OMS Control Panel.



## REFERENCES

- 1  
2  
3 Allison, P. D. (2002). *Missing data*. Thousand Oaks, CA: Sage.  
4  
5 Andersen, P. K., Borgan, Ø., Gill, R. D., & Keiding, N. (1993). *Statistical models*  
6 *based on counting processes*. New York: Springer.  
7  
8 Arbuckle J. L., & Wothke, W. (2006). AMOS 6.0 [Computer software]. Chicago:  
9 Smallwaters.  
10  
11 Azen, S., & Van Guilder, M. (1981). Conclusions regarding algorithms for handling  
12 incomplete data. *1981 Proceedings of the Statistical Computing Section*.  
13 *American Statistical Association*, 53-56.  
14  
15 Bernaards, C. A., & Sijtsma, K. (2000). Influence of imputation and EM methods on  
16 factor analysis when item nonresponse in questionnaire data is nonignorable.  
17 *Multivariate Behavioral Research*, 35, 321-364.  
18  
19 Boyd-Wilson, B. M., Walkey, F. H., McClure, J., Green, D. E. (2000). Do we need  
20 positive illusions to carry out plans? Illusion - and instrumental coping.  
21 *Personality and Individual Differences*, 29, 1141-1152.  
22  
23 Bracken, B. A., & Howell, K. (2004). *Clinical Assessment of Depression:*  
24 *Professional manual*. Odessa, FL: Psychological Assessment Resources.  
25  
26 Burton, A., & Altman, D.G. (2004). Missing covariate data within cancer prognostic  
27 studies: a review of current reporting and proposed guidelines. *British Journal*  
28 *of Cancer*, 91, 4-8.  
29  
30 Cole, D. A., Hoffman, K., Tram, J. M., & Marwell, S. E. (2000). Structural  
31 differences in parent and child reports of children's symptoms of depression  
32 and anxiety. *Psychological Assessment*, 12, 174-184.  
33  
34 Cronbach, J. L. (1951). Coefficient alpha and the internal structure of tests.

1           *Psychometrika*, 16, 297-334.

2

3       De Leeuw, E. D., & Hox, J. J. (1988). Response stimulating factors in mail surveys.  
4           *Journal of Official Statistics*, 4, 241-249.

5

6       Dillman, D. A. (1991). The design and administration of mail surveys. *Annual Review*  
7           *of Sociology*, 17, 225-249.

8

9       Eggen, T. J. H. M., & Verhelst, N. D. (1992). *Item calibration in incomplete testing*  
10           *designs*. (Measurements and Research Department Reports 92-3). Arnhem,  
11           The Netherlands: Cito.

12

13       Fleming, T. R., & Harrington, D. P. (1991). *Counting processes and survival analysis*.  
14           New York: Wiley.

15

16       Foa, E. B., Kozak, M. J., Salkovskis, P. M., Coles, M. E., & Amir, N. (1998). The  
17           validation of a new Obsessive-Compulsive Disorder scale: The Obsessive  
18           Compulsive Inventory. *Psychological Assessment*, 10, 206–214.

19

20       Gough, H. G., Heilbrun, A. B. (1980). *The Adjective Check List, manual 1980 edition*.  
21           Consulting Psychologists Press, Palo Alto, CA: Consulting Psychologists  
22           Press.

23

24       Groves, R. M., & Couper, M. P. (1998). *Nonresponse in household interview surveys*.  
25           New York: Wiley.

26

27       Haitovsky, Y. (1968). Missing data in regression analysis. *Journal of the Royal*  
28           *Statistical Society, Series B*, 67-81.

29

30       Hare, R. D. (2003). *Manual for the Revised Psychopathy Checklist* (2nd ed.). Toronto,  
31           Ontario, Canada: Multi-Health Systems.

32

33       Hishinuma, E. S., Andrade, N. N., Johnson, R. C., McArdle, J. J., Miyamoto, R. H.,  
34           Nahulu, L. B., Makini Jr, G. K., Yuen, N. Y. C., Nishimura, S. T., McDermott

- 1           Jr, J. F., Waldron, J. A., Luke, K. N., & Yates, A. (2000). Psychometric  
2           properties of the Hawaiian culture scale - adolescent version. *Psychological*  
3           *Assessment, 12*, 140-157.  
4
- 5           Huisman, M., Krol, B., & Van Sonderen, F.L.P. (1998). Handling missing data by re-  
6           approaching nonrespondents. *Quality & Quantity, 32*, 77-91.  
7
- 8           Iversen, H., & Rundmo, T. (2002). Personality, risky driving and accident  
9           involvement among Norwegian drivers. *Personality and Individual*  
10          *Differences, 33*, 1251-1263.  
11
- 12          Jacobusse G. W. (2005). WinMICE V1.0 The WinMICE application, a standalone  
13          software tool for multiple imputation when data have a multilevel structure  
14          [Computer software]. Retrieved September 3, 2008 from  
15          <http://web.inter.nl.net/users/S.van.Buuren/mi/html/mice.htm>  
16
- 17          Kim, J. O., & Curry, J. (1977). The treatment of missing data in multivariate analysis,  
18          *Sociological Methods and Research, 6*, 215-240.  
19
- 20          Little, R. J. A., & Rubin, D. B. (2002). *Statistical analysis with missing data* (2nd  
21          ed.). New York: Wiley.  
22
- 23          Little, R. J. A., & Su, H. L. (1989). Item nonresponse in panel surveys, In D.  
24          Kasprzyk, G. Duncan, & M. P. Singh (Eds.), *Panel surveys* (pp. 400-425).  
25          New York: Wiley.  
26
- 27          Mislevy, R. J., & Wu, P. K. (1988). *Inferring examinee ability when some item*  
28          *responses are missing*. (Research Report RR-88-48-ONR). Princeton, NJ:  
29          Educational Testing Service.  
30
- 31          Morin, C. M., Landreville, P., Colecchi, C., McDonald, K., Stone, J. & Ling, W.  
32          (1999). The Beck anxiety inventory: psychometric properties with older  
33          adults. *Journal of Clinical Geropsychology, 5*, 19-29.  
34

- 1 Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63, 581-592.
- 2
- 3 Rubin, D. B. (1987). *Multiple imputation for nonresponse in surveys*. New York:  
4 Wiley.
- 5
- 6 Saggino, A., Kline, P. (1995). Item factor analysis of the Italian version of the Myers-  
7 Briggs Type Indicator. *Personality and Individual Differences*, 19, 243-249.
- 8
- 9 Schafer, J. L. (1997). *Analysis of incomplete multivariate data*. London: Chapman &  
10 Hall.
- 11
- 12 Schafer, J. L. (1998). NORM: Version 2.02 for Windows 95/98/NT. Retrieved,  
13 September 2, 2008, from <http://www.stat.psu.edu/~jls/misoftwa.html>.
- 14
- 15 Segal, D. L., Coolidge, F. L., Cahill, B. S., O'Riley, A. A. (2008). Psychometric  
16 properties of the Beck Depression Inventory—II (BDI-II) among community-  
17 dwelling older adults. *Behavior Modification*, 32, 3-20.
- 18
- 19 Sheviin, M., & Adamson, G. (2005). Alternative factor models and factorial  
20 invariance of the GHQ-12: A large sample analysis using confirmatory factor  
21 analysis. *Psychological Assessment*, 17, 231-236.
- 22
- 23 Sijtsma, K., & Van der Ark, L. A. (2003). Investigation and treatment of missing item  
24 scores in test and questionnaire data. *Multivariate Behavioral Research*, 38,  
25 505-528.
- 26
- 27 S-Plus 8 for Windows [Computer software]. (2007). Seattle, WA: Insightful  
28 Corporation.
- 29
- 30 SPSS Inc. (2008). SPSS 16.0 for Windows [Computer software]. Chicago: author.
- 31
- 32 StataCorp. (2007). Stata Statistical Software: Release 10 [Computer software].  
33 College Station, TX: StataCorp LP.
- 34

- 1 Van Buuren, S. (2009). Item imputation without specifying scale structure.  
2 Methodology, x, xx-xx (this issue).  
3
- 4 Van der Ark, L. A. (2007). Mokken scale analysis in R. *Journal of Statistical*  
5 *Software*, 20, 1-19.  
6
- 7 Van Ginkel, J. R. (2006). MI.sps and MI-mul.sps [Computer code]. Retrieved  
8 September 3, 2008, from <http://www.uvt.nl/mto/software2.html>  
9
- 10 Van Ginkel, J. R., & Van der Ark, L. A. (2005a). SPSS syntax for missing value  
11 imputation in test and questionnaire data. *Applied Psychological*  
12 *Measurement*, 29, 152-153.  
13
- 14 Van Ginkel, J. R. & Van der Ark, L. A. (2005b). TW.SPS and RUNTW.SPS.  
15 [Computer code]. Retrieved September 3, 2008, from  
16 <http://www.uvt.nl/mto/software2.html>  
17
- 18 Van Ginkel, J. R., Van der Ark, L. A., & Sijtsma, K. (2007a). Multiple imputation of  
19 test and questionnaire data and influence on psychometric results. *Multivariate*  
20 *Behavioral Research*, 42, 387-414.  
21
- 22 Van Ginkel, J. R., Van der Ark, L. A., & Sijtsma, K. (2007b). Multiple imputation for  
23 item scores when test data are factorially complex. *British Journal of*  
24 *Mathematical and Statistical Psychology*, 60, 315-337.  
25
- 26 Van Ginkel, J. R., Van der Ark, L. A., Sijtsma, K., & Vermunt, J. K. (2007). Two-  
27 way imputation: A Bayesian method for estimating missing scores in tests and  
28 questionnaires, and an accurate approximation. *Computational Statistics &*  
29 *Data Analysis*, 51, 4013-4027.  
30
- 31 Vermunt, J. K., Van Ginkel, J. R., Van der Ark, L. A., & Sijtsma, K. (in press).  
32 Multiple imputation of incomplete categorical data using latent class analysis.  
33 *Sociological Methodology*.  
34

1 Vorst, H.C.M. (1992). [Responses to the Adjective Checklist] Unpublished raw data.  
2  
3 Watson, D., O'Hara, M. W., Simms, L. J., Kotov, R., Chmielewski, M., McDade-  
4 Montez, E. A., Gamez, W., Stuart, S. (2007). Development and Validation of  
5 the Inventory of Depression and Anxiety Symptoms (IDAS). *Psychological*  
6 *Assessment*, 19, 253-268.  
7  
8 Yuan, Y. C. (2000). *Multiple imputation for missing data: Concepts and new*  
9 *development*. Proceedings of the Twenty-Fifth Annual SAS Users Group  
10 International Conference (Paper, No. 267). Cary, NC: SAS Institute. Retrieved  
11 September 3, 2007, from  
12 <http://www.ats.ucla.edu/stat/sas/library/multipleimputation.pdf>

- 1 Table 1. *Overview of the 22 Subscales in the Adjective Checklist Data (Vorst, 1992)*
- 2 *and Corresponding Item Numbers.*

Scale	Item No.	Scale	Item No.
Communality	1-10	Change	111-119
Achievement	11-20	Succorance	120-129
Dominance	21-30	Abasement	130-139
Endurance	31-40	Deference	140-149
Order	41-50	Personal Adjustment	151-159
Intelligence	51-60	Ideal Self	160-169
Nurturance	61-70	Critical Parent	170-179
Affiliation	71-80	Nurturant Parent	180-189
Exhibition	81-90	Adult	190-199
Autonomy	91-100	Free Child	200-209
Agression	101-110	Adapted Child	210-218

1 Table 2. *Listwise Deletion Results of Statistical Analyses of the ACL data (Vorst,*  
2 *1992) (First Row), and With 5% of the Item Scores Removed According to Either*  
3 *MCAR (Second Row), MAR (Third Row), or NMAR (Fourth Row).*

Data	Alpha	Mean test score	Min. test score	Max. test score	<b>Mean diff.</b>	<i>t</i>	df	<i>p</i>
original	0.807	24.3764	5	40	<b>-1.298</b>	-2.261	431	.024
MCAR	0.802	24.5271	5	40	<b>-0.740</b>	-0.994	256	.321
MAR	0.810	24.2943	5	38	<b>-1.180</b>	-0.768	263	.114
NMAR	0.818	23.4841	5	38	<b>-0.972</b>	-1.254	250	.211



1 Table 3. *Frequency of Occurrence of Missing Data in 24 Issues of Psychological*  
 2 *Assessment, Personality and Individual Differences, and Journal of Personality*  
 3 *Assessment.*

Journal	Vol.	Type of Nonresponse						Total
		UN	AT	IN	PL	Not clear	None reported	
Psychol. Assessment	1995	2	5	14	1	1	21	44
	1997	8	7	17	3	2	25	62
	2000	3	4	17	1	1	12	38
	2002	12	3	18	0	0	9	42
	2005	1	8	13	0	1	18	41
	2007	11	8	22	0	1	9	51
	Total		37	35	101	5	6	94
Pers. Individ. Differ.	1995	3	4	14	0	0	41	62
	1997	9	3	15	0	1	45	73
	2000	4	1	13	0	2	41	61
	2002	10	6	16	0	1	27	60
	2005	7	3	17	0	3	52	82
	2007	10	2	26	0	1	51	90
	Total		43	19	101	0	8	257
J. Pers. Assess.	1995	5	3	19	0	1	25	53
	1997	0	1	8	0	2	27	38
	2000	6	2	7	0	2	14	31
	2002	5	1	14	0	0	15	35
	2005	4	5	13	1	1	11	35
	2007	2	6	11	0	0	10	29
	Total		22	18	72	1	6	102
Total	1995	10	12	47	1	2	87	159
	1997	17	11	40	3	5	97	173
	2000	13	7	37	1	5	67	130
	2002	27	10	48	0	1	51	137
	2005	12	16	43	1	5	81	158
	2007	23	16	59	0	2	70	170
	Total		102	72	274	6	20	453

4  
 5 *Note.* UN = unit nonresponse, AT = attrition, IN = item nonreponse, PL = planned  
 6 missingness.

1 Table 4. *Statistics of the Types of Nonresponse Encountered in 24 Issues of*  
 2 *Psychological Assessment, Personality and Individual Differences, and Journal of*  
 3 *Personality Assessment. For the Studies That Reported Missing Values the Mean (M),*  
 4 *Standard Deviation (SD), Minimum, and Maximum Number of Incomplete Response*  
 5 *Patterns Are Reported.*

Type of nonresponse	<i>N</i>	<i>M</i>	<i>SD</i>	Skewness	Minimum	Maximum
UN	99	0.302	0.219	0.599	0.005	0.856
AT	74	0.186	0.136	1.090	0.016	0.703
IN	186	0.092	0.110	1.970	0.001	0.650
Not clear	10	0.385	0.315	0.326	0.040	0.898

6  
 7 *Note.* *N* = Number of cases where the type of nonresponse was reported. UN = unit  
 8 nonresponse, AT = attrition, IN = item nonreponse.

- 1 Table 5. *Example of a Data Set With Incomplete Item Scores (Sijtsma & Van der Ark,*
- 2 *2003).*

Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
1	2	1	1	.	.
2	3	5	4	5	5
3	4	3	.	3	4
4	1	1	1	3	2
5	.	3	3	.	4
6	5	5	3	.	5
7	1	3	2	2	2
8	3	3	1	2	.

- 1 Table 6. *Example of Deterministic and Stochastic Variable Mean Imputation (Left),*
- 2 *and Deterministic and Stochastic Regression Imputation (Right), in the Data Example*
- 3 *from Sijtsma and Van der Ark (2003).*

Deterministic Variable Mean Imputation						Deterministic Regression Imputation					
Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
1	2	1	1	<b>3</b>	<b>3.67</b>	1	2	1	1	<b>3</b>	<b>2.47</b>
2	3	5	4	5	5	2	3	5	4	5	5
3	4	3	<b>2.14</b>	3	4	3	4	3	<b>2.42</b>	3	4
4	1	1	1	3	2	4	1	1	1	3	2
5	<b>2.71</b>	3	3	<b>3</b>	4	5	<b>2.71</b>	3	3	<b>3.28</b>	4
6	5	5	3	<b>3</b>	5	6	5	5	3	<b>4.13</b>	5
7	1	3	2	2	2	7	1	3	2	2	2
8	3	3	1	2	<b>3.67</b>	8	3	3	1	2	<b>2.61</b>
<i>Mean</i>	2.71	3	2.14	3	3.67						

  

Stochastic Variable Mean Imputation						Stochastic Regression Imputation					
Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
1	2	1	1	<b>0.72</b>	<b>1.38</b>	1	2	1	1	<b>2.97</b>	<b>2.93</b>
2	3	5	4	5	5	2	3	5	4	5	5
3	4	3	<b>2.28</b>	3	4	3	4	3	<b>3.28</b>	3	4
4	1	1	1	3	2	4	1	1	1	3	2
5	<b>4.52</b>	3	3	<b>2.71</b>	4	5	<b>0.83</b>	3	3	<b>2.62</b>	4
6	5	5	3	<b>0.71</b>	5	6	5	5	3	<b>3.93</b>	5
7	1	3	2	2	2	7	1	3	2	2	2
8	3	3	1	2	<b>3.59</b>	8	3	3	1	2	<b>2.55</b>
<i>Mean</i>	2.71	3	2.14	3	3.67						
<i>SD</i>	1.50	1.51	1.21	1.22	1.37						

- 1 Table 7. *Example of Multiple Imputation Using NORM (Schafer, 1998) in the Data*
- 2 *Example From Sijtsma and Van der Ark (2003).*

Imputed Data Set #1						Imputed Data Set #2					
Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
1	2	1	1	<b>3.72</b>	<b>1.93</b>	1	2	1	1	<b>3.29</b>	<b>3.73</b>
2	3	5	4	5	5	2	3	5	4	5	5
3	4	3	<b>2.18</b>	3	4	3	4	3	<b>2.47</b>	3	4
4	1	1	1	3	2	4	1	1	1	3	2
5	<b>5.81</b>	3	3	<b>4.17</b>	4	5	<b>-0.03</b>	3	3	<b>2.86</b>	4
6	5	5	3	<b>5.89</b>	5	6	5	5	3	<b>3.59</b>	5
7	1	3	2	2	2	7	1	3	2	2	2
8	3	3	1	2	<b>1.69</b>	8	3	3	1	2	<b>1.99</b>

Imputed Data Set #3						Imputed Data Set #4					
Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
1	2	1	1	<b>4.82</b>	<b>3.17</b>	1	2	1	1	<b>2.01</b>	<b>2.17</b>
2	3	5	4	5	5	2	3	5	4	5	5
3	4	3	<b>-0.18</b>	3	4	3	4	3	<b>1.87</b>	3	4
4	1	1	1	3	2	4	1	1	1	3	2
5	<b>1.97</b>	3	3	<b>5.74</b>	4	5	<b>2.40</b>	3	3	<b>5.08</b>	4
6	5	5	3	<b>4.43</b>	5	6	5	5	3	<b>3.60</b>	5
7	1	3	2	2	2	7	1	3	2	2	2
8	3	3	1	2	<b>2.52</b>	8	3	3	1	2	<b>4.29</b>

Imputed Data Set #5					
Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
1	2	1	1	<b>0.72</b>	<b>1.38</b>
2	3	5	4	5	5
3	4	3	<b>2.28</b>	3	4
4	1	1	1	3	2
5	<b>4.52</b>	3	3	<b>2.71</b>	4
6	5	5	3	<b>0.71</b>	5
7	1	3	2	2	2
8	3	3	1	2	<b>3.59</b>

1 Table 8. *Frequencies in Which Missing-Data Methods were Used in Studies from 24*  
 2 *Issues of Psychological Assessment, Personality and Individual Differences, and*  
 3 *Journal of Personality Assessment.*

Missing-data method	Type of Nonresponse						Total
	UN	AT	IN	PL	Not clear	None reported	
LD	91	44	164	1	14	0	314
LD-CM	10	8	13	0	0	0	31
LD-CM-R	6	11	8	0	1	0	26
AC	1	11	64	1	4	0	81
AC-CM	0	4	1	0	0	0	5
IMP	0	0	18	1	0	0	19
DMLE	0	1	9	1	1	0	12
VD	1	0	36	0	0	0	37
FU	0	1	2	0	0	0	3
PRO	0	0	10	0	0	0	10
FU-LD-CM	1	1	0	0	0	0	2
Other	0	1	0	0	0	0	1
None reported	0	1	26	2	2	453	484
Total	110	83	351	6	22	453	1025

4  
 5 *Note.* UN = unit nonresponse, AT = attrition, IN = item nonreponse, PL = planned  
 6 missingness. LD = listwise deletion. LD-CM = Listwise deletion with check for  
 7 MCAR, MCAR not rejected. LD-CM-R = Listwise deletion with check for MCAR,  
 8 MCAR rejected. AC = Available case analysis. AC-CM = Available-case analysis  
 9 with check for MCAR, MCAR not rejected. IMP = Imputation. DMLE = Direct  
 10 maximum likelihood estimation. VD = Variable deletion. FU = Follow-up PRO = Pro  
 11 rating. FU-LD-CM = Combination of Follow-up and Listwise deletion with check for  
 12 MCAR.

- 1 Table 9. *Example of Deterministic TW Imputation in the Data Example from Sijtsma*
- 2 *and Van der Ark (2003).*

Case	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$PM_i$
1	2	1	1	<b>1.45</b>	<b>2.12</b>	1.33
2	3	5	4	5	5	4.4
3	4	3	<b>2.76</b>	3	4	3.5
4	1	1	1	3	2	1.6
5	<b>3.17</b>	3	3	<b>3.45</b>	4	3.33
6	5	5	3	<b>4.62</b>	5	4.5
7	1	3	2	2	2	2
8	3	3	1	2	<b>3.04</b>	2.25
$IM_j$	2.71	3	2.14	3	3.67	$OM = 2.88$

1 Table 10. *Results of Statistical Analyses of the ACL data (Vorst, 1992) Without*  
 2 *Missing Data (First Row), and With 5% of the Item Scores Removed According to*  
 3 *Either MCAR (Second Row), MAR (Third Row), or NMAR (Fourth Row). Missing*  
 4 *Data are Imputed Using Method TW.*

Data	alpha	Mean test score	Min. test score	Max. test score	Mean diff.	<i>t</i>	df	<i>p</i>
original	0.807	24.3764	5.00	40.00	-1.298	-2.261	431	.024
MCAR	0.811	24.3982	5.00	40.00	-1.196	-2.039	391	.042
MAR	0.810	24.3473	5.00	39.80	-1.307	-2.136	410	.033
NMAR	0.810	24.1621	5.00	40.00	-1.328	-2.264	402	.024



```

1  GET FILE='C:\imputation\example.sav'.
2
3  DATASET DECLARE deterministic.
4  MVA
5  VARIABLES = X1 X2 X3 X4 X5
6  /EM ( TOLERANCE=0.001 CONVERGENCE=0.0001 ITERATIONS=25 )
7  /REGRESSION ( TOLERANCE=0.001 FLIMIT=4.0 ADDTYPE= NONE OUTFILE=stochastic )
8  .
9
10 GET FILE='C:\imputation\example.sav'.
11
12 SET SEED = 2 .
13 DATASET DECLARE stochastic.
14 MVA
15 VARIABLES = X1 X2 X3 X4 X5
16 /EM ( TOLERANCE=0.001 CONVERGENCE=0.0001 ITERATIONS=25 )
17 /REGRESSION ( TOLERANCE=0.001 FLIMIT=4.0 ADDTYPE=RESIDUAL
18 OUTFILE=stochastic ) .
19 DESCRIPTIVES
20 VARIABLES=X1 X2 X3 X4 X5
21 /STATISTICS=MEAN STDDEV .
22

```

23 *Figure 1: SPSS Syntax for Applying Both Deterministic and Stochastic Regression*  
24 *Imputation in the Example Data Set from Sijtsma and Van der Ark (2003).*