



## A SURVEY ON MISSING DATA AND METHODS TO FIND THE MISSING VALUES

<sup>1</sup>P.Logeshwari, <sup>2</sup>Dr. Antony Selvadoss Thanamani

<sup>1</sup>PhD Research Scholar, <sup>2</sup>Associate Professor,

<sup>1,2</sup>Department of Computer Science,

<sup>1,2</sup>Nallamuthu Gounder Mahalingam College, Pollachi, India.

### ABSTRACT

Missing data plagues almost all surveys, and quite a number of designed experiments. No matter how carefully an investigator tries to have all questions fully responded to in a survey, or how well designed an experiment is; examples of how this can occur are when a question is unanswered in a survey, or a flood has removed a crop planted close to a river. The problem is, how to deal with missing data, once it has been deemed impossible to recover the actual missing values. Traditional approaches include case deletion and mean imputation. These are the default for the major Statistical packages. In the last decade interest has centered on Regression Imputation, and Imputation of values using the Expectation Maximization algorithm, both of which will perform Single Imputation. More recently Multiple Imputation has become available, and is now being included as an option in the mainstream packages.

**Keywords:** - [Missing Data, Imputation techniques, missing methods, Data Identification]

### 1. SURVEY ON MISSING DATA

Historical Development Until the 1970s, missing values were handled primarily by editing. Rubin (1976)

developed a frame-work of inference from incomplete data that remains in use today. The formulation of the EM algorithm (Dempster, Laird, & Rubin, 1977) made it feasible to compute ML estimates in many missing-data problems. Rather than deleting or filling in incomplete cases, ML treats the missing data as random variables to be removed from (i.e., integrated out of) the likelihood function as if they were never sampled. We elaborate on this point later after introducing the notion of MAR. Many examples of EM were described by Little and Rubin (1987). Their book also documented the shortcomings of case deletion and single imputation, arguing for explicit models over informal procedures. About the same time, Rubin (1987) introduced the idea of MI, in which each missing value is replaced with  $m > 1$  simulated values prior to analysis. Creation of MIs was facilitated by computer technology and new methods for Bayesian simulation discovered in the late 1980s (Schafer, 1997). ML and MI are now becoming standard because of implementations in free and commercial software

### 2. MECHANISMS OF MISSING MCAR

The term 'Missing Completely at Random' refers to data where the missingness mechanism does not depend on

the variable of interest, or any other variable, which is observed in the dataset. MCAR is both missing at random, and observed at random (This means the data was collected randomly, and does not depend on any other variable in the data set). This very stringent condition is required in order for case deletion to be valid, and missing data is very rarely MCAR (Rubin, 1976).

### MAR

The term 'Missing at Random' is a misnomer, as the missing data is anything but missing at random. The intuitive meaning of this term is better suited to the term MCAR. What MAR means is missing, but conditional on some other 'X-variable' observed in the data set, although not on the 'Y-variable' of interest (Schafer, 1997).

### NMAR

Not Missing at Random, (or informatively missing, as it is often known) occurs when the Missingness mechanism depends on the actual value of the missing data. This is the most difficult condition to model. Traditional Procedures to find Missing Data Compare the missing and non-missing cases on variables where information is not missing.

Whatever strategy you follow you may be able to add plausibility to your results (or detect potential biases) by comparing sample members on variables that are not missing. For example, in a panel study, some respondents will not be re-interviewed because they could not be found or else refused to participate. You can compare respondents and non-respondents in terms of demographic characteristics such as race, age, income, etc.

If there are noteworthy differences, you can point them out, e.g. lower-income individuals appear to be underrepresented in the sample similarly, and you can compare individuals who answered a question with those who failed to answer. Alternatively, sometimes you may have external

information you can draw on, e.g. you know what percentage of the population is female or black, and you can compare your sample's characteristics with the known population characteristics.

### Dropping variables

When, for one or a few variables, a substantial proportion of cases lack data, the analyst may simply opt to drop the variables. This is no great loss if the variables had little effect on Y anyway. However, you presumably would not have asked the question if you did not think it was important. Still, this is often the best or at least most practical approach. A great deal of missing data for an item might indicate that a question was poorly worded, or perhaps there were problems with collecting the data.

**Dropping subjects**, i.e. list wise (also called case wise) deletion of missing data. Particularly if the missing data is limited to a small number of the subjects, you may just opt to eliminate those cases from the analysis. That is, if a subject is missing data on any of the variables used in the analysis, it is dropped completely. The remaining cases, however, may not be representative of the population. Even if data is missing on a random basis, a list wise deletion of cases could result in a substantial reduction in sample size, if many cases were missing data on at least one variable.

## 3. METHODS TO IDENTIFY MISSING DATA

**(a). Case Deletion:** This can be either list wise (complete case only) or all value (Pairwise-available case), the cases are deleted which contain missing data, for the analysis being carried out.

**(b). Single Imputation:** This can include group means, medians or modes (depending on the data), Regression Imputation, Stochastic Regression Imputation (deterministic regression imputation with an

added random error component), or EM Imputation (this uses the Expectation-Maximization algorithm to predict the missing value), or hot deck imputation, or last value carried forward for longitudinal data, and a variety of other methods (Scheffer, 2000). End users Very often demand a single complete data set.

**(c). Multiple Imputations:** Frequentist MI. This returns  $m$  complete datasheets by imputing  $m$  times. This can be based on propensity scoring, if imputation model fails to converge. Bayesian MI uses MCMC algorithm with a non-informative prior to predict the posterior distribution from which random draws are made, producing  $m$  individual datasheets. Successful multiple imputation may be shunned by an end-user, as the concept of more than one datasheet for a particular survey is daunting to non-statisticians. However, multiple imputation is always better than case deletion, or single ad-hoc methods

**(d) Mean imputation within classes (MC).** This method divides the total sample into imputation classes according to values on the auxiliary variables. The classes may be defined as all the cells in the cross-tabulation of the (categorized) auxiliary variables, but this symmetry is not essential; instead, some auxiliary variables may be used for one part of the sample while others are used for another part, or groups of cells may be combined. If all the cells in the cross-tabulation are used, the linear function can be expressed as a model with the main effects and all levels of interaction for the auxiliary variables. In general, the model can be represented by  $Y_{mi} = \beta_0 + \sum_j \beta_j z_{ji}$ , where the  $z_{ji}$  are dummy variables,  $z_{ji} = 1$  if the  $i$ -th non respondent is in class  $j$ ,  $z_{ji} = 0$  otherwise ( $j = 1, 2, \dots, (H-1)$ ). Since  $e_{mi} = 0$ , the method is a deterministic one.

**(e) Random imputation within classes (RC).** This method corresponds to the random overall method except that it is applied within imputation classes. Each non respondent is assigned the  $y$  value of a

respondent randomly selected from the same imputation class. The method is the stochastic equivalent of the mean within class method, respondent residual selected at random within imputation class  $j$  in which non respondent is located.

**(f) Hot-deck imputation.** The term hot-deck imputation has a variety of meanings, but refers here to the sequential type of procedure used by the Bureau of the Census with the labor force items in the Current Population Survey (CPS) (Brooks and Bailar, 1978). This is sometimes known as the traditional hot-deck procedure. The procedure begins with the specification of imputation classes, and for each class the assignment of a single value for the  $y$ -variable to provide a starting point for the process. These starting values may, for instance, be obtained by taking a respondent value for each class or a representative value such as the class mean from a previous round of the survey. The records of the current survey are then treated sequentially. If a record has a response for the  $y$ -variable; that value replaces the value previously stored for its imputation class. If the record has a missing response, it is assigned the value currently stored for its imputation class. A major attraction of this procedure is its computing economy; since all imputations are made from a single pass through the data file. The hot-deck method is similar to the random within class method in which donors are selected by unrestricted sampling (i.e. SRS with replacement). If the order of the records in the data file were random, the two methods would be equivalent, apart from the start-up process. The sequential hot-deck procedure generally benefits from the non-random order of the data file, since use of the preceding donor in the imputation class yields an additional degree of matching which is advantageous if the file order creates positive autocorrelation. This benefit is unlikely to be substantial, however, when the imputation classes are small and spread throughout the file - as is often the case. A disadvantage of the hot-

deck method is that it may easily give rise to multiple use of donors, a feature which leads to a loss of precision for the survey estimators. This occurs when within a given imputation class a record with a missing response is followed by one or more records with missing responses; all these records are then assigned the value from the last respondent in the class  $s$ . The random within class method with unrestricted sampling of donors shares this disadvantage. With the random within class method, however, the multiple use of donors may be minimized by sampling donors without replacement.

It is impossible to develop a model-free theoretical evaluation for the hot-deck method because of its dependence on the order of the file and its lack of a probability mechanism. For this reason, it will not be examined in the subsequent sections; the results for the random within class method with unrestricted sampling should, however, provide a reasonable guide to its performance.

Useful discussions of the hot-deck procedure are provided by Bailer, Bailey and Corby (1978), Bailer and Bailer (1978, 1979), Ford (1980), Oh and Scheuren (1980), Oh, Scheuren and Nisselson (1980) and I. Sande (1979a,b).

**(g) Flexible matching imputation.** The term flexible matching imputation is used here for the modified hot-deck procedure that has been used since 1976 for the CPS March Income Supplement. The procedure sorts respondents and nonrespondents into a large number of imputation classes, constructed from a detailed categorization of a sizeable set of auxiliary variables. Nonrespondents are then matched with respondents on a hierarchical basis, in the sense that if a nonrespondent cannot be matched with a respondent in the initial imputation class, classes are collapsed and the match is made at a lower level. Three levels are used with the March Income Supplement, the lowest level being such that a match can always be made.

The procedure enables closer matches to be secured for many nonrespondents than does the traditional hot-deck procedure. It also avoids the multiple use of respondents in classes where the number of nonrespondents does not exceed the number of respondents. Further details on the implementation and evaluation of the procedure are given by Coder (1978) and Welniak and Coder (1980).

**(h) Predicted regression imputation (PR).** This method uses respondent data to regress  $y$  on the auxiliary variables. Missing  $y$ -values are then imputed as the predicted values from the regression equation,  $Y_{mi} = b_0 + \sum_j b_j x_{ji}$ . This is a deterministic method with  $e_{mi} = 0$ . The auxiliary variables may be quantitative or qualitative, the latter being incorporated by Means of dummy variables. If the  $y$ -variable is qualitative, log-linear or logistic models may be used. As in any regression analysis, specific interaction terms may be included in the regression equation, and transformations of the variables may be useful.

A special case of the regression model is the ratio model  $Y_{mi} = b_0 x_{zi}$  with a single auxiliary variable and an intercept of zero (Ford, Kleweno and Tortora, 1980). This model may be used in panel surveys with  $z$  representing the same variable as  $y$  measured on the previous wave.

**(i) Random regression imputation (RR)** .This method is the stochastic version of the predicted regression method: the imputed values are the predicted values from the regression equation plus residual terms  $e_{mi}$ . Depending on the assumptions made, the residuals can be determined in various ways, including.

(i) If the residuals are assumed to be homoscedastic and normally distributed, a residual can be chosen at random from a normal distribution with zero mean and variance equal to the residual variance from the regression.

(ii) If the residuals are assumed to come from the same, unspecified distribution, they can

be chosen random from the respondents' residuals.

(iii) As a protection against non-linearity and non-additivity in the regression model, the residuals may be taken from respondents with similar values on the auxiliary variables. If the donor respondent has the identical set of  $z$  values as the nonrespondent, the procedure reduces to assigning the respondent's  $y$ -value to the nonrespondent. This point demonstrates the closer relationship between this procedure and the random within class method. Applications of regression and categorical data models for imputation are described by Schieber (1978), Herzog and Lancaster (1980) and Herzog (1980).

**(j) Distance function matching.** This method assigns the  $y$ -value of the nearest respondent to each nonrespondent, with "nearest" defined by a distance function of the auxiliary variables. The method is primarily concerned with quantitative variables; however, qualitative variables may be included either by using the distance function approach within imputation classes formed by qualitative auxiliary variables or by incorporating these variables into the distance function. With a single auxiliary variable, the sample may be ordered by the variable, and the nearest respondent (donor) to each nonrespondent is taken where "nearest" may be defined as the minimum absolute difference between the nonrespondent's and donor's values in the auxiliary variable or in some transformation of the auxiliary variable. When several auxiliary variables are used, the issue of transformations becomes more critical; one approach is to transform all auxiliary variables to their ranks.

It can be constructed to reduce the multiple use of donors. For instance, distance may be defined as  $D(I + pd)$  where  $D$  is the basic distance,  $d$  is the number of times the donor has already been used and  $p$  is a penalty for each usage (Colledge et al., 1978).

A variant of this method assigns to the nonrespondent the average value of neighboring respondents, for instance the average value of the two adjacent respondents (Ford, 1976). As with other averaging procedures, this procedure suffers the disadvantage of distorting distributions.

**k) Deductive imputation.** This imputation method depends on some redundancy in the data so that a missing response can be deduced from the auxiliary information, i.e.  $y_{mi} = f(z_i)$  exactly. For example, if a record should contain a series of amounts and their total but one of the amounts is missing, the missing value can be deduced by subtraction. The method can be extended to situations where the deduced value is highly likely to be the correct value or at least close to it; for instance, in a panel survey with a variable that remains almost constant over time, a missing response on one wave of the panel may be assigned the record's value for the item on the preceding or succeeding wave.

**(l) Mean imputation overall (MO).** This method assigns the overall respondent mean,  $Y_r$ , to all missing responses. It is the deterministic degenerate form of the linear function with no auxiliary variables **(m) Random imputation overall (RO).** This method assigns each nonrespondent the  $y$ -value of a respondent selected at random from the total respondent sample. The method is the stochastic degenerate form of the linear function with no auxiliary variables,  $Y_{mi} = Y_r + e_{mi}$ , with  $e_{mi} = Y_{rk} - Y_r$ , which reduces to  $Y_{mi} = y_{rk}$ . Given an  $e$  sample initially, the subsample of respondents to act as donors can be selected by any  $e$  sampling scheme (e. g. unrestricted sampling, SRS, proportionate stratified sampling, or systematic sampling).

## CONCLUSIONS

A major attraction of imputation is that it generates a complete data set that may be readily used for many different forms of analysis. As the preceding sections have

shown, however, caution is needed in analyzing a data set that includes imputed values. In the case of univariate analyses, deterministic imputation methods serve well for estimating means and totals, but they distort the distributional properties of the variable; stochastic methods are less efficient for estimating means and totals but they preserve the variability in the respondent data. All methods are likely to attenuate the covariance between the variable subject to imputation and other variables, except for those other variables that are used as auxiliary variables in the imputation scheme. In consequence, when a dataset contains imputed values, special care is needed in studying the interrelationships between variables, whether the interrelationships are examined in terms of cross-tabulations, regression analyses or other forms of multivariate analysis. Alternative ways of handling missing survey data include dropping cases with missing values on the relevant variables from the analysis, direct estimation of the population parameters from a modeling approach, and weighting adjustments. Dropping cases with missing values is a widely used procedure, sometimes adopted on the ground that it avoids assumptions required in procedures which attempt to compensate for missing data.

## REFERENCES

- [1] Little, R.J.A. (1988). A test of missing completely at random for multivariate data with missing values. *Journal of the American Statistical Association*, 83, 1198-1202.
- [2] Little, R.J.A. (1988). Missing data in large surveys. *Journal of Business and Economic Statistics*, 6, 287-301.
- [3] Rubin, D.B. (1976). Inference and missing data. *Biometrika*, 63, 581-592.
- [4] Rubin, D.B. and Schenker, N. (1991). Multiple imputation in health-care data bases: An overview and some applications. *Statistics in Medicine*, 10, 585-598.
- [5] Schafer, J.L. (2000). NORM. Version 2.03. <http://www.stat.psu.edu/~jls/>.
- [6] Schafer, J.L. and Graham, J.W. (2002). Missing data: Our view of the state of the art. *Psychological Methods*, 7, 147-177.
- [7] van Buuren, S., Brands, J.P.L., Groothuis-Oudshoorn, C.G.M., and Rubin, D.B. (2006). Fully conditional specification in multivariate imputation. *Journal of Statistical Computation and Simulation*, 76, 1049-1064.
- [8] van Buuren, S. (2007). Multiple imputation of discrete and continuous data by fully conditional specification. *Statistical Methods in Medical Research*, 16, 219-242.
- [9] van Ginkel, J.R., van der Ark, L.A., and Sijtsma, K. (2007). Multiple imputation of test and questionnaire data and influence on psychometric results. *Multivariate Behavioral Research*, 42, 387-414.