

Dealing with Missing Data: A comparative exploration of approaches utilizing the Integrated City Sustainability Database

Cali Curley, Rachel M. Krause, Richard Feiock, Christopher V. Hawkins

Abstract:

Studies of governments and local organizations using survey data have played a critical role in the development of urban studies and related disciplines. However, missing data pose a daunting challenge for this research. This article seeks to raise awareness about the treatment of missing data in urban studies research by comparing and evaluating three commonly used approaches to deal with missing data – listwise deletion, single imputation, and multiple imputation. Comparative analyses illustrate the relative performance of these approaches using the second generation Integrated City Sustainability Database (ICSD). The results demonstrate the benefit of using an approach to missing data based on multiple imputation, using a theoretically informed and statistically supported set of predictor variables to develop a more complete sample that is free of issues raised by non-response in survey data. The results confirm the usefulness of the ICSD in the study of environmental and sustainability and other policy in U.S. cities. We conclude with a discussion of results and provide a set of recommendations for urban researcher scholars.

Acknowledgement: This material is based upon work supported by the National Science Foundation under Grant Nos. 1461526/1461506/1461460. Any opinions, findings, and conclusions expressed are those of the authors and do not necessarily reflect the views of the National Science Foundation. The University of Kansas Center for Research Methods and Data Analysis provided valuable assistance in the imputation process, as did our graduate research assistants XXX.

Introduction

This article seeks to raise awareness about how to treat missing data in urban studies research. A large proportion of the empirical research on urban politics and policy relies on data collected through surveys of local government or community organization leaders. Surveys provide a relatively efficient way to collect large amounts of consistently measured individual or organizational information needed to conduct comprehensive and accurate statistical analysis. This is particularly important if the aim of research is to produce generalizable findings and contribute to understanding a particular phenomenon by testing theory. However, missing data is a common and significant challenge in survey-based research. It often influences the selection of a statistical method of analysis, and, depending on its severity, can undermine the confidence of analysis. Nonetheless, the problems associated with missing data are among the least acknowledged issues when conducting and reporting analysis.

Missing survey data occurs for three reasons: 1) non-coverage - the observation fell outside of the sample, 2) total nonresponse – the would-be respondent failed to respond to the survey, and 3) item non-response - the respondent skipped a particular survey item (Brick and Kalton, 1996). Although data missing as a result of these different causes presents distinct challenges for the researcher, listwise deletion, the default operation in most statistical software packages, is a common applied remedy for all three. This approach removes any observation from the analysis that has incomplete information, i.e. is missing a value for any variable included in the model for any reason. Peng et al (2006) examined 1,087 published studies in education and psychology, of which 48% contained missing data. Within that subset, they found

that authors used listwise deletion 97% of the time.

This paper demonstrates the impact that different remedies for missing data may have on research findings and offers a rationale for its appropriate treatment. We specifically discuss the classifications of missing data, the specific problems associated with each, and the common approaches that have been developed to address them. This is followed by an illustration of the treatment of missing data using three techniques – listwise deletion, single imputation, and multiple imputation – applied to data from the second generation Integrated City Sustainability Database (ICSD) and comparison of their relative performance in analysis. We use the results of the analysis as the basis for a concluding discussion of the missing data techniques and provide a set of recommendations for researchers using survey data.

Overview of Missing Data

Three classifications of missing data are important to the following discussion: data Missing Completely at Random (MCAR), data Missing at Random (MAR), and data Missing Not at Random (MNAR). This taxonomy provides insight into which tool is appropriate for dealing with the missing data. Table 1 below provides a brief overview. For data that are MCAR the missing values are independent from values of observed or unobserved characteristics in the data set. Therefore, the missing value is not the result of a strategic choice on the part of the respondent nor a function of other captured or uncaptured variables. This means that the observed pattern of missingness is not related to any other data, whether present or missing. For example, MCAR data might result if a survey respondent unintentionally failed to answer a question that the researcher is using as a variable in the analysis. It is difficult to ascertain whether data are truly MCAR; in this situation, the researcher must ask if there is any theoretical reason that the respondent may have wanted to avoid answering that question. Little's (1988)

MCAR test can help inform the assessment as to whether data is truly MCAR or not. When encountering missing data a researcher can calculate a chi-square test to examine patterns of missingness for a number of specified variables (the “mcartest” command in Stata). The null assumption is that the data is MCAR, therefore the researcher hopes to fail to reject the null hypothesis by having a p-value larger than .05. An application of this test is included in the discussion of listwise deletion below. This test is one of several mechanisms that help determine whether the data associated with a particular variable is MCAR and should be utilized along with a logit model – in which the dependent variable takes a value of one if the variable of interest’s value is missing and zero if not – to examine if the values of other observed variables explain its missingness.

[Insert table 1 about here]

If both of these tests suggest that the missing data is MCAR then either listwise deletion or multiple imputations can be used without biasing estimates. Since listwise deletion will impact the power of the analysis, multiple imputations may still be the better approach. However, if the overall number of cases lost is small, listwise deletion is still an appropriate method (Myers, 2011; King et al, 2001). If however, one of the tests fail, the missing data would need to be treated as either “missing at random” (MAR) or “missing not at random” (MNAR). Data that are MAR are characterized by the fact that their presence or absence can be predicted using observed variables. A common example is when an individual intentionally skips the question asking about his/her income in a survey but provides the researcher with values of their employment status, education level, and years of experience at their current job. In this context, the value of the missing data is dependent on the value of observed responses and thus is characterized as being MAR.

On the other hand, there is no available explanation for data that is “missing not at random” (MNAR). When data is MNAR, the researcher cannot approximate the missing values because the values of other relevant variables are also not observed. Consider the previous example, if the observed data did not include employment status, education level, or experience, it would be challenging to determine an expected value of the respondents’ income. Moreover, a respondent’s income itself often determines whether or not (s)he provided a response. Therefore, if the researchers did not capture relevant explanatory variables, the missing data would be considered MNAR. Solutions that handle MAR data, such as multiple imputation relies on responses to other questions and relationships between missing and observable data to determine the value of the missingness. Despite this, multiple imputation and maximum likelihood are often unbiased with MNAR data (Schafer and Graham, 2002). Researchers may also learn about and possibly control for their MNAR data from working through Heckman Selection Models (Little 2016; DeMaris 2014).

It is important to consider the reason why data is missing when determining its treatment in statistical analysis. Since the different approaches – listwise deletion, single imputation, and multiple imputation – each make specific mathematical assumptions, misusing them may invalidate empirical results. Invalid assumptions and incorrect categorizations of missingness may 1) decrease the sample size, decreasing the power to estimate models, 2) increase the potential for biased results, and 3) over or under estimate standard errors. These impacts are important. If a large number of observations are lost, the resulting analysis will lose power and variables that would have otherwise been significant may no longer have enough variation to demonstrate their relationship to the dependent variable. If the subset of observations that were dropped due to missingness is systematically different from those that remain in, then both the

sample and any subsequent estimates generated from it will be biased. These bias related issues and loss in power, creates the potential for standard errors to be over or under estimated which means the model results are unreliable. Table 2 summarizes the advantages, disadvantages, concerns, and missingness assumptions of the different techniques explored in the next sections of this paper.

Approaches to handling missing data:

Scholars utilize a variety of alternative techniques in order to accommodate missing data and minimize its negative effects. Three of the most widely used approaches identified by Little (1988b) are: 1) examining the incomplete cases (Little 2016), 2) replacing values for missing data (Kong et al. 1994), and 3) providing statistical weights to complete cases (Little 2014; Brehm 1993). Within the general category of data replacement, there are specific techniques that vary in complexity. In addition to listwise deletion, two commonly used techniques include single imputation via mean replacement and multiple imputation. The paper proceeds through an examination of these techniques and compares their performance utilizing survey data in an application.

Listwise deletion, the default approach to handle missing data, is a convenient choice in most software packages. Two conditions must be met for listwise deletion to be appropriate for dealing with missing data: the missing data is MCAR and the sample remains large after the deletion. Deleting observations for non-response is less consequential if the values are MCAR, because if missingness is completely random the data deleted would also be random and it would thus not cause the loss of important variation. As previously described, a statistical approach, referred to as Little's test, can help indicate whether data is MCAR (Little 1988a). If the data is instead MAR or MNAR, it is inconsistent with the assumptions of listwise deletion and its use

may result in the sample mean being different from the population mean. It may also affect estimates in a manner similar to selection bias; if a set of respondents systematically choose not to answer a question and those observations are then deleted from the sample, the observations that remain in the analysis may be meaningfully different from the larger population.

The second issue with listwise deletion is that it reduces the sample size and thus the statistical power of the sample may be correspondingly reduced. Smaller samples are more likely to generate false null results that might otherwise not be null with a larger sample. Consider a hypothetical survey sent to a population of 700 respondents that obtained a 50% response rate ($n=350$). Of those respondents, 10% failed to answer a particular question contained in an analysis. If that missing data is MCAR then, by dropping those incomplete cases through listwise deletion, we are essentially taking a random sample of 90% of those respondents. Given the 10% missingness specified, we would only lose 35 cases and respectable sample size remains. Let's now suppose that we have 10% missing on four different variables included in our analysis. If the missingness is completely random then it is unlikely that the same cities skipped those four questions. Therefore, we could lose up to 40% of the total data or 140 responses, which raises concerns about the power of the sample size.

Single imputation is a general term that describes a family of missing data replacement techniques, including value replacement, mean replacement and single regression replacement.

Last value replacement, which can be used with panel or time series data, involves the replication of the most recent value in cases of missingness. Carrying the last known value forward yields a conservative estimate of the treatment effect when a post-test value is missing. For example, if a respondent was asked to rate their health on a scale of 1 to 10 and answered "8" the first time the survey was administered but failed to provide a response the second time it

was administered, the researcher would replace the missing value with 8. A second version of value replacement, sometimes referred to as “hot-decking,” uses information from similar observations to replace missing data. It is built around a premise similar to that of propensity score matching; if observations can be matched with others that look similar across the known values for a set of variables, missing ones can be replaced by the value of its match. This technique works if the data are MCAR or MAR and assumes that otherwise similar respondents are also alike in the category where data for one is missing.

Mean replacement, replaces missing observations with the mean value of that variable from observed responses in the sample. This preserves the overall mean of each variable but reduces the variation of the sample. By holding unobserved variables to the mean, it automatically sets the sum of squared differences for these observations to zero, which causes variance to be underestimated and it may not reflect the true relationship meaning that it is likely to reflect the true relationship between the dependent and independent variables. When the degree of missingness is small and the sample size is large, this technique may be appropriate. The smaller the amount of missingness, the less impact this has on the overall variance estimate. However, in smaller samples, the effect of mean replacement on these relationships will be larger.

An advanced version of single imputation is the single regression replacement method. This approach uses relevant observed variables (i.e. “informing variables”) to predict the value of the missing response via a regression analysis. This technique works well for data that is MAR, because, by definition, the other variables that can inform the missing value are observed in the data. The variable whose missing values are being estimated serves as the dependent variable in a regression and the independent or “informing” variables included in the model are

theoretically or statistically related to it. Once the coefficients of the informing variables are estimated, the missing values of the dependent variable can be calculated for each observation by substituting the associated values of the each informing variable back into the estimated equation. This estimation technique allows the value of missing data to vary by observation based on responses to the informing variables

Consider as an example, a scholar attempting to explain wages for a sample of respondents. However, her data contains several missing responses to a key variable associated with a survey question asking about professional competency. If she knows that age and education level are correlated with the observed values for professional competency, she can use those variables in a regression to develop a best guess of its value for each respondent who failed to provide it. The imputed values for competency can then be used along with all of the observed values for it and other variables in a model to predict wages. This helps illustrate that the point of imputation is not necessarily to pick the “right” value for the missing data, but rather to provide a value that allows all of the other data to be used without hampering the inference of the desired model (Rubin 1987, 1996).

In single regression replacement, the missing value is only measured once, which creates the potential for biasing the standard errors similar to mean replacement since there is no assessment of how likely it is that the imputed value is the true value nor any way to apply weighing based on such an assessment. If the inherent uncertainty in the prediction of the missing values is not accounted for, subsequent analysis may be influenced by the predicted missing values more than the true observed data, creating the potential for included bias and over or under estimated standard errors.

Multiple imputation is an extension of the single imputation regression replacement

method. As its name suggests, missing values are estimated multiple times. Analyzing multiply imputed data follows three steps: 1) the imputation of missing data, 2) the running of independent statistical analysis on the resulting individual data sets, and 3) the pooling of the results across the imputations.

The first step of multiple imputation is similar to that of single regression replacement method described above: variables that are theoretically related or statistically correlated to the target variable are identified and used in an appropriately specified regression model to predict the values of the missing data. However, in multiple imputation, this process is repeated numerous times in order to incorporate the uncertainty in the prediction process. Each missing value is estimated a number of different times and varies by inclusion of randomness. More specifically, the randomness represents a different value of the error term, incorporating the uncertainty in predicting the value of the missingness (Johnson and Young, 2011; White et al, 2010). Therefore, multiple imputation creates numerous data sets, each containing somewhat different estimates of the missing values. Rubin's (1978) formula suggests 3-10 imputations are necessary to produce results that incorporate enough variation in the prediction process; however, others argue the number of imputations should be similar to the percent of missing responses (Graham et al. 2007; Bodner 2008; Royston and White 2011). This ensures that the uncertainty inherent in the prediction of missing values is captured to appropriately the increase standard errors in the actual analysis of interest.

A second key difference between single regression replacement and multiple imputation is in how the data is analyzed as part of a theory-based model once missing values have been imputed. As described above, multiple imputation results in the creation of a number of different data sets. Theory-based models that use multiply imputed data must therefore be estimated

simultaneously with each set of data. Many statistical programs enable data to be specified as imputed, after which the simultaneous estimation is carried out automatically. For example, in STATA multiply imputed data must be specified with the command *miset*, which clearly defines where one data set begins and ends. The analysis is then run as usual, with the only addition in STATA the phrase *mi estimate:* prior to specifying the model.

These designations instruct the statistical software to, in the background, estimate the theory-based model across each of the imputed data sets. For example, if 20 rounds of imputation were used to generate values for the missing data, then 20 distinct data sets are created, and theory-based model is estimated 20 times. Once the analysis is executed, the results are pooled together and the pooled output is reported. This process may take more time than running a typical regression as it has to run that same analysis over 20 different data sets. The pooling process embeds all of the uncertainty from the imputation into the estimates of the standard errors that are presented in the output. The results can be interpreted normally, i.e. as they would be for non-imputed data. There are several different pooling rules, but the specified defaults in statistical packages are usually appropriate. A detailed overview of pooling rules¹ for normally and non-normally distributed parameters can be found in White et al (2011) and Allison (2002), respectively.

¹ For normally distributed parameters, the standard pooling process follows Rubin's Combination Rule, which incorporates the uncertainty generated by the process of imputation into the estimates of the standard errors. Rubin's Combination Rule incorporates the uncertainty or variation due to missing information and the results from just one data set. It does this by essentially averaging the variance over the imputed data sets and incorporating both within-imputation variance and between-imputation variance (White et al, 2011). Allison (2002) provides an overview of pooling methods for non-normally distributed parameters. This pooling typically happens behind the scenes in software packages. Although the model outputs are the pooled coefficients from the individual analyses, the results can be interpreted in the same manner as one would in a normal setting.

In summary, multiple imputation works well when the missing data are MCAR or MAR and is particularly useful with MAR data. It helps to maintain the sample size and eliminate the potential selection that could result if cases with incomplete data were dropped. It also helps to reduce the likelihood of standard error bias. The three steps to analyzing imputed data are: 1) imputing values for the missing data 2) running theoretical analysis using that imputed data and 3) pooling estimates into a single set of results. The first steps involves imputing the missing values to generate an appropriate number of data sets. The number of imputations needed is dependent on the amount of missingness; the greater the percent of data that is missing, the larger the number of imputations are needed. Each imputation results in the creation of complete another data set. The second step is analyzing the imputed data as part of the researcher's theory-based model. This involves running the analysis simultaneously across each imputed data set. In most statistical packages, the researcher does this by specifying the data as imputed and proceeding largely as they otherwise would. The researcher does not typically see the output of this second step. The final step is pooling those results. Pooling generates a single output that incorporates into its standard errors all of the potential uncertainty inherent in the imputation process.

TABLE 2 HERE

Description and Illustration of ICSD Missing Data

The following sections illustrate the relative advantages and disadvantages of each approach by applying it to the awesome new urban studies data resource, the Integrated City Sustainability Database (ICSD) (Feiock et al 2014). We compare listwise deletion, single mean replacement, and multiple imputation techniques to demonstrate the value-added from using multiple imputation when the degree of missingness can have an impact on the outcome of

analysis.

A recent article in this journal by Feiock and colleagues (2014) describes the “Integrated City Sustainability Database” (ICSD) as a solution to the challenges associated with missing data in urban research. The ICSD combines the results of seven national surveys of city sustainability programs that were administered within an 18-month period in 2010-2011 into one comprehensive national data set. Table 3 presents basic information on the seven ICSD component surveys.² The process of survey harmonization yielded a large sample: 2,825 cities completed at least one of the seven surveys. However, the majority of cities did not answer all seven of the surveys meaning that the ICSD contains a considerable amount of missing data.

The first generation of the ICSD utilizes a single regression replacement method to account for missing data (Feiock et al., 2014). The authors deal with missing observations within and across the surveys using a two-stage informed single regression imputation technique, which produced a single unified data set through a two-stage version of single imputation. The first stage imputed missing data within each completed survey and the second used this data to impute across surveys, taking into consideration the different types of missingness. This process generates a single unique value for each missing observation in the original ICSD and results in a single complete data set for the ICSD. This structure facilitates accessibility since users can download and use a single file of imputed data. This “first generation” ICSD represents a meaningful advancement that enables more confident conclusions to be drawn from the results of empirical analyses of local politics, governance and policy (Feiock and Hawkins 2016). It

² The ICSD is a dynamic database that is expanding and anticipate to continue to grow over time as new data on city level sustainability is collected. The original ICSD establishes a 2010/2011 baseline on local sustainability initiatives. As more data is collected by the authors and others it will be added to the ICSD to enable analyses of change over time.

provides single imputed data for an extensive set of cities including both large and small cities that is already being widely used in urban research.³

The two-stage single imputation approach of the first generation database is a significant improvement over listwise deletion, but further improvement is possible through the process of multiple imputation for the cities over 50,000 population. Cities above this population threshold were included in the sample frames for all seven surveys, making their overall levels of missing data lower and making them better candidates for multiple imputation. The second generation ICSD described here compliments the first generation database by providing a multiple imputation version for this subset of ICSD cities.

TABLE 3 HERE

The 683 US cities, which per the 2010 census had populations over 50,000, were included in the sample for each of the seven ICSD component surveys. Their response was particularly strong, with 90 percent of these cities responding to at least one survey. This virtually eliminates self-selection bias among this sub-sample and provides a unique opportunity to examine the sustainability policy, implementation, resources, obstacles, and motivations in medium and large US cities. However, although they all shared a related scope, each survey utilized a somewhat different set of questions and response categories and ended up with a different set of responding cities. This is problematic in a multivariate context where models seek to draw information from across several surveys, because it can drastically reduce the sample size of available data. This reduction in sample size provides an important rationale for utilizing a more advanced method of dealing with missing data, such as multiple imputations.

³ The public release of the ICSD is scheduled for January 2018 <http://localgov.fsu.edu/ICSD/>

[Figure 1 About Here]

Figure 1 summarizes the process used to identify the theoretic and statistically relevant variables that were used as informing variables in the imputation process for the second generation ICSD. The theoretical linkages were determined by developing two “general concepts” – one related to the “activity” and the other to “subject matter” – for every question contained within the seven surveys. For example, the question “Do any of your city’s efforts to encourage retrofits for energy efficiency include: Partnership or collaboration with nonprofit community organizations” is labeled with the activity concept of “Collaboration” and the subject matter concept “Energy”. This develops sets of potentially theoretically related questions – called the concept list. A list of these concepts and how often they are attributed to variables in the ICSD surveys is presented in Table 4.

[Table 4 about here]

The concept lists develop broad groupings of variables that have theoretic relationships and inform one another. In other words, these ‘informing variables’ act almost as independent variables that may provide information to help predict missing values of a particular target variable. In some cases, the theoretically derived list of informing variables is too large to support convergence of the model determining the value of the missing responses and therefore statistical correlations are used to narrow the set. With the objective of identifying a small enough number of informing variables to enable statistical conversion, 0.2 was selected as the minimum correlation⁴ between the variable being imputed and the potential informing variables.

⁴ The 0.2 correlation value selected for this specific data set indicated that a predictor was related to the variable being imputed. Anything below the 0.2 cutoff was deemed unrelated to the variable being imputed. The 0.2 correlation narrowed the related concept list enough to allow convergence and did not eliminate the theoretically related questions to zero in any case.

As a result, only variables that are theoretically and statistically relevant are retained as predictors, resulting in an average of 95 informing variables for each target variable in the data set.

A distribution of the non-missing cases is used to determine the expectation of the distribution for missing responses. For example, if the non-missing responses are normally distributed the imputed responses will maintain a normal distribution. The distribution assigned is variable specific. Twenty imputations are generated for the results of the analysis that determines the value of a missing response. This process is repeated for all missing variables across the seven surveys. For the 683 cities with populations above 50,000, per the 2010 census, complete data is generated for each of the 1,010 variables in the ICSD.

Compared to conducting analysis using either non-imputed or first generation ICSD data, utilizing the multiply imputed data generated from the process described above requires a few additional steps. The STATA code associated with these steps for several different types of analytic techniques have been included in the online appendix. As the code demonstrates, it is quite simple to analyze the imputed data. It primarily requires setting the data as multiply imputed and analyzing using `'mi estimate:'` prior to writing the code as usual.

One complication with analyzing multiply imputed data is the generation of summary statistics. The goal of multiple imputations is to avoid generating a fixed point-estimate for the prediction of the missing value. Generating summary statistics of a single imputed data set, or each independently, would treat each data set as holding a true value for the missing observation. Therefore, traditional summary statistics are an inappropriate match for the technique because they do not account for the uncertainty inherent in the imputation. It may be more appropriate to report either a grand mean, which estimates the average of the multiply imputed data sets

averages, and/or the descriptive statistics from the original, un-imputed data.

A Comparison of Approaches Using the ICSD

We use the ICSD survey data in their raw and two imputed forms to demonstrate the relative performance of each of the three approaches to dealing with missing data: listwise deletion, value replacement, and multiple imputation. For illustration purposes, we examine the factors that influence local action on sustainability in a generic empirical model that corresponds to those typical in the urban affairs literature.

Dependent Variable

The dependent variable is an additive index of the number of environmental sustainability-related policies and actions that cities reported having implemented in their jurisdictions. The additive index is a common dependent variable in quantitative studies of local sustainability (Portney 2003; Krause 2012; Bae and Feiock 2013). We select a dependent variable conducive to analysis using Ordinary Least Squares regression. Sixteen sustainability actions are included in this index and cluster in three primary areas: energy, transportation, and waste disposal.

Independent Variables

The independent variables reflect common operationalization of hypotheses in sustainability studies and relate to cities' motivations to engage in sustainability, obstacles hindering their action, and a series of control variables (Krause 2013; Krause et al., 2016; Hawkins et al., 2016.) The independent variables are intentionally drawn from a limited number of the different ICSD component surveys. The "EECBG Grantee Implementation Survey" supplies the three motivation independent variables: achieving energy cost savings, the desire to

build a sustainable community, and external public pressure. Two of the obstacle variables – lack of staff capacity and lack of information resources – likewise come from the EECBG Grantee Implementation Survey. The third obstacle – a lack of political will – is pulled from the Implementation of Energy Efficiency and Sustainability Programs Survey.⁵

Control variables include population density, per-capita income, form of government, ICLEI membership, percent of racial minority residents, and residents' educational attainment. Each of these control variables have been used in previous studies regarding sustainability policy (Krause 2010; Lubell et al. 2009; Zahran et al. 2008; Feiock et al. 2010; Salon, Murphy & Sciara 2014). The data was collected from the US Census Bureau, the International City/County Management Association, and ICLEI Local Governments for Sustainability, and thus have near complete coverage.

Results

We employ ordinary least squares regression analysis and have examining the tradeoffs between using different approaches to deal with missing data as our primary objective. We use three identical models to estimate the impact of the different missingness treatments. The first model uses listwise deletion to handle the missingness in the survey data, the second uses single imputation mean replacement, and the third uses multiple imputations, which is the approach utilized in the second generation Integrated City Sustainability Database. In order to understand the relevance of the different missing data treatments readers must indulge the cult of statistical significance (Ziliak and McCloskey, 2008).

⁵ We only incorporate variables from three of the seven surveys in this model, which should keep the loss of observations from listwise deletion relatively low. This is done to demonstrate that a more advanced treatment of missing data may be valued even without extreme degrees of missing observations. In other words, we are giving the listwise deletion approach its 'best chance' of success.

Table 5, column two reports the results from the model utilizing *listwise deletion*. Only 111 of the 683 cities with populations over 50,000 remain in the model after listwise deletion removes all observations with incomplete data (a loss of 572). The results using this approach indicate that only one variable – ICLEI membership – has a statistically significant effect on the policy index. The information loss resulting from the drastic reduction in sample size and the potential bias of the complete observations may contribute to the production of null findings in terms of motivations and obstacles to implementing policy.

The third column in Table 5 presents the results of the model using *mean replacement*. For each independent variable in the model, this technique simply replaces the missing observations with the mean value for that variable. This technique increases the size of the sample from 111 to 325. However, it still results in a total loss of 358 observations.⁶ The results generated using mean replacement identify several additional statistically significant relationships compared to listwise deletion. Lack of political will, as well as the control variables population density and education are now statistically significant. ICLEI membership remains significant and the magnitude of its effect is larger. Perhaps the most meaningful change in the results is that, using mean replacement, lack of political will has a negative statistically significant relationship to the policy index dependent variable. Cities lacking political will towards sustainability implement approximately one-half fewer policies than those reporting stronger political will. This suggests that listwise deletion lost a significant amount of variation by deleting observations with incomplete data. However, the concern associated with mean replacement is that the observed significant relationships between the variables may not be true

⁶ This is because utilizing mean replacement for dependent variables is a debated procedure. If the dependent variable were mean replaced, the data would have the full 683 cities.

due to underestimates of the standard deviation and standard error. Ordinary least squares – regression to the mean—is not able to accurately measure variations from the mean (i.e. error) because observations are artificially held at the value of the mean. Therefore, even though these variables are significant, the resulting p-values should be interpreted with caution.

The results from the analysis performed using informed *multiple imputation* are shown in the fourth column and yield a slightly different combination of statistically significant variables in the model, when compared to the other two approaches. Multiple imputation is typically accepted for use in the dependent as well as independent variables (Young and Johnson 2010), which enables the sample size to increase from 325 to 683. In this model, the motivation to build a sustainable community variable is statistically significant and positively associated with the policy index. ICLEI membership and lack of political remain statistically significant, however, the magnitude of both decrease slightly compared to the other models. This model also yields statistically significant relationships for motivation and obstacle variables. Comparing these results to those from the listwise deletion model suggests that null findings in cases with large amounts of missingness may not be null findings after-all. The standard errors in multiple imputation incorporate the uncertainty from the 20 imputation results giving us confidence in the resulting p-values.

Discussion and Conclusion

Listwise deletion, value replacement, and multiple imputation are common approaches for address missing data. Each is associated with particular advantages and disadvantages; and, depending on the nature of the missingness, using the wrong method may provide inaccurate, biased, or inappropriate null findings. This paper elucidated these consequences and specifically described how inaccurate treatment can decrease the power of the sample size, increase the

potential for biased results, and over or under estimate standard errors. This is not to say that multiple imputation is the correct or best solution to dealing with missing data. In fact, this paper suggests that the categorization of missing data should drive the selection of an appropriate approach to dealing with missing data.

Although often a default, listwise deletion is not a blanket solution to missing data problems. Dropping observations from an analysis decreases its power and its overuse may cause variables that help explain the outcome variable to be deemed insignificant. Also problematic is the potential of incorporated bias in the selection process. Listwise deletion might work for data that is Missing Completely at Random (MCAR), but data is very rarely MCAR. It is also possible that techniques such as mean replacement are suitable for use with MCAR data. However, it may result in the effect of these variables being vastly over estimated because the standard errors are made artificially smaller by holding the values to the mean. Multiple imputation, although more complicated, provides theoretically consistent results works for data that is Missing at Random (MAR). Incomplete observations are not dropped from the analysis and, by incorporating the uncertainty of missing responses into the standard errors, the magnitude and significance of the relationships between independent and dependent variables are appropriately measured.

Exploiting the Integrated City Sustainability Database allows us to examine the implications of various treatments of missing data. The second generation ICSD database contains data generated by informed multiple imputation, which enables analysis with larger sample size, less bias, and the ability to interpret the data as though it was not missing. In addition, this technique is applicable to data that is either MAR or MCAR. A large degree of the missingness in the ICSD can be attributed to survey recipient response, which makes multiple

imputation an appropriate choice. However, some variables may not be MAR and therefore should be considered thoughtfully prior to applying this technique. In addition to being more complicated, a disadvantage to using multiply imputed data is that it is not conducive to the generation of standard descriptive statistics, including things like variable means, and basic model fit indicators like R^2 .

In urban studies and across the social sciences there are increasing expectations for rigor and transparency in the management of data including procedures for dealing with missing observations. This is manifested in the Transparency and Openness Promotion (TOP) guidelines that are being adopted by many journals (Nosek et al. 2015). It is our hope that urban scholars begin to treat missing data more explicitly and openly. Included here is an online appendix, with multiple imputation code and description to aid in the utilization process. In 2018, the multiply imputed data included in the second generation ICSD will be made publicly available. In the meantime, select variables from the first generation ICSD are available at <http://localgov.fsu.edu/ICSD/>.

References

- Abayomi, Kobi, Andrew Gelman, and Marc Levy. "Diagnostics for multivariate imputations." *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 57, no. 3 (2008): 273-291.
- Allen, Tammy D., Lillian T. Eby, and Elizabeth Lentz. "Mentorship behaviors and mentorship quality associated with formal mentoring programs: closing the gap between research and practice." *Journal of Applied Psychology* 91, no. 3 (2006): 567.
- Allison, P. D. (2002). *Missing data*. Newbury Park, CA: Sage.
- Andridge, Rebecca R., and Roderick JA Little. "A review of hot deck imputation for survey non-response." *International statistical review* 78, no. 1 (2010): 40-64.
- Bae, Jungah, and Richard Feiock. "Forms of government and climate change policies in US cities." *Urban Studies* 50, no. 4 (2013): 776-788.
- Betsill, Michele M. "Mitigating climate change in US cities: opportunities and obstacles." *Local environment* 6, no. 4 (2001): 393-406.
- Bodner, Todd E. "What improves with increased missing data imputations?." *Structural Equation Modeling* 15, no. 4 (2008): 651-675.
- Brick, J. Michael, and Graham Kalton. "Handling missing data in survey research." *Statistical methods in medical research* 5, no. 3 (1996): 215-238.
- DeMaris, Alfred. "Combating unmeasured confounding in cross-sectional studies: evaluating instrumental-variable and Heckman selection models." *Psychological methods* 19.3 (2014): 380.
- Donders, A. Rogier T., Geert JMG van der Heijden, Theo Stijnen, and Karel GM Moons. "Review: a gentle introduction to imputation of missing values." *Journal of clinical*

- epidemiology* 59, no. 10 (2006): 1087-1091.
- Downey, Ronald G., and Craig V. King. "Missing data in Likert ratings: A comparison of replacement methods." *The Journal of general psychology* 125, no. 2 (1998): 175-191.
- Feiock, Richard C., In Won Lee, Hyung Jun Park, and Keon-Hyung Lee. "Collaboration networks among local elected officials: Information, commitment, and risk aversion." *Urban Affairs Review* 46, no. 2 (2010): 241-262.
- Feiock, Richard C., and Jungah Bae. "Politics, institutions and entrepreneurship: city decisions leading to inventoried GHG emissions." *Carbon Management* 2, no. 4 (2011): 443-453.
- Feiock, Richard C., Rachel M. Krause, Christopher V. Hawkins, and Cali Curley. "The integrated city sustainability database." *Urban Affairs Review* 50, no. 4 (2014): 577-589.
- Fox, James Alan, and Marc L. Swatt. "Multiple imputation of the supplementary homicide reports, 1976–2005." *Journal of Quantitative Criminology* 25, no. 1 (2009): 51-77.
- Gallimore, Jonathan M., Barbara B. Brown, and Carol M. Werner. "Walking routes to school in new urban and suburban neighborhoods: An environmental walkability analysis of blocks and routes." *Journal of Environmental Psychology* 31, no. 2 (2011): 184-191.
- Graham, John W., Allison E. Olchowski, and Tamika D. Gilreath. "How many imputations are really needed? Some practical clarifications of multiple imputation theory." *Prevention Science* 8, no. 3 (2007): 206-213.
- Hawkins, Christopher V., Rachel Krause, Richard Feiock, and Cali Curley. (2016). Making Meaningful Commitments: Accounting for Variation in Cities' Investments of Staff and Fiscal Resources to Sustainability. *Urban Studies*, 53(9): 1902-1924
- Johnson, David R., and Rebekah Young. "Toward best practices in analyzing datasets with missing data: Comparisons and recommendations." *Journal of Marriage and Family* 73.5

- (2011): 926-945.
- Jones, Michael P. "Indicator and stratification methods for missing explanatory variables in multiple linear regression." *Journal of the American statistical association* 91, no. 433 (1996): 222-230.
- King, Gary, James Honaker, Anne Joseph, and Kenneth Scheve. "Analyzing incomplete political science data: An alternative algorithm for multiple imputation." In *American Political Science Association*, vol. 95, no. 01, pp. 49-69. Cambridge University Press, 2001.
- Kong, A., Liu, J. S., & Wong, W. H. (1994). Sequential imputations and Bayesian missing data problems. *Journal of the American statistical association*, 89(425), 278-288.
- Krause, Rachel M. "Policy innovation, intergovernmental relations, and the adoption of climate protection initiatives by US cities." *Journal of urban affairs* 33, no. 1 (2010): 45-60.
- Krause, Rachel M. "Political decision-making and the local provision of public goods: the case of municipal climate protection in the US." *Urban studies* 49, no. 11 (2012): 2399-2417.
- Krause, Rachel M. "The motivations behind municipal climate engagement: An empirical assessment of how local objectives shape the production of a public good." *Cityscape* (2013): 125-141.
- Krause, Rachel M, Richard Feiock, and Christopher Hawkins. "The Administrative Organization of Sustainability Within Local Government." *J Public Adm Res Theory* (2016) 26 (1): 113-127
- Little, Roderick J. "Selection Model (Missing Data)." *Wiley StatsRef: Statistics Reference Online* (2016). Little, R. J., & Rubin, D. B. (2014). *Statistical analysis with missing data*. John Wiley & Sons.
- Little, Roderick JA. "Missing-data adjustments in large surveys." *Journal of Business &*

- Economic Statistics* 6, no. 3 (1988b): 287-296.
- Little, Roderick JA, and Donald B. Rubin. *Statistical analysis with missing data*. John Wiley & Sons, 2014.
- Little, Roderick JA. "A test of missing completely at random for multivariate data with missing values." *Journal of the American Statistical Association* 83, no. 404 (1988a): 1198-1202.
- Lubell, Mark, Richard Feiock, and Susan Handy. "City adoption of environmentally sustainable policies in California's Central Valley." *Journal of the American Planning Association* 75, no. 3 (2009): 293-308.
- Myers, T. A. (2011). Goodbye, listwise deletion: Presenting hot deck imputation as an easy and effective tool for handling missing data. *Communication Methods and Measures*, 5(4), 297-310.
- Miyama, Eriko, and Shunsuke Managi. "Global environmental emissions estimate: application of multiple imputation." *Environmental Economics and Policy Studies* 16, no. 2 (2014): 115-135.
- Nosek, Brian A., George Alter, G. C. Banks, D. Borsboom, S. D. Bowman, S. J. Breckler, S. Buck et al. "Promoting an open research culture." *Science* 348, no. 6242 (2015): 1422-1425.
- Park, Joohyung, and Sejin Ha. "Understanding pro-environmental behavior." *International Journal of Retail & Distribution Management* 40, no. 5 (2012): 388.
- Peng, Chao-Ying Joanne, Michael Harwell, Show-Mann Liou, and Lee H. Ehman. "Advances in missing data methods and implications for educational research." *Real data analysis* (2006): 31-78.
- Portney, Kent E. *Taking sustainable cities seriously: Economic development, the environment,*

- and quality of life in American cities*. Vol. 67. MIT Press, 2003.
- Rose, Roderick A., and Mark W. Fraser. "A simplified framework for using multiple imputation in social work research." *Social Work Research* 32.3 (2008): 171-178.
- Royston, Patrick, and Ian R. White. "Multiple imputation by chained equations (MICE): implementation in Stata." *Journal of Statistical Software* 45, no. 4 (2011): 1-20.
- Rubin, Donald B. *Multiple imputation for nonresponse in surveys*. Vol. 81. John Wiley & Sons, 2004.
- Rubin, Donald B. "Multiple imputation after 18+ years." *Journal of the American statistical Association* 91, no. 434 (1996): 473-489.
- Ryff, Carol D., and Corey Lee M. Keyes. "The structure of psychological well-being revisited." *Journal of personality and social psychology* 69, no. 4 (1995): 719.
- Salon, Deborah, Sinnott Murphy, and Gian-Claudia Sciara. "Local climate action: motives, enabling factors and barriers." *Carbon Management* 5.1 (2014): 67-79.
- Schafer, Joseph L. *Analysis of incomplete multivariate data*. CRC press, 1997.
- Schafer, Joseph L., and John W. Graham. "Missing data: our view of the state of the art." *Psychological methods* 7, no. 2 (2002): 147.
- Schneider, Tapio. "Analysis of incomplete climate data: Estimation of mean values and covariance matrices and imputation of missing values." *Journal of Climate* 14, no. 5 (2001): 853-871.
- Van der Heijden, Geert JMG, A. Rogier T. Donders, Theo Stijnen, and Karel GM Moons. "Imputation of missing values is superior to complete case analysis and the missing-indicator method in multivariable diagnostic research: a clinical example." *Journal of clinical epidemiology* 59, no. 10 (2006): 1102-1109.

- White, Ian R., Rhian Daniel, and Patrick Royston. "Avoiding bias due to perfect prediction in multiple imputation of incomplete categorical variables." *Computational statistics & data analysis* 54.10 (2010): 2267-2275.
- White, I. R., Royston, P., & Wood, A. M. (2011). Multiple imputation using chained equations: issues and guidance for practice. *Statistics in medicine*, 30(4), 377-399.
- Young, Rebekah, and David R. Johnson. "Imputing the missing Y's: Implications for survey producers and survey users." In *Proceedings of the AAPOR conference abstracts*, pp. 6242-6248. 2010.
- Zahran, Sammy, Himanshu Grover, Samuel D. Brody, and Arnold Vedlitz. "Risk, stress, and capacity: Explaining metropolitan commitment to climate protection." *Urban affairs review* 43, no. 4 (2008): 447-474.
- Zhang, Paul. "Multiple imputation: theory and method." *International Statistical Review* 71, no. 3 (2003): 581-592.
- Ziliak, Stephen Thomas, and Deirdre N. McCloskey. *The cult of statistical significance: How the standard error costs us jobs, justice, and lives*. University of Michigan Press, 2008.

Appendix: STATA Multiple Imputation Code

***The following is Multiple Imputation code as related to using the ICSD imputed data for STATA.

***Please see <http://XXXXXX/> for details on what is currently available for public use.

**Read in data as usual.

**Import the data as an imputed file or ice object

```
mi import ice, automatic
```

** Get a list of all commands for mi estimation, any of these commands can be used to analyze data as you normally would.

```
help mi estimation
```

** In order to use linear regression with continuous DV and an X variable. Options are typically added before the colon

```
mi estimate : regress Y_variablename X_variablename
```

**Logistic regression with dichotomous DV and an X variable and code to set a variables value to dichotomous

```
recode variablename 1 = 0 2 = 1
label define variablename 1 "Yes" 0 "no", replace
mi estimate : logistic Y_variablename X_variablename
```

**Ordinal-response regression

```
mi estimate : ologit Y_variablename X_variablename
```

**Multinomial logistic regression, items with more than 2 response options that are not ordered.

```
mi estimate : mlogit Y_variablename X_variablename
```

**In order to look at means across imputations or proportion of responses across imputations use the following code. These statistics are how to calculate the variance across imputations (level of uncertainty).

```
mi estimate : mean variablename
```

```
mean variablename if _mi_m == 0
```

```
mi estimate : proportion variablename
```

```
proportion variablename if _mi_m == 0
```

* Here's some code to run the individual regressions, save the

* R-squares, and summarize them for you.

```

* Define loop
qui sum _mi_m, detail
local imax = r(max)

* Create empty matrix for R-squared values
mata: R = J(`imax',1,.)

* Run regressions, save R-squared
foreach j of numlist 1/`imax' {
    qui reg Y_variablename X_variablename if _mi_m==`j' // the only thing to change is the
    regression variables in this line //
    local r2 = e(r2)
    mata: R[`,1] = `r2'
}

mata: mean = mean(R)
mata: median = mm_quantile(R,1,.5)
mata: st_numscalar("r2mean", mean[1,1])
mata: st_numscalar("r2med", median[1,1])

di "The mean R-squared is: " r2mean
di "The median R-squared is: " r2med

```

Missing Completely at Random (MCAR)	Missing at Random (MAR)	Missing Not at Random (MNAR, non-ignorable)
Missingness is independent from characteristics of either the observed data or the unobserved values in the data set	Missingness is entirely explained by the observed data, i.e. after observed values are accounted for, missingness is randomly distributed.	Missing observations are dependent upon unobserved values; missingness cannot be accounted for by controlling for observed data.

Table 2: Techniques of Imputation*

TECHNIQUES	Listwise Deletion (Complete Case Analysis)	Single Imputation		Multiple Imputation
		Mean Replacement (Mean Substitution)	Single Regression Replacement	
Technique Summary	Remove any entries with missing values; perform analysis without these observations	For variable "a" with missing values, take the mean of all included observations. Substitute the mean of "a" for missing values of "a."	Estimate the distribution of the missing variable(s) given covariates; take a random draw from this distribution for each value; perform analysis as usual**	Estimate the distribution (Bayesian posterior distribution) of the missing variable, given covariates; take random draws from this distribution to produce multiple versions (usually 3-10) of an imputed data set; Perform analysis on each imputed data set and pool the results
Missingness Assumption	MCAR, occasionally MAR	MCAR	MCAR or MAR	MCAR or MAR
Advantages	Easiest, simplest	Preserves the mean of the dataset; Simple; allows use of all observations	Avoids bias in estimating; simpler than multiple imputation	Accounts for the extra uncertainty produced by imputing data; produces better estimates of missing values
Disadvantages	Loses valuable information; potentially contributes to bias	Artificially reduces standard deviation of data set, distorts relationships between variables	Misrepresents uncertainty of estimates; more complicated than listwise deletion or mean replacement	Requires complicated statistical methods or complicated software; harder to understand; takes extra steps
Impacts on Interpretation	Statistical analysis loses power; estimates could be biased if data is not missing completely at random	Estimate could be biased, Standard errors will be artificially low; Could produce results that are highly statistically significant, but inaccurate	Although theoretically unbiased, reduces confidence intervals of estimates;	Because the method accounts for extra uncertainty, results can be interpreted as if data was not missing.
References				
Method Exploration	Jones 1996, 223; Schafer and Graham 2002, 155.	Downey and King 1998; Shafer and Graham 2002, 159.	Donders et al. 2006, 1088- 1089; Schneider 2001; van der Heijden et al. 2006;***	Donders et al. 2006, 1089; King et al. 2001; Rubin 1987; Schafer 1997; Zhang 2003;
Application	Park and Ha 2012, 394; Ryff and Keyes 1995, 722.	Allen et al. 2006, 572; Gallimore et al 2011, 186- 187		Abayomi et al. 2008; Fox and Swatt 2009; Miyama and Managi 2014;

*Additional missingness reference can be found in Schafer and Graham 2002, 151.

**Single Imputation, defined more broadly, includes any method that replaces missing data with a single value. This would include mean replacement and hot deck imputation; the latter is summarized by Andridge and Little 2010.

***Applications of the single imputation technique are limited; these are primarily theoretical explorations of the technique.

Table 3. Characteristics of the Surveys Comprising the Integrated City Sustainability Database.			
Survey Name	Sampling Frame	Respondents	Response Rate (%)
ICMA Local Government Sustainability Policies and Programs Survey	8,569 local governments with a population of 10,000 or more residents	2,176	25.4
NLC Sustainability Survey	1,708 mayors in cities over 10,000	442	26.6
EECBG Grantee Implementation Survey	970 municipal governments receiving EECBG awards, including all cities over 30,000	747	77
Implementation of Energy Efficiency and Sustainability Programs	1,180 cities: all with populations over 50,000 and a random sample of 500 cities with populations between 20,000 and 50,000	679	57.5
National Survey of Sustainability Management in U.S. Cities	601 cities with populations over 50000	263	44
Municipal Climate Protection Survey	664 cities with populations over 50000	329	49.5
Municipal Government Questionnaire	425 cities with populations over 50,000 that have indicated explicit involvement in climate protection	255	60
Note. ICMA = International City/County Management Association; NLC = The National League of Cities; EECBG = Energy Efficiency and Conservation Block Grant.			

General Concept	Category	Description/Keywords	Count*
Climate	Subject Matter	Climate change, climate protection, adaptation	71
Economic	Subject Matter	Green business, green jobs, buy local programs, farmers' market	50
EECBG	Subject Matter	Energy Efficiency Conservation Block Grant, American Resource and Recovery Act (ARRA), stimulus	109
Energy	Subject Matter	Energy, energy efficiency, energy conservation	306
Environment	Subject Matter	Land use, water, recycling, trees, community gardens, food	122
Social	Subject Matter	Low-income, population, health, equity	32
Sustainability	Subject Matter	Sustainability	172
Transportation	Subject Matter	Vehicles, car-pooling, telework, condensed/flexible work days	69
Collaboration	Activity	Collaboration in general, partnership, cooperation	70
Community action	Activity	Any policy or programmatic action (loan program, tax credit, rebates, regulation, retrofit) that targets the community at large	114
Community planning	Activity	inventory from community-wide emissions,	7
Contracting	Activity	Contracting, outsourcing	29
General action	Activity	Any policy or programmatic action that does NOT specify target groups	93
General Planning	Activity	planning, adopted planning goals, adopted policy	36
Government Action	Activity	Any policy or programmatic action targeting government operations (publicly-owned building, purchase (credits), incentives, utility retrofit)	128
Government Planning	Activity	goal, inventory from city government operations	9
Infrastructure	Activity	own operate, facility	46
Inter-department	Activity	Coordinate within the city	46
Inter-governmental	Activity	Collaborate with other localities, state/federal government, cross-influence	59
Motivation	Activity	Why, What are the drivers of action?	45
Obstacle	Activity	Why not, Barriers	46
Performance measures	Activity	measurement, resulting from efforts, indicators, evaluation	58
Priority	Activity	How important?	47

Public Engagement	Activity	Public education, info center, engage with...	31
Resources	Activity	Designated staff, money, funding	73
*Represent number of variables characterized as general concept			

		Listwise Deletion		Mean Replacement		Multiple Imputation	
		Policy Index	Standard Error	Policy Index	Standard Error	Policy Index	Standard Error
M	Reduced energy cost	-0.066	0.678	-0.402	0.385	-0.211	0.150
M	Sustainable Communities	0.251	0.395	0.396	0.277	0.328**	0.136
M	Public Pressure	0.446	0.354	0.394	0.255	0.145	0.129
O	Staff Capacity	0.355	0.384	0.067	0.283	0.070	0.188
O	Lack of Information	0.085	0.461	-0.165	0.320	-0.080	0.198
O	Lack of Political Will	-0.3027	0.377	-0.627**	0.297	-0.550***	0.177
C	Population Psq mile	0.000	0.000	0.0001*	0.000	0.000	0.000
C	Percapita income	0.000	0.000	0.000	0.000	0.000	0.000
C	Iclei member 2010	1.149**	0.582	1.814***	0.332	1.013***	0.255
C	Council Manager	-3.372	2.763	-2.872	2.705	-0.384	0.497
C	Mayor Council	-3.112	2.754	-2.973	2.702	-0.336	0.514
C	Percent Minority	0.012	0.015	-0.008	0.008	-0.004	0.006
C	Percent bachelors+	0.051	0.034	0.033*	0.020	0.013	0.014
	Constant	8.473**	3.464	9.757***	3.015	8.170***	0.886
	Sample Size	111		325		683	
		Adj R2	0.0906	Adj R2	0.1814	Prob >F	0
Motivation Variable (M), Obstacle Variable (O), Control Variable (C)							

Figure 1: Process Flow of Informed Multiple Imputation

