



Comparison of Modern Imputation Methods for Missing Laboratory Data in Medicine

Journal:	<i>BMJ Open</i>
Manuscript ID:	bmjopen-2013-002847
Article Type:	Research
Date Submitted by the Author:	07-Mar-2013
Complete List of Authors:	Waljee, Akbar; VA Center for Clinical Management Research, Ann Arbor VA Medical Center; University of Michigan, Department of Internal Medicine Mukherjee, Ashin; University of Michigan, Department of Statistics Singal, Amit; UT Southwestern Medical Center, Department of Internal Medicine Zhang, Yiwei; University of Michigan, Department of Statistics Warren, Jeffrey; University of Michigan, Pathology Balis, Ulysses; University of Michigan, Pathology Marrero, Jorge; UT Southwestern Medical Center, Department of Internal Medicine Zhu, Ji; University of Michigan, Department of Statistics Higgins, Peter; University of Michigan, Department of Internal Medicine
Primary Subject Heading:	Gastroenterology and hepatology
Secondary Subject Heading:	Diagnostics, Gastroenterology and hepatology, Health informatics
Keywords:	Gastroduodenal disease < GASTROENTEROLOGY, Information technology < BIOTECHNOLOGY & BIOINFORMATICS, Inflammatory bowel disease < GASTROENTEROLOGY

SCHOLARONE™
Manuscripts

Comparison of Modern Imputation Methods for Missing Laboratory Data in Medicine

Ashin Mukherjee MS¹, Akbar K. Waljee MD MS^{2,3}, Amit G. Singal MD MS^{4,5}, Yiwei Zhang MS¹, Jeffrey Warren MD⁶, Ulysses Balis MD⁶, Jorge Marrero⁴, Ji Zhu PhD¹, Peter DR Higgins MD PhD²

¹Department of Statistics, University of Michigan, Ann Arbor, MI

²Department of Internal Medicine, University of Michigan, Ann Arbor, MI

³Veterans Affairs Center for Clinical Management Research, Ann Arbor, MI

⁴Department of Internal Medicine, UT Southwestern Medical Center, Dallas, TX

⁵Department of Clinical Sciences, UT Southwestern, Dallas, TX

⁶Department of Pathology, University of Michigan, Ann Arbor, MI

Correspondence:

Akbar K. Waljee, M.D., M.Sc.
VA Center for Clinical Management Research,
Ann Arbor VA Medical Center,
2215 Fuller Road, IID,
Ann Arbor, MI 48105
e-mail: awaljee@med.umich.edu

Financial Support: Dr. Waljee's research is funded by a VA HSR&D CDA-2 Career Development Award 1IK2HX000775. Dr. Singal's research is funded by an ACG Junior Faculty Development Award and grant number KL2 RR024983-05. Dr. Higgins' research is supported by NIH R01 GM097117. The content is solely the responsible of the authors and does not necessarily represent the official views of UT-STAR, the University of Texas Southwestern Medical Center at Dallas and its affiliated academic and health care centers, the National Center for Advancing Translational Sciences, or the National Institutes of Health.

1
2 **Acknowledgements:**
3
4

5 Ashin Mukherjee and Akbar Waljee contributed equally to this work and should be considered co-first
6 authors.
7

8
9 The content is solely the responsibility of the authors and does not necessarily reflect the official views
10 of the Department of Veterans Affairs, the National Center for Research Resources, or the National
11 Institutes of Health.
12
13

14
15
16 **Conflicts of Interest:** The authors disclose no conflicts.
17
18
19

20
21 Abstract Word Count: 191
22

23 Manuscript Word Count: 2671
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Article Summary:

1) Article Focus

- Multi-analyte Assays with Algorithmic Analyses (MAAAs) are a relatively new approach to leveraging value from laboratory data to predict clinical outcomes. It is not known how robust MAAA models are when individual laboratory data points are missing.
- Recent developments in machine learning have used laboratory data to build MAAA models to predict health care outcomes - These models can be sensitive to missing laboratory data. It is not known whether modern imputation methods can robustly address the problem of missing data, and whether predictive models will remain accurate when imputed values are used.
- Multiple methods have been developed in order to deal with missing data, including single imputation, multiple imputation, multiple imputation by chained equations (MICE), nearest neighbor estimation, and missForest. Although these have been shown to be effective with other types of missing data, little data exists regarding the absolute or comparative effectiveness of these methods in accurately imputing missing laboratory data in predictive models. The aims of our study were: (1) to compare the accuracy of four different imputation methods for missing laboratory data in two large data sets; and (2) to compare the effect of imputed values from each method on the accuracy of predictive models based on these data sets.

2) Key Message

- We found that missForest methods consistently produced the lowest imputation error and had the smallest prediction difference when models used imputed laboratory values.
- The small absolute changes in predictions with these models, despite 10-30% missing laboratory data, speak to the robustness of these multi-analyte assays with algorithmic analysis (MAAAs).
- With increasing complexity of these models, and increasing numbers of analytes, the risk of missing values increases, and methods to cope with missing values and preserve the accuracy of the model are needed. missForest appears to be a robust and accurate approach to the issue of missing laboratory values when used in these two MAAAs.

3) Limitations and Strengths

- The main limitations of missForest as a solution to missing laboratory data for predictive modeling applications are: a requirement for skilled R programming for implementation, and slightly more demanding computational needs, compared to NN or MICE methods, with minimal increases in processing time (generally in the range of 10-20 sec).
- The strength of this is that missForest methods consistently produced the lowest imputation error and had the smallest prediction difference when models used imputed laboratory values and that it is a readily available freeware R package making it a very convenient solution for any practical missing value problems.

Abstract

Background: Missing laboratory data is a common issue, but the optimal method of imputation of missing values has not been determined. The aims of our study were to compare the accuracy of four imputation methods for missing laboratory data and to compare the effect of the imputed values on the accuracy of predictive models.

Methods: Non-missing laboratory data was randomly removed with varying frequencies from two large data sets, and we then compared the ability of four methods – missForest, mean imputation, nearest neighbor imputation, and multivariate imputation by chained equations (MICE) – to impute the simulated missing data. We characterized the accuracy of the imputation and the effect of the imputation on predictive ability.

Results: MissForest had the least imputation error for both continuous and categorical variables at each frequency of missingness, and it had the smallest prediction difference when models used imputed laboratory values. In both data sets, MICE had the second least imputation error and prediction difference, followed by nearest neighbor imputation and mean imputation.

Conclusion: MissForest is a highly accurate method of imputation for missing laboratory data and outperforms other common imputation techniques in terms of imputation error and maintaining predictive ability with imputed values.

Introduction

“ You can have data without information, but you cannot have information without data ”

- Daniel Keys Moran

Missing data is a nearly ubiquitous problem when conducting research, particularly when using large data sets. Missing data can occur in the form of random or non-random patterns. Non-random missing data can introduce systematic error and make the study population less representative of the general population. Although random missing data does not introduce systematic error, it leads to significant loss in statistical power and predictive ability. Missing data are rarely completely at random and must be carefully managed. Multi-analyte Assays with Algorithmic Analyses (MAAAs) are a relatively new approach to leveraging value from laboratory data to predict clinical outcomes. Several of these are now available with implemented CPT codes (i.e. FibroSure, Risk of Ovarian Malignancy (ROMA), PreDx Diabetes Risk Score), but it is not known how robust MAAA models are when individual laboratory data points are missing.

Recent developments in machine learning have used laboratory data to build MAAA models to predict health care outcomes [1-3] - These models can be sensitive to missing laboratory data, which may result from hemolyzed samples, clumped platelets, or other uncommon sample or processing problems. It is not known whether modern imputation methods can robustly address the problem of missing data, and whether predictive models will remain accurate when imputed values are used.

Multiple methods have been developed in order to deal with missing data, including single imputation, multiple imputation, multiple imputation by chained equations (MICE)[3] , nearest neighbor estimation [4], and missForest[5] . Although these have been shown to be effective with other types of missing data, little data exists regarding the absolute or comparative effectiveness of these methods in accurately imputing missing laboratory data in predictive models. The aims of our study were: (1) to compare the accuracy of four different imputation methods for missing laboratory data in two large

1
2 data sets; and (2) to compare the effect of imputed values from each method on the accuracy of
3
4 predictive models based on these data sets.
5
6
7

8 9 **Methods**

10 **The University of Michigan (UM) Cirrhosis Cohort and Predictive Model for Hepatocellular**

11 **Carcinoma**

12
13
14
15
16 Between January 2004 and September 2006, consecutive patients with cirrhosis but no detectable HCC
17
18 were prospectively identified and entered into a surveillance program using ultrasound and alpha
19
20 fetoprotein (AFP), as has been previously described in greater detail [6]. Patients were enrolled if they
21
22 had Child-Pugh class A or B cirrhosis and absence of known HCC at the time of initial evaluation.
23
24

25
26 Patients diagnosed with HCC within the first six months of enrollment (prevalent cases) were excluded.

27
28 Other exclusion criteria included clinical evidence of significant hepatic decompensation (refractory
29
30 ascites, grade 3-4 encephalopathy, active variceal bleeding, or hepatorenal syndrome), co-morbid
31
32 medical conditions with a life expectancy of less than one year, prior solid organ transplant, and a
33
34 known extrahepatic primary tumor. Patients were followed until the time of HCC diagnosis, liver
35
36 transplantation, death, or until the study was terminated on July 31, 2010.
37
38

39
40 The following demographic and clinical data were collected at the time of enrollment: age, gender,
41
42 race, body mass index (BMI), past medical history, lifetime alcohol use, and lifetime tobacco use. Data
43
44 regarding their liver disease included the underlying etiology and presence of ascites, encephalopathy,
45
46 or esophageal varices. Laboratory data of interest at the time of enrollment included: platelet count,
47
48 aspartate aminotransferase (AST), alanine aminotransferase (ALT), alkaline phosphatase, bilirubin,
49
50 albumin, international normalized ratio (INR), and AFP. This data set was used as the basis of a
51
52 published predictive model to identify patients with hepatocellular carcinoma with a c statistic of 0.70
53
54
55
56
57 [2].
58
59
60

The UM Inflammatory Bowel Disease Cohort and Predictive Model for Thiopurine Clinical Response.

The study sample included all patients who had thiopurine metabolite analysis, CBC, and a comprehensive chemistry panel drawn within a 24-hour period at the University of Michigan between May 1, 2004, and August 31, 2006 and is described in greater detail in the manuscript [1]. This study was approved by the University of Michigan Medical Institutional Review Board with a waiver of explicit consent from the subjects. The patient sample included 774 cases, in a total of 346 individuals. For the analysis of the outcome of clinical response to thiopurines, 5 exclusion criteria were applied: exclusion of patients who did not have IBD, exclusion of patients who had not started on thiopurines at the time when metabolites were measured, exclusion of patients on biologic anti-tumor necrosis factor therapy, exclusion of patients without documentation of their clinical status at the time of laboratory measurement, and exclusion of patients who had an infection that confounded assessment of clinical response. This data set was used as the basis of a predictive model to identify patients with clinical response to thiopurine immune suppressant medication with a c statistic of 0.86[1].

Description of Imputation Techniques

We compared missForest with three other commonly used imputation methods that can handle both continuous and categorical variables, namely, mean Imputation, Nearest Neighbor Imputation, and Multivariate Imputation by Chained Equations (MICE). We briefly introduce these methods below.

The recently proposed missForest method makes use of highly flexible and versatile random forest models [7, 8] to achieve missing value imputation. It creates a random forest model for each variable using the rest of the variables in the data set and uses that to predict the missing values for that variable. This is done in a cyclic fashion for all variables and the entire process is iteratively repeated until a stopping criterion is attained. The advantages of using the random forest model are that it can handle

1
2 both continuous and categorical responses, requires very little tuning, and provides an internally cross-
3 validated error estimate. This was implemented via the 'missForest' package available in R.
4

5
6 Nearest Neighbor algorithms were originally proposed in the supervised pattern recognition literature.
7
8 Troyanskaya et al. proposed an imputation method based on nearest neighbor search [4]. The basic idea
9 is to compute a distance measure between each pair of observations based on the non-missing
10 variables. Then the k-nearest observations that have non-missing values for that particular variable are
11 used to impute a missing value via a weighted mean of the neighboring values. In order to
12 accommodate both continuous and categorical variables the Gower distance is used [9]. For the
13 categorical variables we imputed the missing values by weighted mode instead of a weighted mean as
14 used for continuous variables. Cross-validation error measures are used to select the optimal number of
15 nearest neighbors denoted by k. The function 'kNN' in R package 'VIM' was used to implement this
16 method.
17
18
19
20
21
22
23
24
25
26
27
28
29

30 Mean imputation is one of the most naïve and easiest methods for imputing missing values. The mean
31 (for continuous variables) or mode (for categorical variables) of the non-missing values of each
32 variable were used to impute the missing values. This method does not take advantage of any
33 correlation among the variables and therefore can perform rather poorly when such correlations are
34 present.
35
36
37
38
39
40
41

42 MICE was proposed by Van Buuren et al. [3]. It requires the user to specify a conditional model for
43 each variable, using the other variables as predictors. By default we used a linear regression model for
44 continuous variables, a logistic regression model for binary variables, and a polytomous logistic
45 regression for categorical variables with more than two levels. The algorithm works by iteratively
46 imputing the missing values based on the fitted conditional models until a stopping criterion is
47 satisfied. In that way it is very similar to the missForest algorithm; the main difference being that
48 missForest uses more flexible decision trees for each conditional model. We implemented this in R
49 using the package 'mice'[10].
50
51
52
53
54
55
56
57
58
59
60

Statistical Analysis

We used two separate studies to perform the comparison between the methods of imputation described in the previous section. We describe the studies, implementation details, and our results below. The structure of the statistical analysis is the same for both studies. We start with a published predictive model built with the training data set. The test set refers to observations that were not part of the training set; these were solely used for assessing the performance of the model. The test sets did not have any missing values, so we randomly removed a proportion of values to simulate data missing at random. We then imputed the missing values by the four previously discussed methods, and the imputed laboratory results were compared to the actual values that were removed from the data set. We then used the imputed data to make clinical outcome predictions with the published models, and the results were compared with the predictions made using the complete test data with no missing data. We use the average relative error (for continuous variables) and misclassification error (for categorical variables) to assess the imputation performance. To quantify the effect of the imputation on predictive models we compare the predicted classes from non-missing test data to predicted classes from imputed test data and compute the misclassification error. This is important because if a particular variable had very little influence on the predictive model, then larger imputation errors are tolerable, resulting in negligible loss of prediction accuracy. On the other hand, small imputation errors in very important variables might lead to significantly different predicted class, which is of greater clinical concern. We varied the frequency of missing values to change the difficulty of the imputation problem. We report the average results over multiple random runs. We found that the Nearest Neighbor results are quite robust to the choice of number of nearest neighbors (k) if k is moderately large, therefore we fixed the number of nearest neighbors at 5 in both studies.

Results

Cirrhosis Cohort and HCC model

This study evaluated the effect of imputation on a published predictive model for HCC based on 21 predictor variables that included demographic, clinical and laboratory values using random forest modeling [2]. The random forest model was developed on a data set of 446 patients collected at University of Michigan (UM cohort). It proved to be more accurate than traditional logistic regression models. Of the 21 variables, 10 of them were categorical in nature while 11 were continuous. We used the first 200 observations from the publicly available data set from the HALT-C trial as our test set and randomly replaced 10%, 20%, or 30% of the observations with missing values. The process was repeated for 30 replications and we report the average results.

The accuracy of the four imputation methods is compared in Figure 1. The vertical axis plots the percentage relative error for continuous variables and percentage misclassification error for categorical variables, while the horizontal axis groups the results according to the proportion of missing values. Each boxplot represent the error measure over 30 random replications. As expected, the imputation error increases on an average as we increase the proportion of missing values in the test data but the variation tends to reduce slightly which is due to averaging over many more missing observations. MissForest has the least imputation error for both continuous and categorical variables at each level of missing proportion, followed by MICE, NN, and mean imputation of continuous laboratory values. MICE and NN have similar imputation accuracy for categorical variables.

In Figure 2 the vertical axis plots the error measure for imputation on predictive model, obtained by comparing the predicted classes (Low Risk/High Risk) for each test observation with no-missing values against the predicted class after imputing the artificial missing values. Therefore an error measure of 5 on the vertical axis implies that 5% of the test observations had their predicted classes wrongly switched (either low risk => high risk or high risk => low risk) due to the imputation. As above, each

1
2 boxplot reflects the results of 30 random runs. It is clear from the figure that missForest performs the
3
4 best with NN and MICE following closely. The gap increases as we increase the proportion of missing
5
6 observations making the problem harder.
7
8
9

10 11 **Inflammatory Bowel Disease Cohort and Thiopurine Clinical Response Model**

12
13
14 Waljee et al. showed that random forest models using laboratory values outperform 6-thioguanine (6-
15
16 TGN) metabolite tests as well as traditional logistic regression models in predicting clinical response to
17
18 thiopurines [1]. The analysis was carried out on a data set collected at University of Michigan that
19
20 included 395 patients. 26 variables, which included 25 laboratory values and age, were used to predict
21
22 the clinical response of each patient. In this study all of the variables were continuous in nature. To
23
24 create a separate test set, we split the data set randomly into a training set consisting of 250
25
26 observations and a test set of 145 observations, using stratified sampling to keep the ratio of clinical
27
28 responder to non-responders fixed. We then introduced random missing values into the test set as
29
30 before and performed the same comparative study of the four imputation methods. The whole process
31
32 was replicated 30 times to obtain stable results. Below we summarize our findings via boxplots.
33
34
35

36
37 Again in this study, missForest beats its competitors in both imputation accuracy and the effect of
38
39 imputed values on the accuracy of clinical predictions. The trends remain the same for imputation error
40
41 with MICE coming out second best followed by NN and mean imputation. For predictive accuracy we
42
43 find the relative order becomes missForest > MICE > mean imputation > NN. This also shows that best
44
45 method with respect to imputation error need not be the best when we consider the effect of imputation
46
47 on predictive models. The performance gap between missForest and MICE is considerably lower than
48
49 in the previous study. This might be explained by the fact that in the thiopurine study, both the training
50
51 and test sets came from the same cohort, as we generated the training and test sets by random splits,
52
53 while in HCC study the training and test sets were completely different cohorts leading to an extra
54
55 degree of variation.
56
57
58
59
60

Discussion

We have performed an extensive simulation study using two actual datasets and two published predictive models to compare the performance of four methods of missing value imputation. We compared four popular methods namely, missForest, Nearest Neighbor, MICE and mean imputation in two studies simulating data missing at random. We found that missForest methods consistently produced the lowest imputation error and had the smallest prediction difference when models used imputed laboratory values. In addition, the ready availability of the freeware R package makes missForest a very convenient solution for any practical missing value problems. The main limitations of missForest as a solution to missing laboratory data for predictive modeling applications are: a requirement for skilled R programming for implementation, and slightly more demanding computational needs, compared to NN or MICE methods, with minimal increases in processing time (generally in the range of 10-20 sec)

The small absolute changes in predictions with these models, despite 10-30% missing laboratory data, speak to the robustness of these multi-analyte assays with algorithmic analysis (MAAAs). MAAAs are currently a hot topic, and several have been released with CPT codes in 2012. One example is the HCV FibroSure (LabCorp, code 0001M) which uses ALT, alpha 2 macroglobulin, apolipoprotein A1, total bilirubin, GGT, and haptoglobin to estimate fibrosis and necroinflammatory activity in the liver in patients with hepatitis C. With increasing complexity of these models, and increasing numbers of analytes, the risk of missing values increases, and methods to cope with missing values and preserve the accuracy of the model are needed. missForest appears to be a robust and accurate approach to the issue of missing laboratory values when used in these two MAAAs.

Contributorship:

Ashin Mukherjee – study concept and design, statistical analysis and interpretation of the data, drafting of the manuscript, critical revision of the manuscript

Akbar Waljee - study concept and design, acquisition of the data, statistical analysis and interpretation of the data, drafting of the manuscript, critical revision of the manuscript and study supervision.

Amit Singal - acquisition of data, drafting of the manuscript and critical revision of the manuscript.

Yiwei Zhang - statistical analysis and interpretation of the data and critical revision of the manuscript.

Jeffrey Warren - study concept and design, critical revision of the manuscript.

Ulysses Balis - study concept and design, critical revision of the manuscript.

Jorge Marrero - study concept and design, critical revision of the manuscript.

Ji Zhu - statistical analysis and interpretation of the data and critical revision of the manuscript.

Peter Higgins - study concept and design, acquisition of the data, statistical analysis and interpretation of the data, critical revision of the manuscript and study supervision.

Competing Interests: None**Ethics approval:** University of Michigan IRB

Figure Legends

FIGURE 1. Imputation error comparison for categorical and continuous variables for four competing imputation methods at three levels of the proportion of missing values for the HCC study.

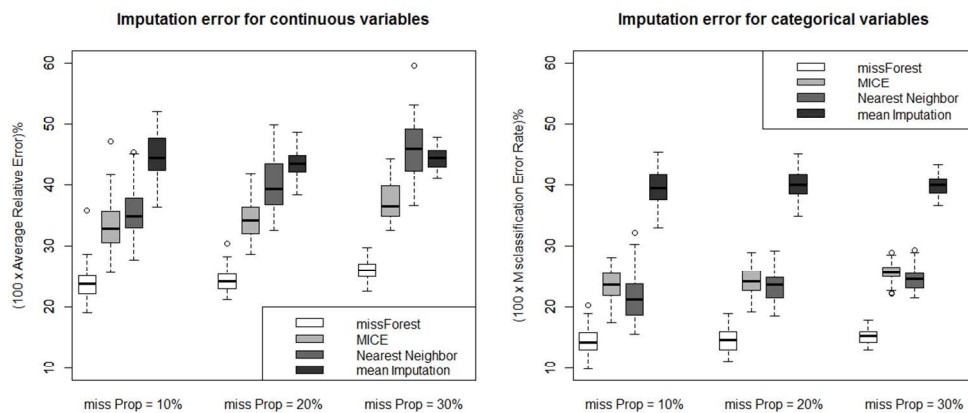
FIGURE 2. Percentage of wrongly predicted observations after missing value imputation by the four competing methods at three levels of missing value proportions in the test data.

FIGURE 3. Left: Imputation error comparison, Right: Misclassification error in predicted classes due to imputation for four competing imputation methods at three levels of the proportion of missing values for the Thiopurine study.

References

- 1 Waljee AK, Joyce JC, Wang S, et al. Algorithms outperform metabolite tests in predicting response of patients with inflammatory bowel disease to thiopurines. *Clinical gastroenterology and hepatology : the official clinical practice journal of the American Gastroenterological Association*. 2010 Feb;**8**(2):143-50.
- 2 Singal AG, Waljee AK, Mukherjee A, Higgins PD, Zhu J, Marrero JA. Machine Learning Algorithms Outperform Conventional Regression Models in Identifying Risk Factors for Hepatocellular Carcinoma in Patients With Cirrhosis. *Gastroenterology*. 2012 May;**142**(5):S984-S.
- 3 van Buuren S, Boshuizen HC, Knook DL. Multiple imputation of missing blood pressure covariates in survival analysis. *Stat Med*. 1999 Mar 30;**18**(6):681-94.
- 4 Troyanskaya O, Cantor M, Sherlock G, et al. Missing value estimation methods for DNA microarrays. *Bioinformatics*. 2001 Jun;**17**(6):520-5.
- 5 Stekhoven DJ, Buhlmann P. MissForest--non-parametric missing value imputation for mixed-type data. *Bioinformatics*. 2012 Jan 1;**28**(1):112-8.
- 6 Singal AG, Conjeevaram HS, Volk ML, et al. Effectiveness of hepatocellular carcinoma surveillance in patients with cirrhosis. *Cancer epidemiology, biomarkers & prevention : a publication of the American Association for Cancer Research, cosponsored by the American Society of Preventive Oncology*. 2012 May;**21**(5):793-9.
- 7 Liaw A, Wiener M. Classification and Regression by randomForest. 2002;**2**(3):18-22.
- 8 Breiman L. Random forests. 2001;**45**(1):5-32.
- 9 Gower JC. A General Coefficient of Similarity and Some of Its Properties. *Biometrics*. 1971;**27**(4):857-71.
- 10 Buuren van S, Groothuis-Oudshoorn K. mice: Multivariate Imputation by Chained Equations in R. *Journal of Statistical Software*. 2011;**45**(3).

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
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

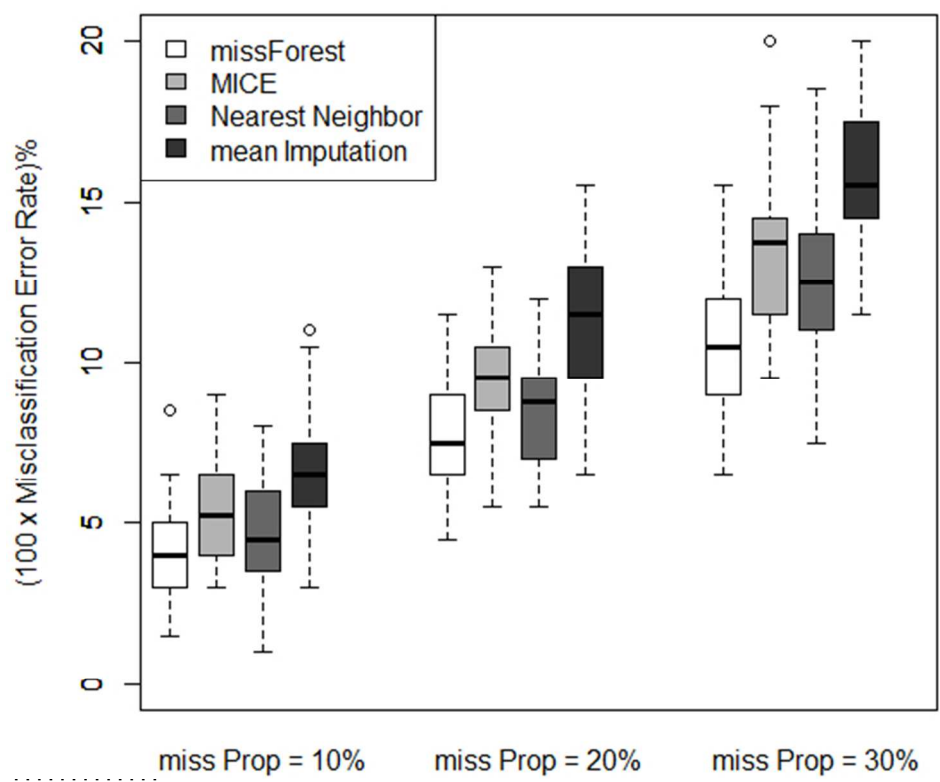


430x204mm (72 x 72 DPI)

er review only

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
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Misclassification Error due to Imputation

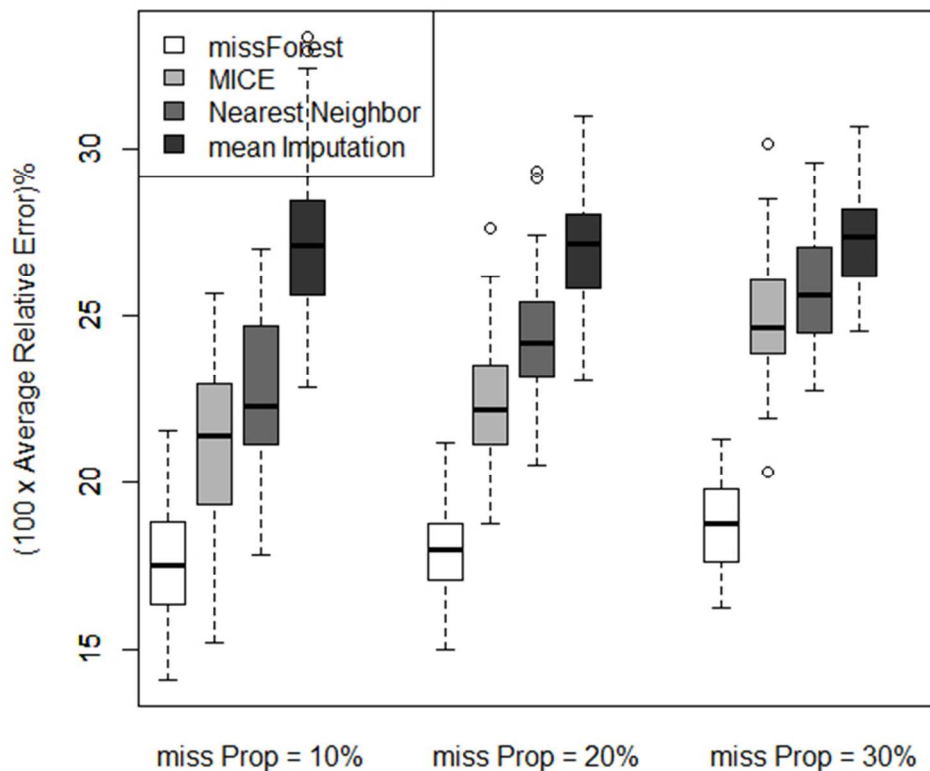


204x204mm (72 x 72 DPI)



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
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Imputation Error Comparison

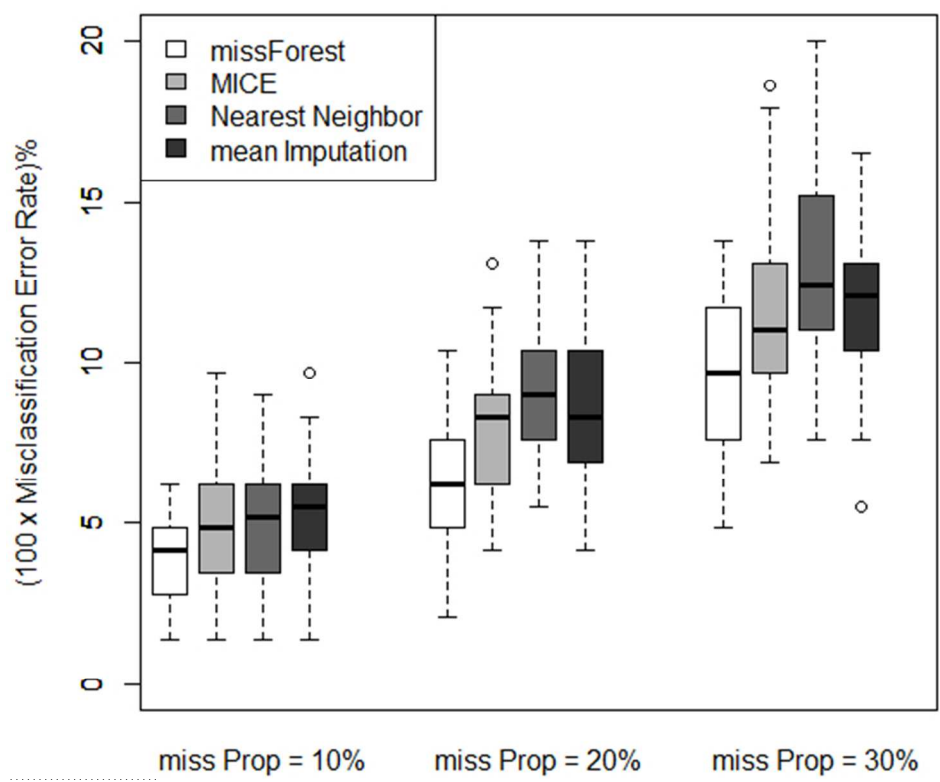


204x204mm (72 x 72 DPI)



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
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Misclassification Error due to Imputation



204x204mm (72 x 72 DPI)





Comparison of Modern Imputation Methods for Missing Laboratory Data in Medicine

Journal:	<i>BMJ Open</i>
Manuscript ID:	bmjopen-2013-002847.R1
Article Type:	Research
Date Submitted by the Author:	15-May-2013
Complete List of Authors:	Waljee, Akbar; VA Center for Clinical Management Research, Ann Arbor VA Medical Center; University of Michigan, Department of Internal Medicine Mukherjee, Ashin; University of Michigan, Department of Statistics Singal, Amit; UT Southwestern Medical Center, Department of Internal Medicine Zhang, Yiwei; University of Michigan, Department of Statistics Warren, Jeffrey; University of Michigan, Pathology Balis, Ulysses; University of Michigan, Pathology Marrero, Jorge; UT Southwestern Medical Center, Department of Internal Medicine Zhu, Ji; University of Michigan, Department of Statistics Higgins, Peter; University of Michigan, Department of Internal Medicine
Primary Subject Heading:	Gastroenterology and hepatology
Secondary Subject Heading:	Diagnostics, Gastroenterology and hepatology, Health informatics
Keywords:	Gastroduodenal disease < GASTROENTEROLOGY, Information technology < BIOTECHNOLOGY & BIOINFORMATICS, Inflammatory bowel disease < GASTROENTEROLOGY

SCHOLARONE™
Manuscripts

Comparison of Modern Imputation Methods for Missing Laboratory Data in Medicine

Akbar K. Waljee MD MS^{2,3}, Ashin Mukherjee MS¹, Amit G. Singal MD MS^{4,5}, Yiwei Zhang MS¹,
Jeffrey Warren MD⁶, Ulysses Balis MD⁶, Jorge Marrero⁴, Ji Zhu PhD¹, Peter DR Higgins MD PhD²

¹Department of Statistics, University of Michigan, Ann Arbor, MI

²Department of Internal Medicine, University of Michigan, Ann Arbor, MI

³Veterans Affairs Center for Clinical Management Research, Ann Arbor, MI

⁴Department of Internal Medicine, UT Southwestern Medical Center, Dallas, TX

⁵Department of Clinical Sciences, UT Southwestern, Dallas, TX

⁶Department of Pathology, University of Michigan, Ann Arbor, MI

Correspondence:

Akbar K. Waljee, M.D., M.Sc.

VA Center for Clinical Management Research,

Ann Arbor VA Medical Center,

2215 Fuller Road, IIID,

Ann Arbor, MI 48105

e-mail: awaljee@med.umich.edu

Financial Support: Dr. Waljee's research is funded by a VA HSR&D CDA-2 Career Development Award 1IK2HX000775. Dr. Singal's research is funded by an ACG Junior Faculty Development Award and grant number KL2 RR024983-05. Dr. Higgins' research is supported by NIH R01 GM097117. The content is solely the responsible of the authors and does not necessarily represent the official views of

1 UT-STAR, the University of Texas Southwestern Medical Center at Dallas and its affiliated academic
2
3 and health care centers, the National Center for Advancing Translational Sciences, or the National
4
5 Institutes of Health.
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
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Acknowledgements:

Ashin Mukherjee and Akbar Waljee contributed equally to this work and should be considered co-first authors.

The content is solely the responsibility of the authors and does not necessarily reflect the official views of the Department of Veterans Affairs, the National Center for Research Resources, or the National Institutes of Health.

Conflicts of Interest: The authors disclose no conflicts.

Abstract Word Count: 211

Manuscript Word Count: 2843

Article Summary:

1) Article Focus

- Multi-analyte Assays with Algorithmic Analyses (MAAAs) are a relatively new approach to leveraging value from laboratory data to predict clinical outcomes. It is not known how robust MAAA models are when individual laboratory data points are missing.
- Recent developments in machine learning have used laboratory data to build MAAA models to predict health care outcomes - These models can be sensitive to missing laboratory data. It is not known whether modern imputation methods can robustly address the problem of missing data, and whether predictive models will remain accurate when imputed values are used.
- Multiple methods have been developed in order to deal with missing data, including single imputation, multiple imputation, multivariate imputation by chained equations (MICE), nearest neighbor estimation, and missForest. Although these have been shown to be effective with other types of missing data, little data exists regarding the absolute or comparative effectiveness of these methods in accurately imputing missing completely at random laboratory data in predictive models. The aims of our study were: (1) to compare the accuracy of four different imputation methods for missing laboratory data in two large data sets; and (2) to compare the effect of imputed values from each method on the accuracy of predictive models based on these data sets.

2) Key Message

- We found that missForest methods consistently produced the lowest imputation error and had the smallest prediction difference when models used imputed laboratory values.
- The small absolute changes in predictions with these models, despite 10-30% missing laboratory data, speak to the robustness of these multi-analyte assays with algorithmic analysis (MAAAs).
- With increasing complexity of these models, and increasing numbers of analytes, the risk of missing values increases, and methods to cope with missing values and preserve the accuracy of the model are needed. missForest appears to be a robust and accurate approach to the issue of missing laboratory values when used in these two MAAAs.

3) Limitations and Strengths

- The main limitations of missForest as a solution to missing laboratory data for predictive modeling applications are: a requirement for skilled R programming for implementation, and slightly more demanding computational needs, compared to NN or MICE methods.
- The simulations in this manuscript use data missing at random. The results presented here may not be generalizable to situations in which laboratory values are missing in a bias, non-random way.
- The strength of this is that missForest methods consistently produced the lowest imputation error and had the smallest prediction difference when models used imputed laboratory values and that it is a readily available freeware R package making it a very convenient solution for any practical missing value problems.

Abstract

Background: Missing laboratory data is a common issue, but the optimal method of imputation of

1 missing values has not been determined. The aims of our study were to compare the accuracy of four
2
3
4 imputation methods for missing completely at random laboratory data and to compare the effect of the
5
6 imputed values on the accuracy of two clinical predictive models.
7

8 **Methods:** Non-missing laboratory data was randomly removed with varying frequencies from two
9
10 large data sets, and we then compared the ability of four methods – missForest, mean imputation,
11
12 nearest neighbor imputation, and multivariate imputation by chained equations (MICE) – to impute the
13
14 simulated missing data. We characterized the accuracy of the imputation and the effect of the
15
16 imputation on predictive ability in two large datasets.
17
18

19
20 **Results:** MissForest had the least imputation error for both continuous and categorical variables at each
21
22 frequency of missingness, and it had the smallest prediction difference when models used imputed
23
24 laboratory values. In both data sets, MICE had the second least imputation error and prediction
25
26 difference, followed by nearest neighbor imputation and mean imputation.
27
28

29
30 **Conclusion:** MissForest is a highly accurate method of imputation for missing laboratory data and
31
32 outperforms other common imputation techniques in terms of imputation error and maintaining
33
34 predictive ability with imputed values in two clinical predicative models.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Introduction

“ You can have data without information, but you cannot have information without data”

- Daniel Keys Moran

Missing data is a nearly ubiquitous problem when conducting research, particularly when using large data sets. Missing data can occur in the form of random or non-random patterns. Non-random missing data can introduce systematic error and make the study population less representative of the general population. Although random missing data does not introduce systematic error, it leads to significant loss in statistical power and predictive ability. Missing data are rarely completely at random and must be carefully managed. Multi-analyte Assays with Algorithmic Analyses (MAAAs) are a relatively new approach to leveraging value from laboratory data to predict clinical outcomes. Several of these are now available with implemented CPT codes (i.e. FibroSure, Risk of Ovarian Malignancy (ROMA), PreDx Diabetes Risk Score), but it is not known how robust MAAA models are when individual laboratory data points are missing.

Recent developments in machine learning have used laboratory data to build MAAA models to predict health care outcomes [1-3] - These models can be sensitive to missing completely at random laboratory data, which may result from hemolyzed samples, clumped platelets, or other uncommon sample or processing problems. It is not known whether modern imputation methods can robustly address the problem of missing data, and whether predictive models will remain accurate when imputed values are used.

Multiple methods have been developed in order to deal with missing data, including single imputation, multiple imputation, multivariate imputation by chained equations (MICE)[3] , nearest neighbor estimation [4], and missForest[5] . Although these have been shown to be effective with other types of missing data, little data exists regarding the absolute or comparative effectiveness of these methods in accurately imputing missing laboratory data in predictive models. The aims of our study were: (1) to compare the accuracy of four different imputation methods for missing completely at random

1 laboratory data in two large data sets; and (2) to compare the effect of imputed values from each
2 method on the accuracy of predictive models based on these data sets.
3
4
5
6
7

8 **Methods**

9
10 The University of Michigan (UM) Cirrhosis Cohort and Predictive Model for Hepatocellular
11 Carcinoma
12

13
14 Between January 2004 and September 2006, consecutive patients with cirrhosis but no detectable HCC
15 were prospectively identified and entered into a surveillance program using ultrasound and alpha
16 fetoprotein (AFP), as has been previously described in greater detail [6]. Patients were enrolled if they
17 had Child-Pugh class A or B cirrhosis and absence of known HCC at the time of initial evaluation.
18 Patients diagnosed with HCC within the first six months of enrollment (prevalent cases) were excluded.
19 Other exclusion criteria included clinical evidence of significant hepatic decompensation (refractory
20 ascites, grade 3-4 encephalopathy, active variceal bleeding, or hepatorenal syndrome), co-morbid
21 medical conditions with a life expectancy of less than one year, prior solid organ transplant, and a
22 known extrahepatic primary tumor. Patients were followed until the time of HCC diagnosis, liver
23 transplantation, death, or until the study was terminated on July 31, 2010.
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

39 The following demographic and clinical data were collected at the time of enrollment: age, gender,
40 race, body mass index (BMI), past medical history, lifetime alcohol use, and lifetime tobacco use. Data
41 regarding their liver disease included the underlying etiology and presence of ascites, encephalopathy,
42 or esophageal varices. Laboratory data of interest at the time of enrollment included: platelet count,
43 aspartate aminotransferase (AST), alanine aminotransferase (ALT), alkaline phosphatase, bilirubin,
44 albumin, international normalized ratio (INR), and AFP. This data set was used as the basis of a
45 published predictive model to identify patients with hepatocellular carcinoma with a c statistic of 0.70
46
47
48
49
50
51
52
53
54
55
56 [2].
57
58
59
60

1 The UM Inflammatory Bowel Disease Cohort and Predictive Model for Thiopurine Clinical Response.
2
3 The study sample included all patients who had thiopurine metabolite analysis, CBC, and a
4 comprehensive chemistry panel drawn within a 24-hour period at the University of Michigan between
5
6 May 1, 2004, and August 31, 2006 and is described in greater detail in the manuscript [1]. This study
7
8 was approved by the University of Michigan Medical Institutional Review Board with a waiver of
9
10 explicit consent from the subjects. The patient sample included 774 cases, in a total of 346 individuals.
11
12 For the analysis of the outcome of clinical response to thiopurines, 5 exclusion criteria were applied:
13
14 exclusion of patients who did not have IBD, exclusion of patients who had not started on thiopurines at
15
16 the time when metabolites were measured, exclusion of patients on biologic anti-tumor necrosis factor
17
18 therapy, exclusion of patients without documentation of their clinical status at the time of laboratory
19
20 measurement, and exclusion of patients who had an infection that confounded assessment of clinical
21
22 response. This data set was used as the basis of a predictive model to identify patients with clinical
23
24 response to thiopurine immune suppressant medication with a c statistic of 0.86[1].
25
26
27
28
29
30
31
32
33

34 **Description of Imputation Techniques**

35
36 We compared missForest with three other commonly used imputation methods that can handle both
37
38 continuous and categorical variables, namely, mean Imputation, Nearest Neighbor Imputation, and
39
40 Multivariate Imputation by Chained Equations (MICE). We briefly introduce these methods below.
41
42 The recently proposed missForest method makes use of highly flexible and versatile random forest
43
44 models [7, 8] to achieve missing value imputation. It creates a random forest model for each variable
45
46 using the rest of the variables in the data set and uses that to predict the missing values for that variable.
47
48 This is done in a cyclic fashion for all variables and the entire process is iteratively repeated until a
49
50 stopping criterion is attained. The advantages of using the random forest model are that it can handle
51
52 both continuous and categorical responses, requires very little tuning, and provides an internally cross-
53
54 validated error estimate. This was implemented via the 'missForest' package available in R.
55
56
57
58
59
60

1 Nearest Neighbor algorithms were originally proposed in the supervised pattern recognition literature.
2
3 Troyanskaya et al. proposed an imputation method based on nearest neighbor search [4]. The basic idea
4 is to compute a distance measure between each pair of observations based on the non-missing
5 variables. Then the k-nearest observations that have non-missing values for that particular variable are
6 used to impute a missing value via a weighted mean of the neighboring values. In order to
7 accommodate both continuous and categorical variables the Gower distance is used [9]. For the
8 categorical variables we imputed the missing values by weighted mode instead of a weighted mean as
9 used for continuous variables. Cross-validation error measures are used to select the optimal number of
10 nearest neighbors denoted by k. The function 'kNN' in R package 'VIM' was used to implement this
11 method.
12
13
14
15
16
17
18
19
20
21
22
23

24 Mean imputation is one of the most naïve and easiest methods for imputing missing values. The mean
25 (for continuous variables) or mode (for categorical variables) of the non-missing values of each
26 variable were used to impute the missing values. This method does not take advantage of any
27 correlation among the variables and therefore can perform rather poorly when such correlations are
28 present.
29
30
31
32
33
34
35

36 MICE was proposed by Van Buuren et al. [3]. It requires the user to specify a conditional model for
37 each variable, using the other variables as predictors. By default we used a linear regression model for
38 continuous variables, a logistic regression model for binary variables, and a polytomous logistic
39 regression for categorical variables with more than two levels. The algorithm works by iteratively
40 imputing the missing values based on the fitted conditional models until a stopping criterion is
41 satisfied. In that way it is very similar to the missForest algorithm; the main difference being that
42 missForest uses more flexible decision trees for each conditional model. We implemented this in R
43 using the package 'mice'[10].
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Statistical Analysis

We used two separate studies to perform the comparison between the methods of imputation described in the previous section. We describe the studies, implementation details, and our results below. The structure of the statistical analysis is the same for both studies. We start with a published predictive model built with the training data set. The test set refers to observations that were not part of the training set; these were solely used for assessing the performance of the model. The test sets did not have any missing values, so we randomly removed a proportion of values to simulate data missing completely at random. We then imputed the missing values by the four previously discussed methods, and the imputed laboratory results were compared to the actual values that were removed from the data set. We then used the imputed data to make clinical outcome predictions with the published models, and the results were compared with the predictions made using the complete test data with no missing data. The prediction models were created using both logistic models and random forest models as some might argue that using random forest imputation methods might favor the random forest prediction models. We use the average relative error (for continuous variables) and misclassification error (for categorical variables) to assess the imputation performance for both the logistic and random forest models. To quantify the effect of the imputation on predictive models we compare the predicted classes from non-missing test data to predicted classes from imputed test data and compute the misclassification error for both the logistic and random forest models. This is important because if a particular variable had very little influence on the predictive model, then larger imputation errors are tolerable, resulting in negligible loss of prediction accuracy. On the other hand, small imputation errors in very important variables might lead to significantly different predicted class, which is of greater clinical concern. We varied the frequency of missing values to change the difficulty of the imputation problem. We report the average results over multiple random runs. We found that the Nearest Neighbor results are quite robust to the choice of number of nearest neighbors (k) if k is moderately large, therefore we fixed the number of nearest neighbors at 5 in both studies.

Results

Cirrhosis Cohort and HCC model

This study evaluated the effect of imputation on a published predictive model for HCC based on 21 predictor variables that included demographic, clinical and laboratory values using random forest modeling [2]. The random forest model was developed on a data set of 446 patients collected at University of Michigan (UM cohort). It proved to be more accurate than traditional logistic regression models. Of the 21 variables, 10 of them were categorical in nature while 11 were continuous. We used the first 200 observations from the publicly available data set from the HALT-C trial as our test set and randomly replaced 10%, 20%, or 30% of the observations with missing values. The process was repeated for 1000 replications and we report the average results.

The accuracy of the four imputation methods for both continuous and categorical variables are compared in Figure 1 for the cirrhosis cohort and HCC model. Figure 1A and 1B represents the logistic model and Figure 1C and 1D reflect the random forest prediction model. The vertical axis plots the percentage relative error for continuous variables and percentage misclassification error for categorical variables, while the horizontal axis groups the results according to the proportion of missing values. Each boxplot represent the error measure over 1000 random replications. As expected, the imputation error increases on an average as we increase the proportion of missing values in the test data but the variation tends to reduce slightly which is due to averaging over many more missing observations. MissForest has the least imputation error for both continuous and categorical variables at each level of missing proportion, followed by MICE, NN, and mean imputation of continuous laboratory values. MICE and NN have similar imputation accuracy for categorical variables. MissForest works well using both logistic and random forest prediction models.

In Figure 2A and 2B, which represents the logistic and random forest prediction models respectively, the vertical axis plots the error measure for imputation on predictive model, obtained by comparing the predicted classes (Low Risk/High Risk) for each test observation with no-missing values against the

1 predicted class after imputing the artificial missing values. Therefore an error measure of 5 on the
2
3 vertical axis implies that 5% of the test observations had their predicted classes wrongly switched
4
5 (either low risk => high risk or high risk => low risk) due to the imputation. As above, each boxplot
6
7 reflects the results of 1000 random runs. It is clear from the figure that missForest performs the best
8
9 with NN and MICE following closely. The gap increases as we increase the proportion of missing
10
11 observations making the problem harder.
12
13
14
15
16
17

18 Inflammatory Bowel Disease Cohort and Thiopurine Clinical Response Model

19
20 Waljee et al. showed that random forest models using laboratory values outperform 6-thioguanine (6-
21
22 TGN) metabolite tests as well as traditional logistic regression models in predicting clinical response to
23
24 thiopurines [1]. The analysis was carried out on a data set collected at University of Michigan that
25
26 included 395 patients. 26 variables, which included 25 laboratory values and age, were used to predict
27
28 the clinical response of each patient. In this study all of the variables were continuous in nature. To
29
30 create a separate test set, we split the data set randomly into a training set consisting of 250
31
32 observations and a test set of 145 observations, using stratified sampling to keep the ratio of clinical
33
34 responder to non-responders fixed. We then introduced random missing values into the test set as
35
36 before and performed the same comparative study of the four imputation methods. The whole process
37
38 was replicated 1000 times to obtain stable results. Below we summarize our findings via boxplots.
39
40
41
42

43 Again in this study, missForest beats its competitors in both imputation accuracy and the effect of
44
45 imputed values on the accuracy of clinical predictions based on the logistic model (Figure 3A) and the
46
47 random forest model (Figure 3B). The trends remain the same for imputation error, Figure 4A and B
48
49 representing logistic models and random forest models respectively, with MICE coming out second
50
51 best followed by NN and mean imputation. For predictive accuracy we find the relative order becomes
52
53 missForest > MICE > mean imputation > NN. This also shows that best method with respect to
54
55 imputation error need not be the best when we consider the effect of imputation on predictive models.
56
57
58
59
60

1 The performance gap between missForest and MICE is considerably lower than in the previous study.
2
3 This might be explained by the fact that in the thiopurine study, both the training and test sets came
4
5 from the same cohort, as we generated the training and test sets by random splits, while in HCC study
6
7 the training and test sets were completely different cohorts leading to an extra degree of variation.
8
9

10 **Discussion**

11
12 We have performed an extensive simulation study using two clinical datasets and two published
13
14 predictive models to compare the performance of four methods of missing value imputation for missing
15
16 data completely at random. We included both local (randomForest) and global (logistic) modeling
17
18 approaches to avoid bias that might favor a local (MissForest) imputation approach. While the
19
20 superiority of MissForest for imputation of missing lab values this will not be generalizable to all
21
22 predictive models or datasets, this manuscript highlights the value of missForest to impute missing
23
24 data. We compared four popular methods namely, missForest, Nearest Neighbor, MICE and mean
25
26 imputation, in two studies simulating data missing completely at random. We found that these
27
28 simulation methods consistently produced the lowest imputation error and had the smallest prediction
29
30 difference when models used imputed laboratory values. In addition, the ready availability of the
31
32 freeware R package makes missForest and its simulations a very convenient solution for any practical
33
34 missing value problems. The main limitations of these simulations as a solution to missing laboratory
35
36 data for predictive modeling applications are: a requirement for skilled R programming for
37
38 implementation, and slightly more demanding computational needs, compared to NN or MICE
39
40 methods. An additional limitation in this study is that these simulations did not address the issue of data
41
42 missing for non-random reasons. There could be an association between the clinical outcome of interest
43
44 and the missingness of certain predictors. At this point, we cannot generalize these results to situation
45
46 in which data is missing for non-random reasons.
47
48
49
50
51
52
53
54

55 The small absolute changes in predictions with these models, despite 10-30% missing laboratory data,
56
57 speak to the robustness of these multi-analyte assays with algorithmic analysis (MAAAs). MAAAs are
58
59
60

1 currently a hot topic, and several have been released with CPT codes in 2012. One example is the HCV
2
3 FibroSure (LabCorp, code 0001M) which uses ALT, alpha 2 macroglobulin, apolipoprotein A1, total
4
5 bilirubin, GGT, and haptoglobin to estimate fibrosis and necroinflammatory activity in the liver in
6
7 patients with hepatitis C. With the increasing complexity of these models, and increasing numbers of
8
9 analytes, the risk of missing completely at random values increases, and methods to cope with missing
10
11 values and preserve the accuracy of the model are needed. missForest appears to be a robust and
12
13 accurate approach to the issue of missing laboratory values when used in these two MAAAs and may
14
15 be applicable to other datasets with missing completely at random datasets.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Contributions:

Ashin Mukherjee – study concept and design, statistical analysis and interpretation of the data, drafting of the manuscript, critical revision of the manuscript

Akbar Waljee - study concept and design, acquisition of the data, statistical analysis and interpretation of the data, drafting of the manuscript, critical revision of the manuscript and study supervision.

Amit Singal - acquisition of data, drafting of the manuscript and critical revision of the manuscript.

Yiwei Zhang - statistical analysis and interpretation of the data and critical revision of the manuscript.

Jeffrey Warren - study concept and design, critical revision of the manuscript.

Ulysses Balis - study concept and design, critical revision of the manuscript.

Jorge Marrero - study concept and design, critical revision of the manuscript.

Ji Zhu - statistical analysis and interpretation of the data and critical revision of the manuscript.

Peter Higgins - study concept and design, acquisition of the data, statistical analysis and interpretation of the data, critical revision of the manuscript and study supervision.

Competing Interests: None

Ethics approval: University of Michigan, IRB

Data Sharing: Agreement: A prediction model using the cross-sectional data from the hepatocellular carcinoma cohort had currently been submitted to the American Journal of Gastroenterology and we are have a revise and resubmit

Figure Legends

FIGURE 1. Imputation error comparison for categorical and continuous variables for four competing imputation methods at three levels of the proportion of missing values for the logistic prediction model (Figure 1A and 1B) and random forest prediction model (Figure 1C and 1D) in the HCC study

FIGURE 2. Percentage of wrongly predicted observations after missing value imputation by the four competing methods at three levels of missing value proportions in the test data for the logistic prediction model (Figure 2A) and the random forest prediction model (Figure 2B) in the HCC study.

FIGURE 3. Imputation error for four competing imputation methods at three levels of the proportion of missing values for the logistic prediction model (Figure 3A) and random forest prediction model (Figure 3B) in the Thiopurine response model

Figure 4. Percentage of wrongly predicted observations after missing value imputation by the four competing methods at three levels of missing value proportions in the test data for the logistic prediction model (Figure 4A) and the random forest prediction model (Figure 4B) in the Thiopurine response model.

References

- 1 Waljee AK, Joyce JC, Wang S, et al. Algorithms outperform metabolite tests in predicting response of patients with inflammatory bowel disease to thiopurines. *Clinical gastroenterology and hepatology : the official clinical practice journal of the American Gastroenterological Association*. 2010 Feb;**8**(2):143-50.
- 2 Singal AG, Waljee AK, Mukherjee A, et al. Machine Learning Algorithms Outperform Conventional Regression Models in Identifying Risk Factors for Hepatocellular Carcinoma in Patients With Cirrhosis. *Gastroenterology*. 2012 May;**142**(5):S984-S.
- 3 van Buuren S, Boshuizen HC, Knook DL. Multiple imputation of missing blood pressure covariates in survival analysis. *Stat Med*. 1999 Mar 30;**18**(6):681-94.
- 4 Troyanskaya O, Cantor M, Sherlock G, et al. Missing value estimation methods for DNA microarrays. *Bioinformatics*. 2001 Jun;**17**(6):520-5.
- 5 Stekhoven DJ, Buhlmann P. MissForest--non-parametric missing value imputation for mixed-type data. *Bioinformatics*. 2012 Jan 1;**28**(1):112-8.
- 6 Singal AG, Conjeevaram HS, Volk ML, et al. Effectiveness of hepatocellular carcinoma surveillance in patients with cirrhosis. *Cancer epidemiology, biomarkers & prevention : a publication of the American Association for Cancer Research, cosponsored by the American Society of Preventive Oncology*. 2012 May;**21**(5):793-9.
- 7 Liaw A, Wiener M. Classification and Regression by randomForest. 2002;**2**(3):18-22.
- 8 Breiman L. Random forests. 2001;**45**(1):5-32.
- 9 Gower JC. A General Coefficient of Similarity and Some of Its Properties. *Biometrics*. 1971;**27**(4):857-71.
- 10 Buuren van S, Groothuis-Oudshoorn K. mice: Multivariate Imputation by Chained Equations in R. *Journal of Statistical Software*. 2011;**45**(3).

Comparison of Modern Imputation Methods for Missing Laboratory Data in Medicine

Ashin Mukherjee MS¹, Akbar K. Waljee MD MS^{2,3}, Amit G. Singal MD MS^{4,5}, Yiwei Zhang MS¹,
Jeffrey Warren MD⁶, Ulysses Balis MD⁶, Jorge Marrero⁴, Ji Zhu PhD¹, Peter DR Higgins MD PhD²

¹Department of Statistics, University of Michigan, Ann Arbor, MI

²Department of Internal Medicine, University of Michigan, Ann Arbor, MI

³Veterans Affairs Center for Clinical Management Research, Ann Arbor, MI

⁴Department of Internal Medicine, UT Southwestern Medical Center, Dallas, TX

⁵Department of Clinical Sciences, UT Southwestern, Dallas, TX

⁶Department of Pathology, University of Michigan, Ann Arbor, MI

Correspondence:

Akbar K. Waljee, M.D., M.Sc.

VA Center for Clinical Management Research,

Ann Arbor VA Medical Center,

2215 Fuller Road, IID,

Ann Arbor, MI 48105

e-mail: awaljee@med.umich.edu

Financial Support: Dr. Waljee's research is funded by a VA HSR&D CDA-2 Career Development Award 1IK2HX000775. Dr. Singal's research is funded by an ACG Junior Faculty Development Award and grant number KL2 RR024983-05. Dr. Higgins' research is supported by NIH R01 GM097117. The content is solely the responsible of the authors and does not necessarily represent the official views of

1
2 UT-STAR, the University of Texas Southwestern Medical Center at Dallas and its affiliated academic
3
4 and health care centers, the National Center for Advancing Translational Sciences, or the National
5
6 Institutes of Health.
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
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Acknowledgements:

Ashin Mukherjee and Akbar Waljee contributed equally to this work and should be considered co-first authors.

The content is solely the responsibility of the authors and does not necessarily reflect the official views of the Department of Veterans Affairs, the National Center for Research Resources, or the National Institutes of Health.

Conflicts of Interest: The authors disclose no conflicts.

Abstract Word Count: 211

Manuscript Word Count: 2843

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
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Article Summary:

1) Article Focus

- Multi-analyte Assays with Algorithmic Analyses (MAAAs) are a relatively new approach to leveraging value from laboratory data to predict clinical outcomes. It is not known how robust MAAA models are when individual laboratory data points are missing.
- Recent developments in machine learning have used laboratory data to build MAAA models to predict health care outcomes - These models can be sensitive to missing laboratory data. It is not known whether modern imputation methods can robustly address the problem of missing data, and whether predictive models will remain accurate when imputed values are used.
- Multiple methods have been developed in order to deal with missing data, including single imputation, multiple imputation, **multivariate imputation** by chained equations (MICE), nearest neighbor estimation, and missForest. Although these have been shown to be effective with other types of missing data, little data exists regarding the absolute or comparative effectiveness of these methods in accurately imputing missing **completely at random** laboratory data in predictive models. The aims of our study were: (1) to compare the accuracy of four different imputation methods for missing laboratory data in two large data sets; and (2) to compare the effect of imputed values from each method on the accuracy of predictive models based on these data sets.

2) Key Message

- 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
- We found that missForest methods consistently produced the lowest imputation error and had the smallest prediction difference when models used imputed laboratory values.
 - The small absolute changes in predictions with these models, despite 10-30% missing laboratory data, speak to the robustness of these multi-analyte assays with algorithmic analysis (MAAAs).
 - With increasing complexity of these models, and increasing numbers of analytes, the risk of missing values increases, and methods to cope with missing values and preserve the accuracy of the model are needed. missForest appears to be a robust and accurate approach to the issue of missing laboratory values when used in these two MAAAs.

26 3) Limitations and Strengths

- 27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- The main limitations of missForest as a solution to missing laboratory data for predictive modeling applications are: a requirement for skilled R programming for implementation, and slightly more demanding computational needs, compared to NN or MICE methods.
 - The simulations in this manuscript use data missing at random. The results presented here may not be generalizable to situations in which laboratory values are missing in a bias, non-random way.
 - The strength of this is that missForest methods consistently produced the lowest imputation error and had the smallest prediction difference when models used imputed laboratory values and that it is a readily available freeware R package making it a very convenient solution for any practical missing value problems.

Abstract

Background: Missing laboratory data is a common issue, but the optimal method of imputation of missing values has not been determined. The aims of our study were to compare the accuracy of four imputation methods for missing **completely at random** laboratory data and to compare the effect of the imputed values on the accuracy of **two clinical** predictive models.

Methods: Non-missing laboratory data was randomly removed with varying frequencies from two large data sets, and we then compared the ability of four methods – missForest, mean imputation, nearest neighbor imputation, and multivariate imputation by chained equations (MICE) – to impute the simulated missing data. We characterized the accuracy of the imputation and the effect of the imputation on predictive ability in **two large datasets**.

Results: MissForest had the least imputation error for both continuous and categorical variables at each frequency of missingness, and it had the smallest prediction difference when models used imputed laboratory values. In both data sets, MICE had the second least imputation error and prediction difference, followed by nearest neighbor imputation and mean imputation.

Conclusion: MissForest is a highly accurate method of imputation for missing laboratory data and outperforms other common imputation techniques in terms of imputation error and maintaining predictive ability with imputed values **in two clinical predicative models**.

Introduction

“ You can have data without information, but you cannot have information without data”

- Daniel Keys Moran

Missing data is a nearly ubiquitous problem when conducting research, particularly when using large data sets. Missing data can occur in the form of random or non-random patterns. Non-random missing data can introduce systematic error and make the study population less representative of the general population. Although random missing data does not introduce systematic error, it leads to significant loss in statistical power and predictive ability. Missing data are rarely completely at random and must be carefully managed. Multi-analyte Assays with Algorithmic Analyses (MAAAs) are a relatively new approach to leveraging value from laboratory data to predict clinical outcomes. Several of these are now available with implemented CPT codes (i.e. FibroSure, Risk of Ovarian Malignancy (ROMA), PreDx Diabetes Risk Score), but it is not known how robust MAAA models are when individual laboratory data points are missing.

Recent developments in machine learning have used laboratory data to build MAAA models to predict health care outcomes [1-3] - These models can be sensitive to missing **completely at random** laboratory data, which may result from hemolyzed samples, clumped platelets, or other uncommon sample or processing problems. It is not known whether modern imputation methods can robustly address the problem of missing data, and whether predictive models will remain accurate when imputed values are used.

Multiple methods have been developed in order to deal with missing data, including single imputation, multiple imputation, **multivariate** imputation by chained equations (MICE)[3] , nearest neighbor estimation [4], and missForest[5] . Although these have been shown to be effective with other types of missing data, little data exists regarding the absolute or comparative effectiveness of these methods in accurately imputing missing laboratory data in predictive models. The aims of our study were: (1) to compare the accuracy of four different imputation methods for missing **completely at random**

laboratory data in two large data sets; and (2) to compare the effect of imputed values from each method on the accuracy of predictive models based on these data sets.

Methods

The University of Michigan (UM) Cirrhosis Cohort and Predictive Model for Hepatocellular Carcinoma

Between January 2004 and September 2006, consecutive patients with cirrhosis but no detectable HCC were prospectively identified and entered into a surveillance program using ultrasound and alpha fetoprotein (AFP), as has been previously described in greater detail [6]. Patients were enrolled if they had Child-Pugh class A or B cirrhosis and absence of known HCC at the time of initial evaluation. Patients diagnosed with HCC within the first six months of enrollment (prevalent cases) were excluded. Other exclusion criteria included clinical evidence of significant hepatic decompensation (refractory ascites, grade 3-4 encephalopathy, active variceal bleeding, or hepatorenal syndrome), co-morbid medical conditions with a life expectancy of less than one year, prior solid organ transplant, and a known extrahepatic primary tumor. Patients were followed until the time of HCC diagnosis, liver transplantation, death, or until the study was terminated on July 31, 2010.

The following demographic and clinical data were collected at the time of enrollment: age, gender, race, body mass index (BMI), past medical history, lifetime alcohol use, and lifetime tobacco use. Data regarding their liver disease included the underlying etiology and presence of ascites, encephalopathy, or esophageal varices. Laboratory data of interest at the time of enrollment included: platelet count, aspartate aminotransferase (AST), alanine aminotransferase (ALT), alkaline phosphatase, bilirubin, albumin, international normalized ratio (INR), and AFP. This data set was used as the basis of a published predictive model to identify patients with hepatocellular carcinoma with a c statistic of 0.70 [2].

1
2 The UM Inflammatory Bowel Disease Cohort and Predictive Model for Thiopurine Clinical Response.
3
4 The study sample included all patients who had thiopurine metabolite analysis, CBC, and a
5
6 comprehensive chemistry panel drawn within a 24-hour period at the University of Michigan between
7
8 May 1, 2004, and August 31, 2006 and is described in greater detail in the manuscript [1]. This study
9
10 was approved by the University of Michigan Medical Institutional Review Board with a waiver of
11
12 explicit consent from the subjects. The patient sample included 774 cases, in a total of 346 individuals.
13
14 For the analysis of the outcome of clinical response to thiopurines, 5 exclusion criteria were applied:
15
16 exclusion of patients who did not have IBD, exclusion of patients who had not started on thiopurines at
17
18 the time when metabolites were measured, exclusion of patients on biologic anti-tumor necrosis factor
19
20 therapy, exclusion of patients without documentation of their clinical status at the time of laboratory
21
22 measurement, and exclusion of patients who had an infection that confounded assessment of clinical
23
24 response. This data set was used as the basis of a predictive model to identify patients with clinical
25
26 response to thiopurine immune suppressant medication with a c statistic of 0.86[1].
27
28
29
30
31
32
33
34

35 **Description of Imputation Techniques**

36
37 We compared missForest with three other commonly used imputation methods that can handle both
38
39 continuous and categorical variables, namely, mean Imputation, Nearest Neighbor Imputation, and
40
41 Multivariate Imputation by Chained Equations (MICE). We briefly introduce these methods below.
42
43

44 The recently proposed missForest method makes use of highly flexible and versatile random forest
45
46 models [7, 8] to achieve missing value imputation. It creates a random forest model for each variable
47
48 using the rest of the variables in the data set and uses that to predict the missing values for that variable.
49
50

51 This is done in a cyclic fashion for all variables and the entire process is iteratively repeated until a
52
53 stopping criterion is attained. The advantages of using the random forest model are that it can handle
54
55 both continuous and categorical responses, requires very little tuning, and provides an internally cross-
56
57 validated error estimate. This was implemented via the 'missForest' package available in R.
58
59
60

1
2 Nearest Neighbor algorithms were originally proposed in the supervised pattern recognition literature.
3
4 Troyanskaya et al. proposed an imputation method based on nearest neighbor search [4]. The basic idea
5
6 is to compute a distance measure between each pair of observations based on the non-missing
7
8 variables. Then the k-nearest observations that have non-missing values for that particular variable are
9
10 used to impute a missing value via a weighted mean of the neighboring values. In order to
11
12 accommodate both continuous and categorical variables the Gower distance is used [9]. For the
13
14 categorical variables we imputed the missing values by weighted mode instead of a weighted mean as
15
16 used for continuous variables. Cross-validation error measures are used to select the optimal number of
17
18 nearest neighbors denoted by k. The function 'kNN' in R package 'VIM' was used to implement this
19
20 method.
21
22
23
24

25 Mean imputation is one of the most naïve and easiest methods for imputing missing values. The mean
26
27 (for continuous variables) or mode (for categorical variables) of the non-missing values of each
28
29 variable were used to impute the missing values. This method does not take advantage of any
30
31 correlation among the variables and therefore can perform rather poorly when such correlations are
32
33 present.
34
35
36

37 MICE was proposed by Van Buuren et al. [3]. It requires the user to specify a conditional model for
38
39 each variable, using the other variables as predictors. By default we used a linear regression model for
40
41 continuous variables, a logistic regression model for binary variables, and a polytomous logistic
42
43 regression for categorical variables with more than two levels. The algorithm works by iteratively
44
45 imputing the missing values based on the fitted conditional models until a stopping criterion is
46
47 satisfied. In that way it is very similar to the missForest algorithm; the main difference being that
48
49 missForest uses more flexible decision trees for each conditional model. We implemented this in R
50
51 using the package 'mice'[10].
52
53
54
55
56
57
58
59
60

Statistical Analysis

We used two separate studies to perform the comparison between the methods of imputation described in the previous section. We describe the studies, implementation details, and our results below. The structure of the statistical analysis is the same for both studies. We start with a published predictive model built with the training data set. The test set refers to observations that were not part of the training set; these were solely used for assessing the performance of the model. The test sets did not have any missing values, so we randomly removed a proportion of values to simulate data missing **completely at random**. We then imputed the missing values by the four previously discussed methods, and the imputed laboratory results were compared to the actual values that were removed from the data set. We then used the imputed data to make clinical outcome predictions with the published models, and the results were compared with the predictions made using the complete test data with no missing data. **The prediction models were created using both logistic models and random forest models as some might argue that using random forest imputation methods might favor the random forest prediction models.** We use the average relative error (for continuous variables) and misclassification error (for categorical variables) to assess the imputation performance for **both the logistic and random forest models**. To quantify the effect of the imputation on predictive models we compare the predicted classes from non-missing test data to predicted classes from imputed test data and compute the misclassification error **for both the logistic and random forest models**. This is important because if a particular variable had very little influence on the predictive model, then larger imputation errors are tolerable, resulting in negligible loss of prediction accuracy. On the other hand, small imputation errors in very important variables might lead to significantly different predicted class, which is of greater clinical concern. We varied the frequency of missing values to change the difficulty of the imputation problem. We report the average results over multiple random runs. We found that the Nearest Neighbor results are quite robust to the choice of number of nearest neighbors (k) if k is moderately large, therefore we fixed the number of nearest neighbors at 5 in both studies.

Results

Cirrhosis Cohort and HCC model

This study evaluated the effect of imputation on a published predictive model for HCC based on 21 predictor variables that included demographic, clinical and laboratory values using random forest modeling [2]. The random forest model was developed on a data set of 446 patients collected at University of Michigan (UM cohort). It proved to be more accurate than traditional logistic regression models. Of the 21 variables, 10 of them were categorical in nature while 11 were continuous. We used the first 200 observations from the publicly available data set from the HALT-C trial as our test set and randomly replaced 10%, 20%, or 30% of the observations with missing values. The process was repeated for 1000 replications and we report the average results.

The accuracy of the four imputation methods for both continuous and categorical variables are compared in Figure 1 for the cirrhosis cohort and HCC model. Figure 1A and 1B represents the logistic model and Figure 1C and 1D reflect the random forest prediction model. The vertical axis plots the percentage relative error for continuous variables and percentage misclassification error for categorical variables, while the horizontal axis groups the results according to the proportion of missing values. Each boxplot represent the error measure over 1000 random replications. As expected, the imputation error increases on an average as we increase the proportion of missing values in the test data but the variation tends to reduce slightly which is due to averaging over many more missing observations. MissForest has the least imputation error for both continuous and categorical variables at each level of missing proportion, followed by MICE, NN, and mean imputation of continuous laboratory values. MICE and NN have similar imputation accuracy for categorical variables. MissForest works well using both logistic and random forest prediction models.

In Figure 2A and 2B, which represents the logistic and random forest prediction models respectively, the vertical axis plots the error measure for imputation on predictive model, obtained by comparing the predicted classes (Low Risk/High Risk) for each test observation with no-missing values against the

1
2 predicted class after imputing the artificial missing values. Therefore an error measure of 5 on the
3
4 vertical axis implies that 5% of the test observations had their predicted classes wrongly switched
5
6 (either low risk => high risk or high risk => low risk) due to the imputation. As above, each boxplot
7
8 reflects the results of 1000 random runs. It is clear from the figure that missForest performs the best
9
10 with NN and MICE following closely. The gap increases as we increase the proportion of missing
11
12 observations making the problem harder.
13
14

15 16 17 18 Inflammatory Bowel Disease Cohort and Thiopurine Clinical Response Model 19

20
21 Waljee et al. showed that random forest models using laboratory values outperform 6-thioguanine (6-
22
23 TGN) metabolite tests as well as traditional logistic regression models in predicting clinical response to
24
25 thiopurines [1]. The analysis was carried out on a data set collected at University of Michigan that
26
27 included 395 patients. 26 variables, which included 25 laboratory values and age, were used to predict
28
29 the clinical response of each patient. In this study all of the variables were continuous in nature. To
30
31 create a separate test set, we split the data set randomly into a training set consisting of 250
32
33 observations and a test set of 145 observations, using stratified sampling to keep the ratio of clinical
34
35 responder to non-responders fixed. We then introduced random missing values into the test set as
36
37 before and performed the same comparative study of the four imputation methods. The whole process
38
39 was replicated 1000 times to obtain stable results. Below we summarize our findings via boxplots.
40
41
42

43
44 Again in this study, missForest beats its competitors in both imputation accuracy and the effect of
45
46 imputed values on the accuracy of clinical predictions based on the logistic model (Figure 3A) and the
47
48 random forest model (Figure 3B). The trends remain the same for imputation error, Figure 4A and B
49
50 representing logistic models and random forest models respectively, with MICE coming out second
51
52 best followed by NN and mean imputation. For predictive accuracy we find the relative order becomes
53
54 missForest > MICE > mean imputation > NN. This also shows that best method with respect to
55
56 imputation error need not be the best when we consider the effect of imputation on predictive models.
57
58
59
60

1
2 The performance gap between missForest and MICE is considerably lower than in the previous study.
3
4 This might be explained by the fact that in the thiopurine study, both the training and test sets came
5
6 from the same cohort, as we generated the training and test sets by random splits, while in HCC study
7
8 the training and test sets were completely different cohorts leading to an extra degree of variation.
9

10 11 **Discussion**

12
13 We have performed an extensive simulation study using two clinical datasets and two published
14
15 predictive models to compare the performance of four methods of missing value imputation for missing
16
17 data completely at random. We included both local (randomForest) and global (logistic) modeling
18
19 approaches to avoid bias that might favor a local (MissForest) imputation approach. While the
20
21 superiority of MissForest for imputation of missing lab values this will not be generalizable to all
22
23 predictive models or datasets, this manuscript highlights the value of missForest to impute missing
24
25 data. We compared four popular methods namely, missForest, Nearest Neighbor, MICE and mean
26
27 imputation, in two studies simulating data missing completely at random. We found that these
28
29 simulation methods consistently produced the lowest imputation error and had the smallest prediction
30
31 difference when models used imputed laboratory values. In addition, the ready availability of the
32
33 freeware R package makes missForest and its simulations a very convenient solution for any practical
34
35 missing value problems. The main limitations of these simulations as a solution to missing laboratory
36
37 data for predictive modeling applications are: a requirement for skilled R programming for
38
39 implementation, and slightly more demanding computational needs, compared to NN or MICE
40
41 methods. An additional limitation in this study is that these simulations did not address the issue of data
42
43 missing for non-random reasons. There could be an association between the clinical outcome of interest
44
45 and the missingness of certain predictors. At this point, we cannot generalize these results to situation
46
47 in which data is missing for non-random reasons.
48
49
50
51
52
53
54

55
56 The small absolute changes in predictions with these models, despite 10-30% missing laboratory data,
57
58 speak to the robustness of these multi-analyte assays with algorithmic analysis (MAAAs). MAAAs are
59
60

1
2 currently a hot topic, and several have been released with CPT codes in 2012. One example is the HCV
3
4 FibroSure (LabCorp, code 0001M) which uses ALT, alpha 2 macroglobulin, apolipoprotein A1, total
5
6 bilirubin, GGT, and haptoglobin to estimate fibrosis and necroinflammatory activity in the liver in
7
8 patients with hepatitis C. With the increasing complexity of these models, and increasing numbers of
9
10 analytes, the risk of **missing completely at random** values increases, and methods to cope with missing
11
12 values and preserve the accuracy of the model are needed. missForest appears to be a robust and
13
14 accurate approach to the issue of missing laboratory values when used in these two MAAAs **and may**
15
16 **be applicable to other datasets with missing completely at random datasets.**
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Contributions:

Ashin Mukherjee – study concept and design, statistical analysis and interpretation of the data, drafting of the manuscript, critical revision of the manuscript

Akbar Waljee - study concept and design, acquisition of the data, statistical analysis and interpretation of the data, drafting of the manuscript, critical revision of the manuscript and study supervision.

Amit Singal - acquisition of data, drafting of the manuscript and critical revision of the manuscript.

Yiwei Zhang - statistical analysis and interpretation of the data and critical revision of the manuscript.

Jeffrey Warren - study concept and design, critical revision of the manuscript.

Ulysses Balis - study concept and design, critical revision of the manuscript.

Jorge Marrero - study concept and design, critical revision of the manuscript.

Ji Zhu - statistical analysis and interpretation of the data and critical revision of the manuscript.

Peter Higgins - study concept and design, acquisition of the data, statistical analysis and interpretation of the data, critical revision of the manuscript and study supervision.

Competing Interests: None

Ethics approval: University of Michigan, IRB

Data Sharing: Agreement: There is no additional data available.

Figure Legends

FIGURE 1. Imputation error comparison for categorical and continuous variables for four competing imputation methods at three levels of the proportion of missing values for the logistic prediction model (Figure 1A and 1B) and random forest prediction model (Figure 1C and 1D) in the HCC study

FIGURE 2. Percentage of wrongly predicted observations after missing value imputation by the four competing methods at three levels of missing value proportions in the test data for the logistic prediction model (Figure 2A) and the random forest prediction model (Figure 2B) in the HCC study.

FIGURE 3. Imputation error for four competing imputation methods at three levels of the proportion of missing values for the logistic prediction model (Figure 3A) and random forest prediction model (Figure 3B) in the Thiopurine response model

Figure 4. Percentage of wrongly predicted observations after missing value imputation by the four competing methods at three levels of missing value proportions in the test data for the logistic prediction model (Figure 4A) and the random forest prediction model (Figure 4B) in the Thiopurine response model.

References

- 1 Waljee AK, Joyce JC, Wang S, et al. Algorithms outperform metabolite tests in predicting response of patients with inflammatory bowel disease to thiopurines. *Clinical gastroenterology and hepatology : the official clinical practice journal of the American Gastroenterological Association*. 2010 Feb;**8**(2):143-50.
- 2 Singal AG, Waljee AK, Mukherjee A, Higgins PD, Zhu J, Marrero JA. Machine Learning Algorithms Outperform Conventional Regression Models in Identifying Risk Factors for Hepatocellular Carcinoma in Patients With Cirrhosis. *Gastroenterology*. 2012 May;**142**(5):S984-S.
- 3 van Buuren S, Boshuizen HC, Knook DL. Multiple imputation of missing blood pressure covariates in survival analysis. *Stat Med*. 1999 Mar 30;**18**(6):681-94.
- 4 Troyanskaya O, Cantor M, Sherlock G, et al. Missing value estimation methods for DNA microarrays. *Bioinformatics*. 2001 Jun;**17**(6):520-5.
- 5 Stekhoven DJ, Buhlmann P. MissForest--non-parametric missing value imputation for mixed-type data. *Bioinformatics*. 2012 Jan 1;**28**(1):112-8.
- 6 Singal AG, Conjeevaram HS, Volk ML, et al. Effectiveness of hepatocellular carcinoma surveillance in patients with cirrhosis. *Cancer epidemiology, biomarkers & prevention : a publication of the American Association for Cancer Research, cosponsored by the American Society of Preventive Oncology*. 2012 May;**21**(5):793-9.
- 7 Liaw A, Wiener M. Classification and Regression by randomForest. 2002;**2**(3):18-22.
- 8 Breiman L. Random forests. 2001;**45**(1):5-32.
- 9 Gower JC. A General Coefficient of Similarity and Some of Its Properties. *Biometrics*. 1971;**27**(4):857-71.
- 10 Buuren van S, Groothuis-Oudshoorn K. mice: Multivariate Imputation by Chained Equations in R. *Journal of Statistical Software*. 2011;**45**(3).

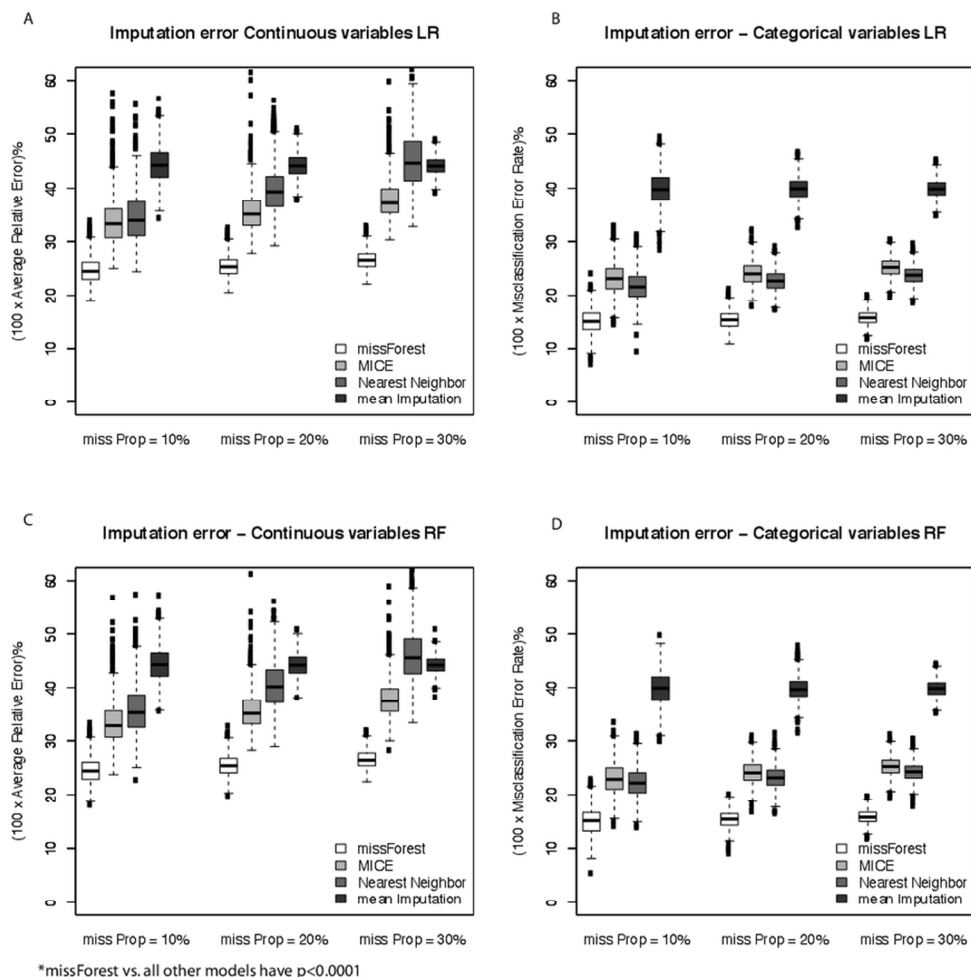


FIGURE 1. Imputation error comparison for categorical and continuous variables for four competing imputation methods at three levels of the proportion of missing values for the logistic prediction model (Figure 1A and 1B) and random forest prediction model (Figure 1C and 1D) in the HCC study. 90x90mm (300 x 300 DPI)

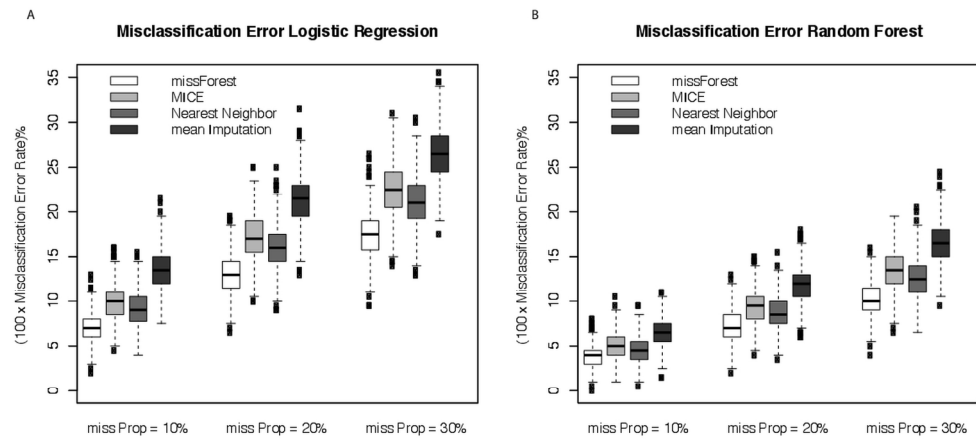


FIGURE 2. Percentage of wrongly predicted observations after missing value imputation by the four competing methods at three levels of missing value proportions in the test data for the logistic prediction model (Figure 2A) and the random forest prediction model (Figure 2B) in the HCC study. 185x90mm (300 x 300 DPI)

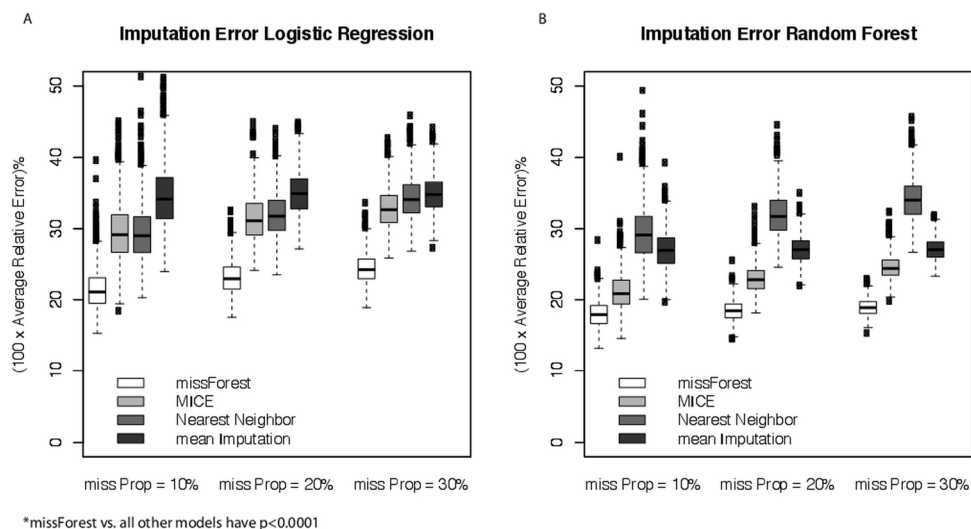
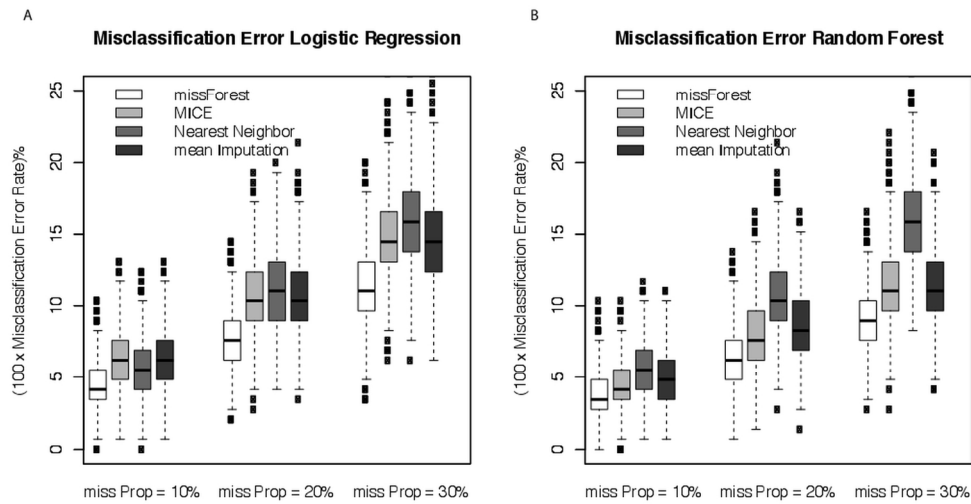


FIGURE 3. Imputation error for four competing imputation methods at three levels of the proportion of missing values for the logistic prediction model (Figure 3A) and random forest prediction model (Figure 3B) in the Thiopurine response model.
162x90mm (300 x 300 DPI)

review only



26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 4. Percentage of wrongly predicted observations after missing value imputation by the four competing methods at three levels of missing value proportions in the test data for the logistic prediction model (Figure 4A) and the random forest prediction model (Figure 4B) in the Thiopurine response model.
161x90mm (300 x 300 DPI)