

No solution? Evaluating Multiple Imputation

Anthony S. Chapman, Dr Steven Turner, Dr. Wei Pang

Abstract

Blah

1. Introduction

Data collection has been increasing Missing data is inevitable (human and computing reasons, i.e. people not putting it in or computer corrupting it,) Non-computing people either impute willy-neely or ignore missing data - need to use as much data as you can. (Ask Graham about theory about using as much data as possible for better analysis) Many imputation algorithms out there with many parameters, which is best? Need

2. Background

Talk about something [1] [5]

3. The Problems

3.1. Incompleteness

With so many new data being collected daily [5], it was inevitable that some of the data would have missing values [3], whether they be through human error or computational inefficiency. Although there are ways to combat missing data such as mean-value imputation or multiple imputation [3, 4, 2], many researchers whom are not very computational or statistically confident would rather disregard any records with missing values [1, 2, 3, 4, 5]. As an example, in [1], the authors decided to use 2,758 records for analysis out of the possible 44,261 mainly due to missing data, this is a mere 6.2% out of the records available. There must be a way for even non-computing or non-statistical researchers to benefit from the tools available.

3.2. Will it work on my data?

This next problems arises when a researchers does decide to use the data with missing values but does not have sufficient knowledge to apply the available

methods, like Multiple Imputation by Chained Equations (MICE [2]) using the computational language R [3] or the Impute Missing Values function in the statistical software SPSS [4]. The problem is, how do you know if the imputed values are representative to the truth, how do you know whether record 2,754 column 5 is male or not after you apply the imputation method.

Even if the imputation method has been proven to work on someone else's dataset such as [6], this no indication it will work for yours, unless you have the exact dataset as them, which is unrealistic. This is due to the many reasons and ways that missing data is creates, for example there might be a relationship between one missing value and another one.

In order to test whether an imputation method works on your dataset, you need something to compare the results to, a benchmark, like this one would be able to analyse what effect of the methods. Unfortunately, it is very difficult to find a complete dataset which contains the same characteristics of your own dataset, there will always be differences.

3.3. Which imputation is best for me

There is nothing to easily compare different imputation techniques, post researches haven't got computing background... (need a way to formally say that statement, maybe number of disciplines VS computing??)

You can't compare an evaluation of one imputation on data A and a different imputation on data B, is chocolate better than bacon?

4. Possible Solutions

4.1. Incompleteness

Imputations solves incompleteness, you just have to be careful how you use it.

Can't blindly impute something as it might result in bias and unreliable results.

4.2. Testing your own data

Use your own level or missingness as a benchmark and create mini-me's as bench-mark. You are the closest thing to yourself. Group theory stuff, multidimensional-mixed data distance measurements, Gower, medoids, widths and dissimilarities.

Just because it worked on someone else, doesn't mean it works for you, cite papers who test specific datasets.

4.3. Comparing Imputations

Will now be able to compare different imputations on your own dataset with "normalised " results for comparison.

5. Conclusion

It's better to use all the data you can but can't blindly imputation. This framework indicates whether your data

6. Discussion

Working on implementing this, CIEMI, any researcher regardless the computing ability will be able to use it.

References

- [1] Amy M. Branum, Jennifer D. Parker, Keim Sarah A., and Schempf Ashley H. Prepregnancy body mass index and gestational weight gain in relation to child body mass index among siblings. *American Journal of Epidemiology*, 174(10):1159–1165, 2011.
- [2] ALAN C. COCK. Working with missing values. *Journal of Marriage and Family*, 174(67):10121028, 2005.
- [3] Therese D. Pigott. A review of methods for missing data. *Educational Research and Evaluation*, 7(4):. 353–383, 2001.
- [4] Donald B. Rubin. An overview of multiple imputation.
- [5] ScienceDaily. Big data, for better or worse. 2013 (accessed: January 18, 2016). <http://www.sciencedaily.com/releases/2013/05/130522085217.htm>.
- [6] Anoop D. Shah and Jonathan W. Bartlett. Comparison of parametric and random forest mice in imputation of missing data in survival analysis. 2014.