

# Ordinal Variables and Missing Data

February 16, 2020

## General

- This paper should be an intersection of ordinal variables and missing data. Jeff thinks there is room there for a contribution
- The general idea is that general treatments of missing data (listwise deletion, multiple imputation etc.) lead to different results – depending on the type of variable. How you treat **Don't Know** and **Refused** and how this treatment affects the results depends greatly on the type of variable that you use to impute the missing data. The idea is that you would need to approach things differently depending on whether you use a nominal, interval, or ordinal variable to fill in values for **Don't Know/Refused**. So I am developing a method to specifically treat **Don't Know/Refused** with ordinal variables that improves over current uses that are generic for all types of variables
- This method should be a customized form of multiple imputation that is closer to the data than general multiple imputation. Jeff and Skyler developed affinity scores and hot decking in their BJPS paper. They used the number of exact matches (in the form of other participants) to calculate the affinity score. Instead, I should use a weighted distance solution between ordinal variable categories. I would use the OP model from the blocking paper to weight the distances between the categories in matches (in the form of other participants). In other words, I would use the underlying ordered probit numbers to create the weights. So this would be a specific ordinal variable adjustment of the affinity score building
- Set up chapter structure

## Code

- My code functions for one variable with NAs, for 10,000 iterations
  - I have saved results for `inc`, `age`, `Dem`, `Rep` for `hd.ord` and `hd.norm`
  - For all 4, `na.omit` performs the best. This isn't surprising, since deleting observations with missing values shouldn't be a problem with MCAR
- My code does not function When run for several variables with NAs
  - It throws up a replacement error somewhere down the line when run for 10,000 iterations. It's never in the same place, so it must be something random
  - Always the same: # of NAs overall. Not always the same between some runs: # of NAs per column, # of rows with NAs
  - Find out why my code doesn't work for NAs in several variables
    - \* The number of NAs inserted wasn't actually the issue. It was in the `OPMord` code
    - \* I painfully ran 1,000 iterations for two variables, gradually adding one line at a time to the function. I discovered the error was where I assign values to `df.cases`, specifically where the columns are all combined and the resulting vector added to the data as `educ.new`

- \* I saved the vector to an empty vector, ran everything for 1,000 iterations, then looked at the vectors and the corresponding versions of `df.cases`
- \* There were rows with all NAs in `df.cases`, which makes the vector shorter than the data, which means it can't be assigned to it as a column. I reran the `df.cases` creation code to find the culprit
- \* I had overlooked that `int.df$Values` has one value fewer than there are variable levels, since it lists cutpoints between the levels. For 6 levels, I had assigned `int.df$length(levels(...))`, which picks the 6th value of `int.df$Values` – but it only goes to 5. I had to add `-1`. Now it's working
- Method to insert NAs
  - I tried `prodNA`, which works for MCAR but gives me nothing for MAR. And MCAR works well for `na.omit`, since by definition it's a random sample of missing data, which doesn't matter
  - Use `mice ampute()`
    - \* It has MCAR and MAR options. MAR is what I need
    - \* It doesn't work on just one variable; it needs at least two
    - \* `ampute()` works very well
- With the code adjusted (`-1`), `OPMord` now works. However, some of the `int.dfs` don't have all the education levels: `int.df` with all levels has 6 rows (one fewer than the levels, since each row shows a cutoff), but some have 5. This means `OPMcut` now doesn't work
  - I could adjust `OPMcut` to work with fewer levels, but then I would be, in the end, taking means of most data with all levels and some data with fewer levels. That's problematic
  - It's better to discard the iterations where `int.df` has 5 rows after `OPMord` has run and then continue with the other functions. I set the code up to do that. However, after about 40 percent of running 12,500 iterations, there was an error: I hadn't factored in that the `int.dfs` could also have fewer rows than that. A few of them had 4 rows, which caused the error. I adjusted the code to now discard the iterations where `int.df` has anything other than 6 rows
- The data ran for 12,500 iterations, 80 percent NAs, `Rep` and `inc`, `hot.deck.ord`, `hot.deck.norm`, `na.omit`. The results looked promising: `hd.ord` performed better than `hd.norm` and `na.omit`
- More methods
  - Include `mice`
  - Include `Amelia`
  - Use the defaults for both since that's what over 90% of people do (even if they shouldn't)
  - Include `hot.deck.norm` on the original education values (I ran it on the cutpoints)
- `hot.deck impContinuous`
  - No need to use it, since no variable is continuous (age has more values than `sdCutoff` but it's nominal, not continuous)
- `hot.deck sdCutoff`
  - No need to change the default
- Add more variables/NAs

- Add NAs to 6 of the variables
- Demographics I still have: Age(numeric), employment (5 levels), ideology (liberal, conservative, neither), following public affairs (ordinal, 4 levels), media interest (numeric, accumulative count of activities), participation (numeric, accumulative count of activities)
- I can't add them all, since `mice` and `Amelia` don't work when there is collinearity
- I used code to find and discard variables with high correlation
- Percentage of NAs
  - I'm currently using 80 percent. Jeff used 20, 50, and 80
  - I did one run each for 20, 50, and 80 percent
- Preliminary results
  - Across the board, `hd.ord` performs better than normal `hot.deck`, but worse than `Amelia` and `mice`
  - Across the board, `Amelia` is better than `mice` for `inc` and `age`. It's also better for `interest`, except for 20 percent NAs, where `mice` wins
  - For `Dem`, `Female`, and `White`, `Amelia` also edges first, though `mice` sometimes is ahead by one unit or so in the fourth decimal
  - The difference is less pronounced for the binary variables and most visible in the nominal ones. It is worst for "age", which has the most unique values
  - As the percentage of NAs increases, the estimates are further off from the true means. This is to be expected for all methods, but it affects `hd.ord` and `hot.deck` much more than `Amelia` and `mice`
- Jeff likens this to the search for the pot of gold at the end of the rainbow. There will be something there, somewhere. I just need to find one specific scenario, one specific set of circumstances where I do particularly well or better. He gave me a bunch of pointers/ideas where to go from here. I'm using `ampute()` to generate NAs MAR. `ampute()` is part of `mice` – it's not a great surprise that `mice` then performs really well. `amelia()` performs even better – Jeff thinks `amelia()` calls `mice` at some point. `ampute()` has lots of default settings that are not immediately obvious. Jeff says to mess with those
- Run it with `bycases = FALSE` (TRUE introduces a very specific pattern of missingness), `cont = FALSE`, `type = "MID"`
  - I had to take out `White` because `ampute()` complained there were too many variables for `bycases = FALSE`
  - I ran it for 20 percent NAs and saved `framing.nas` and `ampute()` output – no significant change in results: `hd.ord` is better than `hot.deck` across the board. `Amelia` is better than `mice` for `inc`, `age`, and `interest`. `mice` is better for `Dem` and `Female`
  - The same for 10 percent
- Running it with ANES data on Code Ocean
  - Previous small tests have shown that the results should not be different when compared to the framing data
  - I have results from running things for 2,363 iterations, but the `amelia` runtime results were much faster than expected here. I did this run with `lapply`. I want to rerun this with a loop, as before, before I write anything up here

- Runtimes
  - `hd.ord()` outperforms `hotdeck()` (minimally) and massively outperforms `mice` and `Amelia` on the framing data
  - I have some ANES runtime results for fewer iterations (1,307)
  - I also have ANES runtime results for more iterations (2,363) obtained through `lapply`. Weirdly, `amelia` is almost as fast as `hd.ord` here ... I want to rerun this with a loop before I write anything up here
- Own function to create NAs as MAR
  - Jeff: “It’s MAR if you ‘delete’ values in column 1 based on values in column 2”, e.g. ‘delete’ the values for age where income is 1 and where income is 5
  - I wrote a function that samples the ID numbers for a specified percentage of each selected variable’s unique values, then replaces the corresponding values of a different variable with NAs, based on the sampled ID numbers
  - I have run this function on the framing data
  - Seems to be more MCAR than MAR, since `na.omit` performs really well now
- Continue with paper notes

## Theory

- Step by step fill in the sections on Missing Data, Deletion, and Imputation
- Rework the introduction
  - Current stuff in there is very broad and nowhere near detailed enough (taken from the Kerwin application)