

# Milestone 8: To Turn In

Maria Burzillo

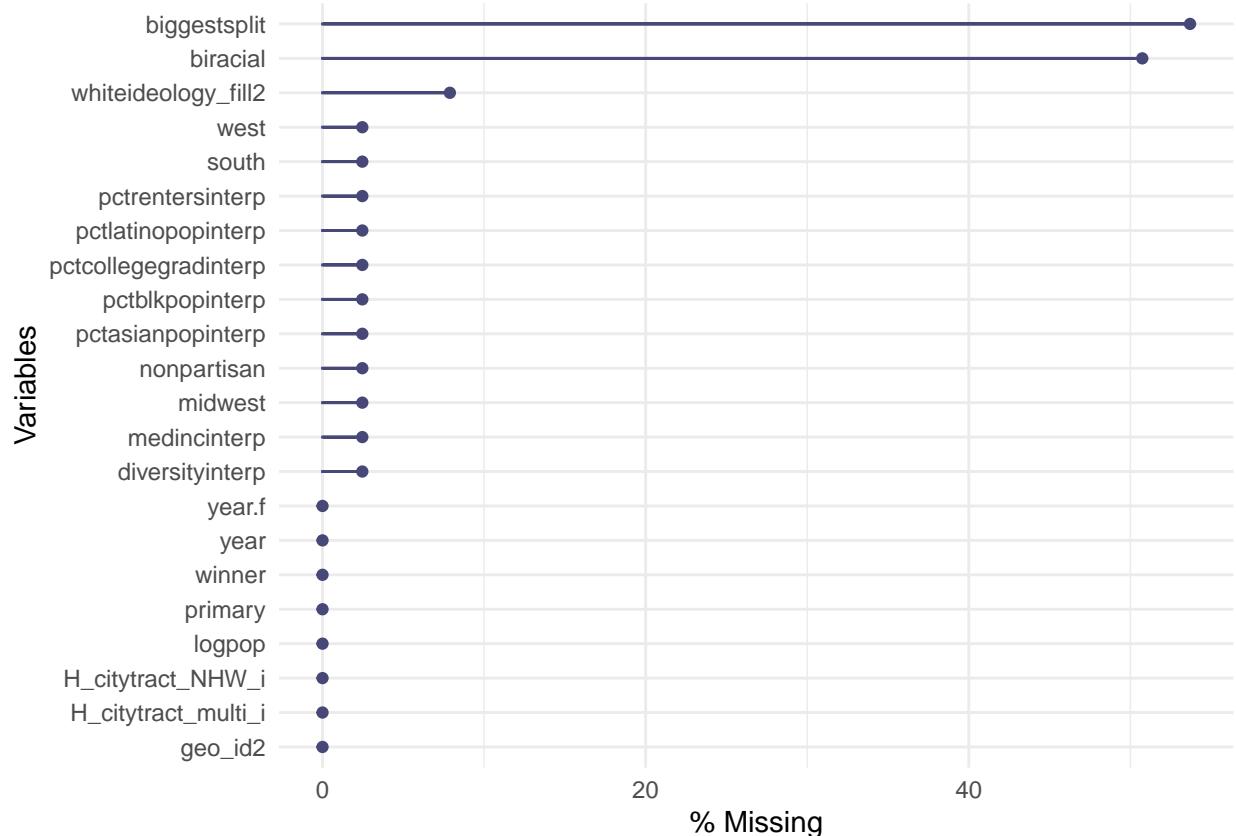
3/28/2020

## **Abstract**

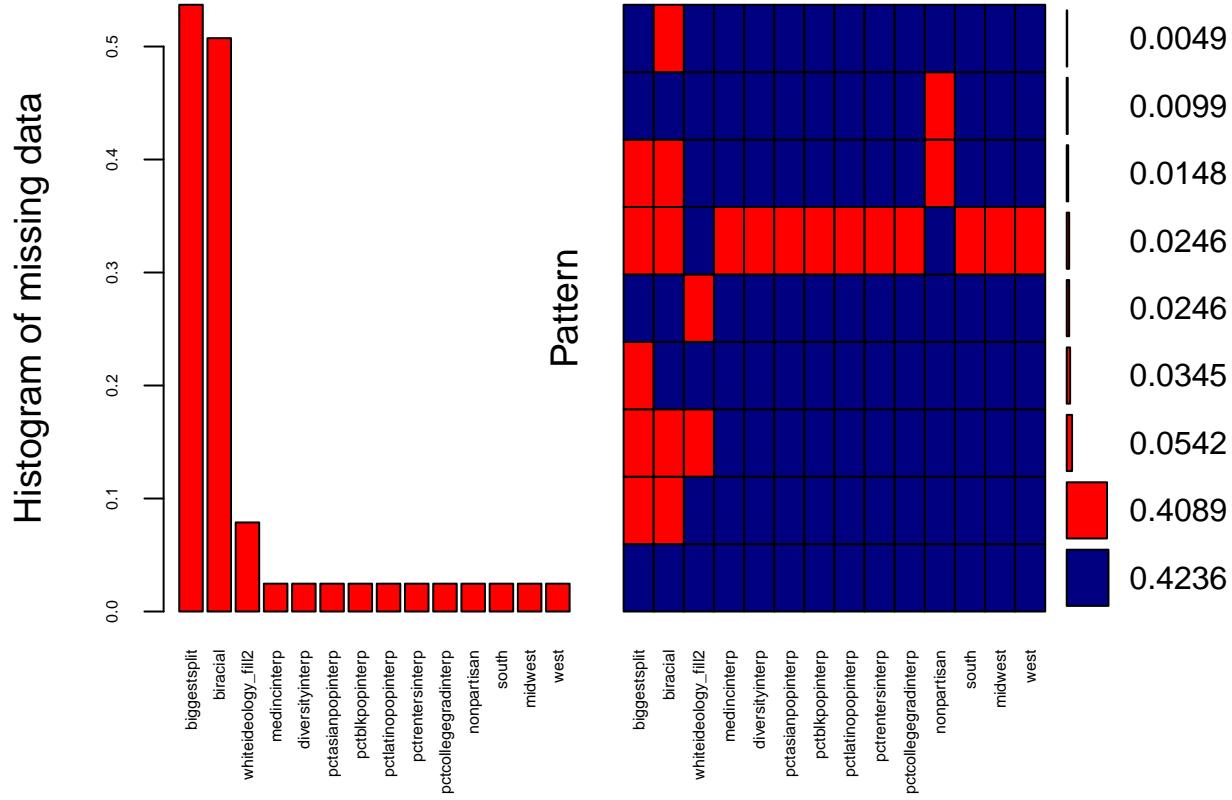
This is an extension of Jessica Trounstine’s “Segregation and Inequality in Public Goods” (2016). I was able to replicate the main results of Trounstine’s paper in R to suggest that racial segregation contributes to political polarization and decreased spending on public goods. Additionally, I extend the analysis by imputing missing data and rerunning Trounstine’s original model as a robustness check.

## **Extension 1**

One aspect of Trounstine’s paper with room for improvement is that there is a large amount of missing data in her datasets upon which she bases her analyses. Using R’s mice package (add citation), we can perform multiple imputation. By performing imputation multiple times, this helps account for the uncertainty inherent in the individual imputations. Before performing the multiple imputations, we will first look at the missing data to see if there are any patterns.



The above plot reveals that `biggestsplit` and `biracial` are the two variables with by far the largest percentage of missing values. Both of these are missing approximately 50% of their observations. Since `biggestsplit` is the dependent variable in our analysis with this dataset, this is an important fact.



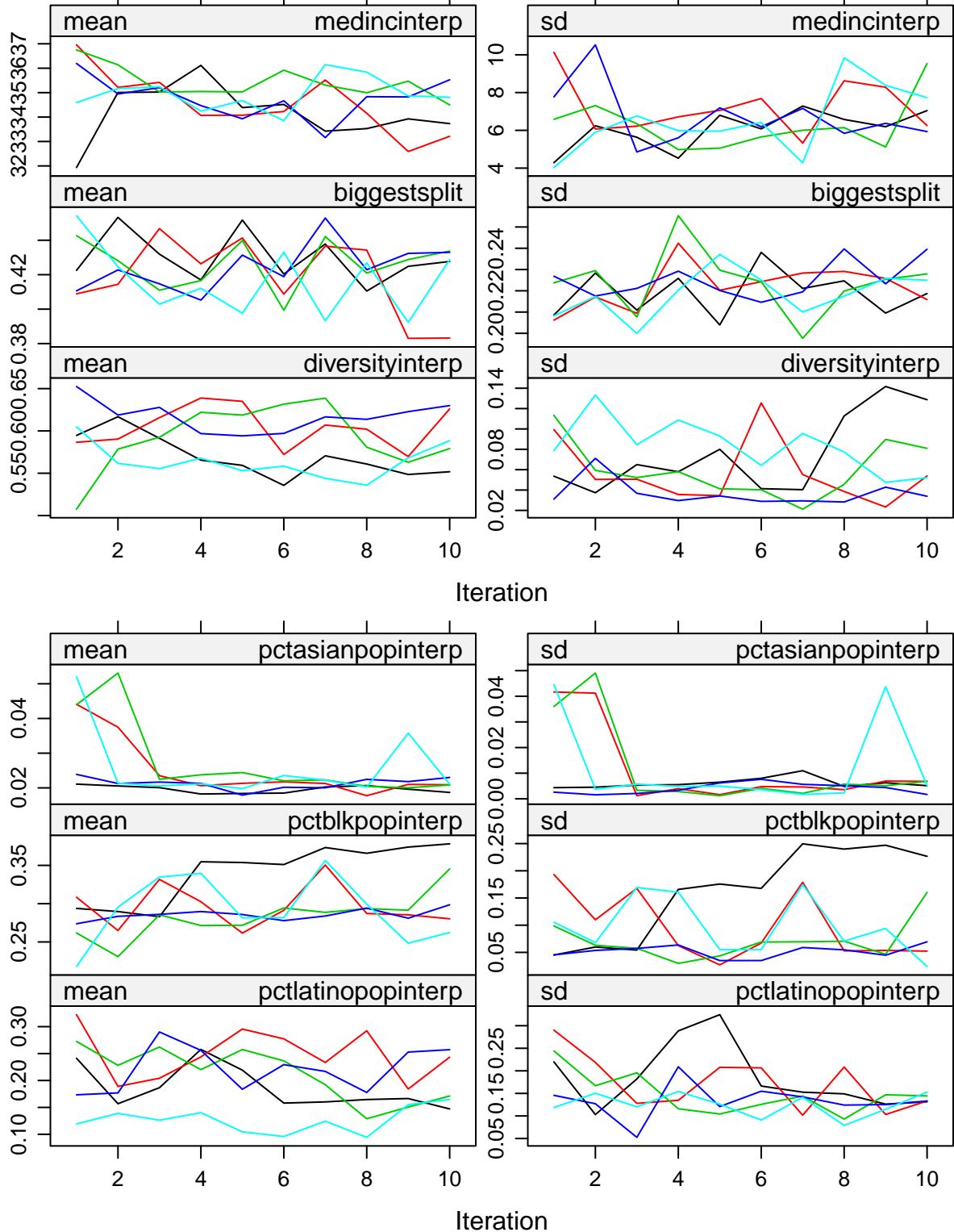
Variables sorted by number of missings:

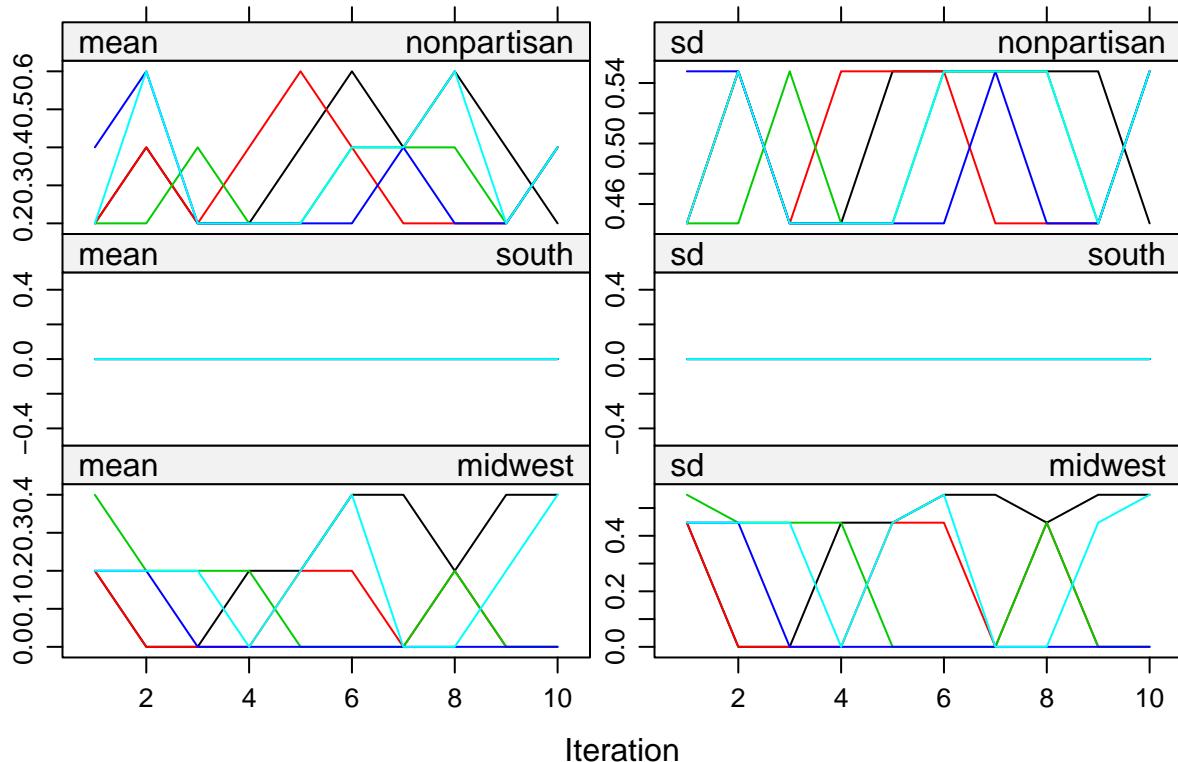
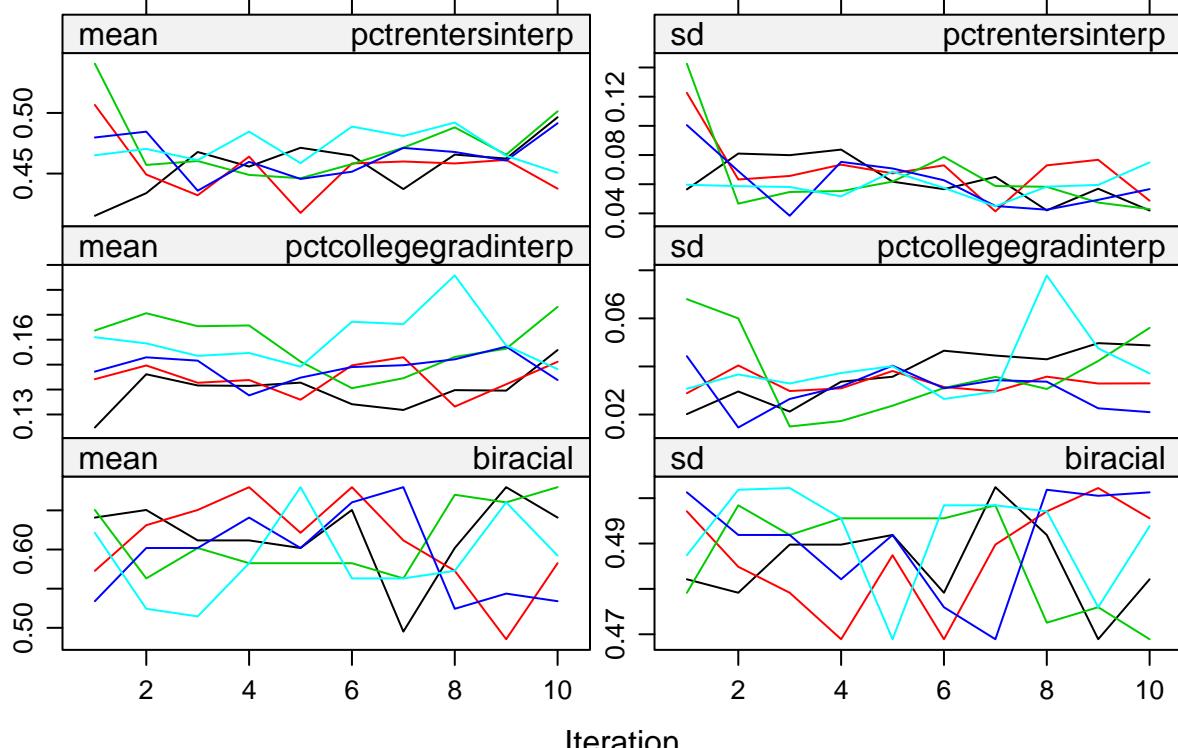
Variable	Count
biggestsplit	0.53694581
biracial	0.50738916
whiteideology_fill2	0.07881773
medincinterp	0.02463054
diversityinterp	0.02463054
pctasianpopinterp	0.02463054
pctblkpopinterp	0.02463054
pctlatinopopinterp	0.02463054
pctrentersinterp	0.02463054
pctcollegegradinterp	0.02463054
nonpartisan	0.02463054
south	0.02463054
midwest	0.02463054
west	0.02463054

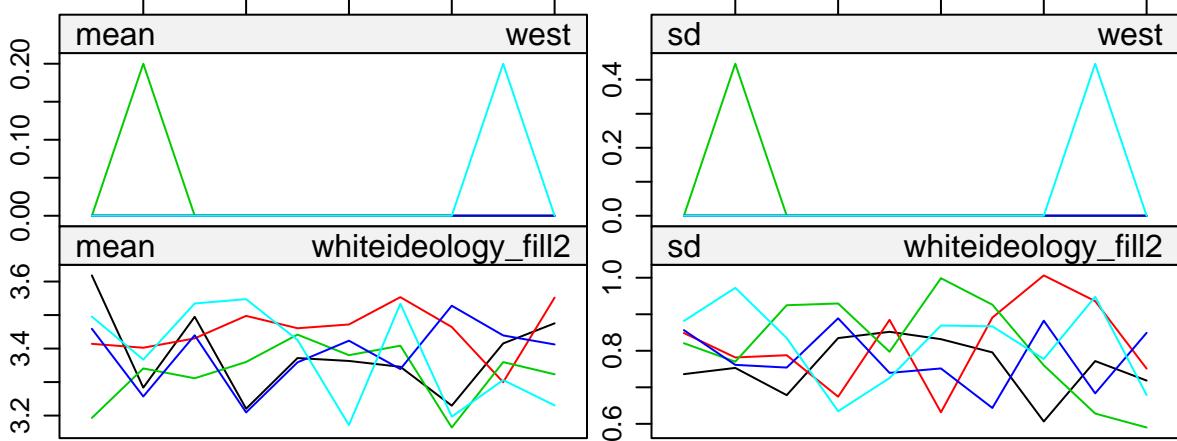
To better understand any potential patterns in missing data, we then plotted the pattern of missingness for only those variables missing values. From the plot on the right, we can see that approximately 42% of observations are complete. There seem to be a correspondence between missing a value for biggest split and missing biracial. There also seem to be about 2% of values for which most of the variables are missing. However, most observations are not missing more than 2-3 values.

Finding no clear patterns in the missing data, I next performed multiple imputations (with 10 iterations) on the dataset. A non-stochastic imputation method, Classification and Regression Trees (CART), was used instead of the default because of an error with matrix inversion caused by the data. Before examining the results of Trounstein's model using the imputed data, I first run some diagnostic tests of the imputation results to make sure that everything is running as expected.

First, I check the convergence of the algorithm used within mice() for each of the variables. For the most part, the fits intertwine and do not exhibit any trends at later iterations, as desired. There are some issues with the results for some of the variables with very few missing values (such as West and South), which makes sense as we would expect the mean to be less reliable due to the law of large numbers, and convergence is more difficult. Nevertheless, since there are so few of these values in the actual dataset, this is not a major concern.

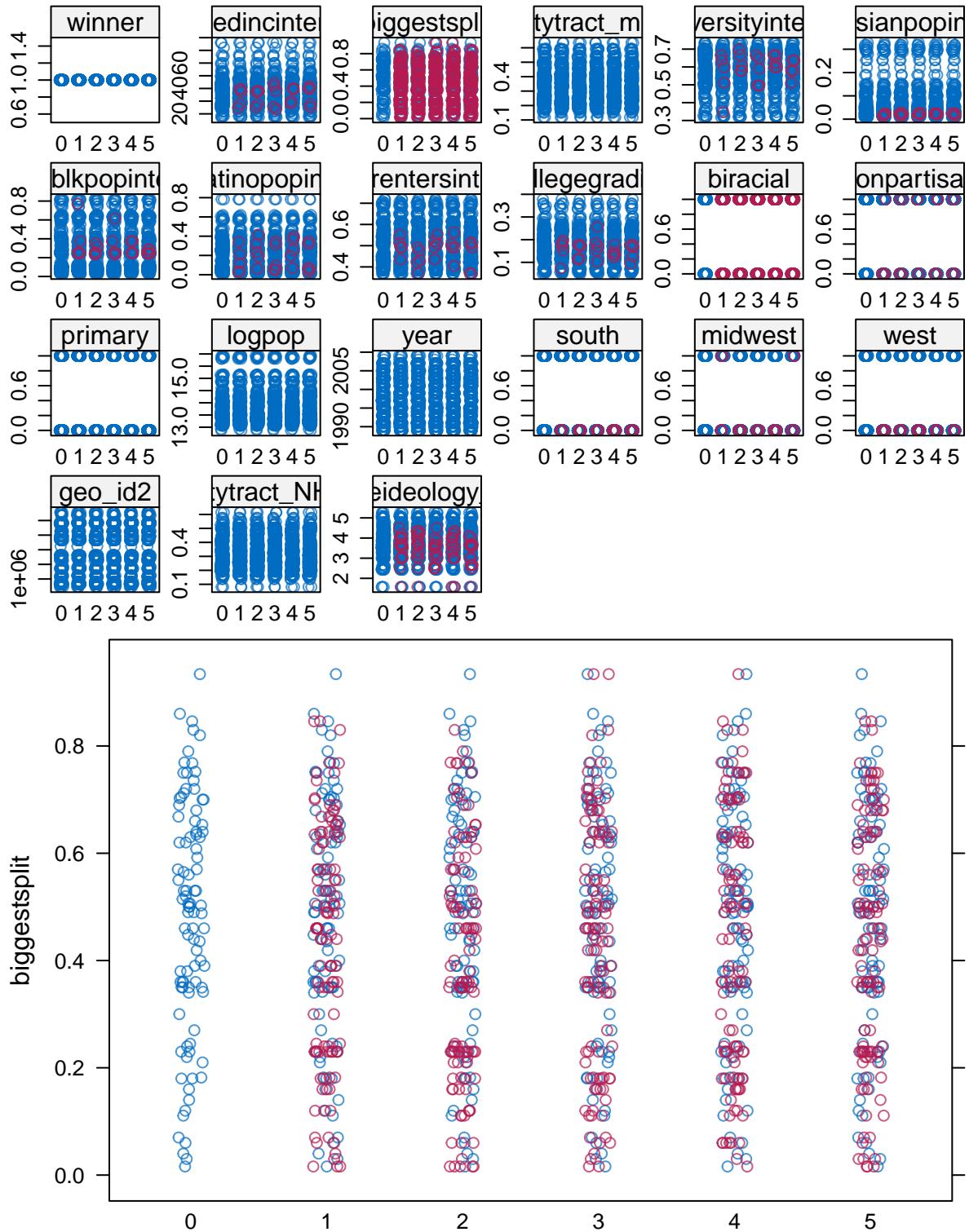






### Iteration

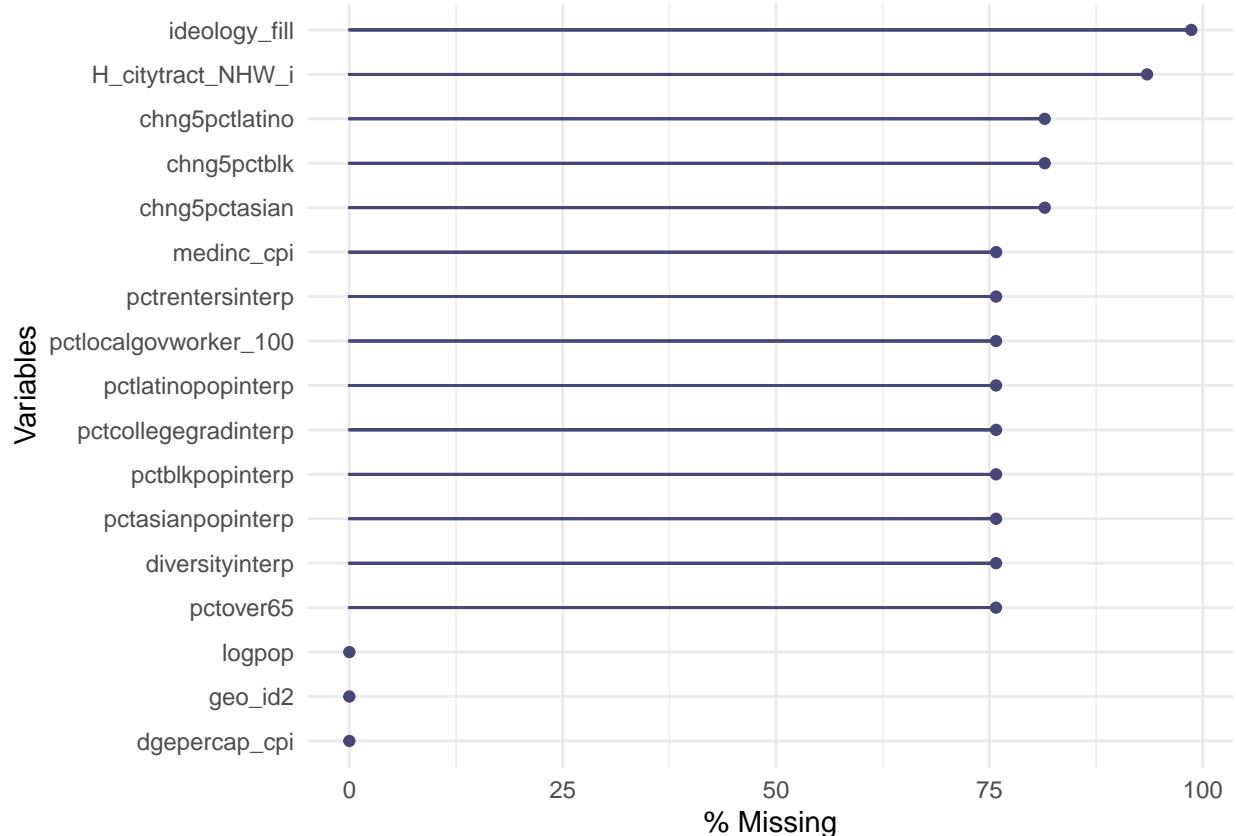
We can also check the imputed values against the original values using `stripplot()`. Each column in each subplot represents a separate iteration. The magenta points represent the imputed data. The values of the variable in question are along the y axis. We expect the spread of the data to be similar if the imputations were done well. If the data were missing completely at random, then the imputed data should have the same distribution as the original data. In particular, we want to be sure that the imputations are within a plausible range of the data. This is the case for all of our imputed variables, and there does not seem to be cause for alarm from these results.



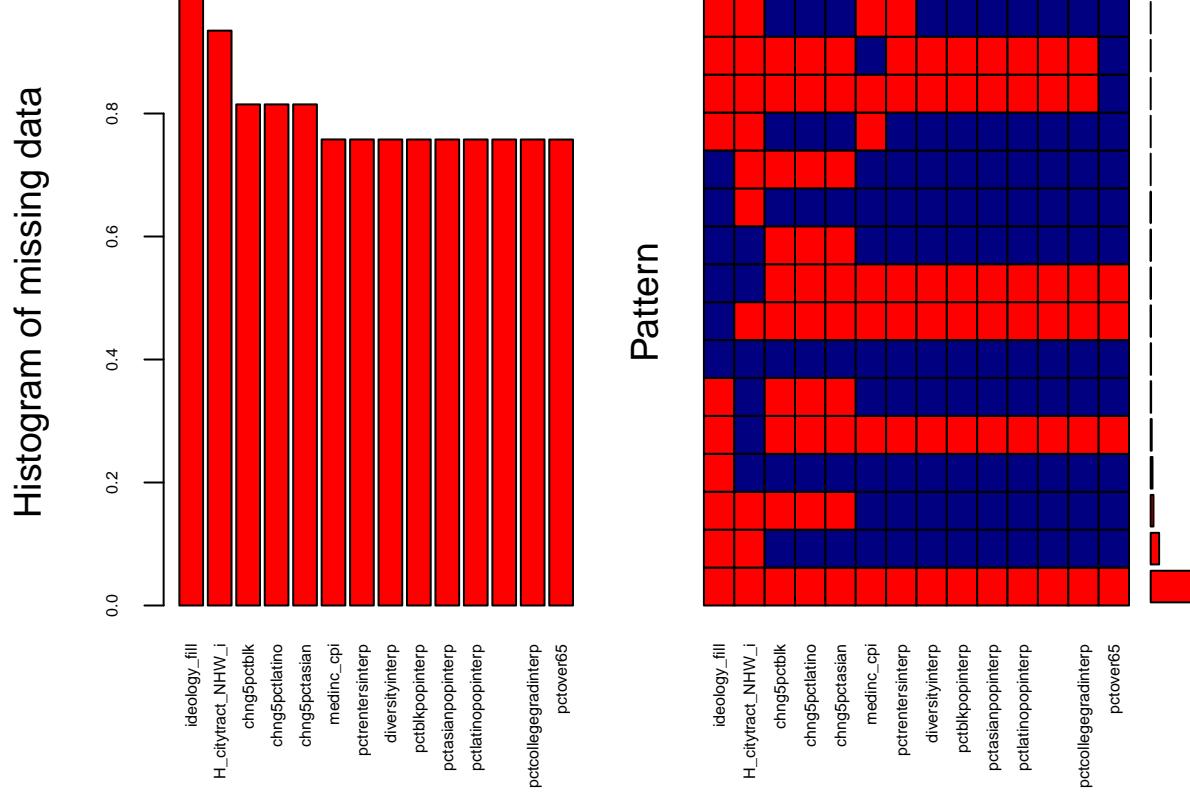
Finally, we can also look at the density plots for each variable's actual data and for their imputed data from each of the iterations, which are represented in magenta. Overall, the density plots align quite well for the variables with the most missing data, `biracial` and `biggestsplit`, and relatively well for most of the other variables with missing data. Again, the fit is less good for variables with fewer missing data points.







Warning in plot.aggr(res, ...): not enough horizontal space to display frequencies



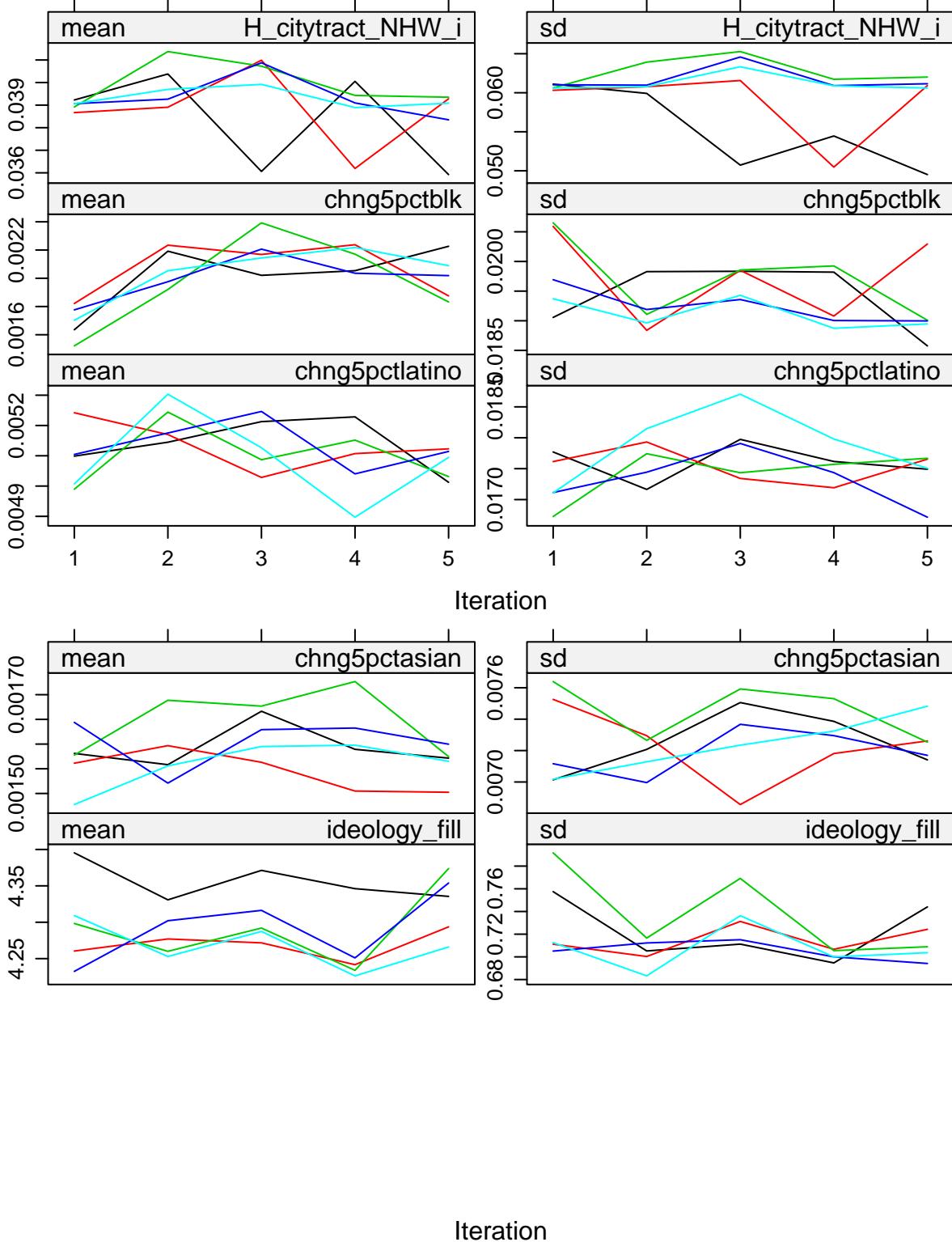
Variables sorted by number of missings:

Variable	Count
ideology_fill	0.9866041
H_citytract_NHW_i	0.9347787
chng5pctblk	0.8148547
chng5pctlatino	0.8148547
chng5pctasian	0.8148547
medinc_cpi	0.7578741
pctrentersinterp	0.7577980
diversityinterp	0.7577914
pctblkpopinterp	0.7577914
pctasianpopinterp	0.7577914
pctlatinopopinterp	0.7577914
pctlocalgovworker_100	0.7577914
pctcollegegradinterp	0.7577914
pctover65	0.7576722

[1] 1751.791

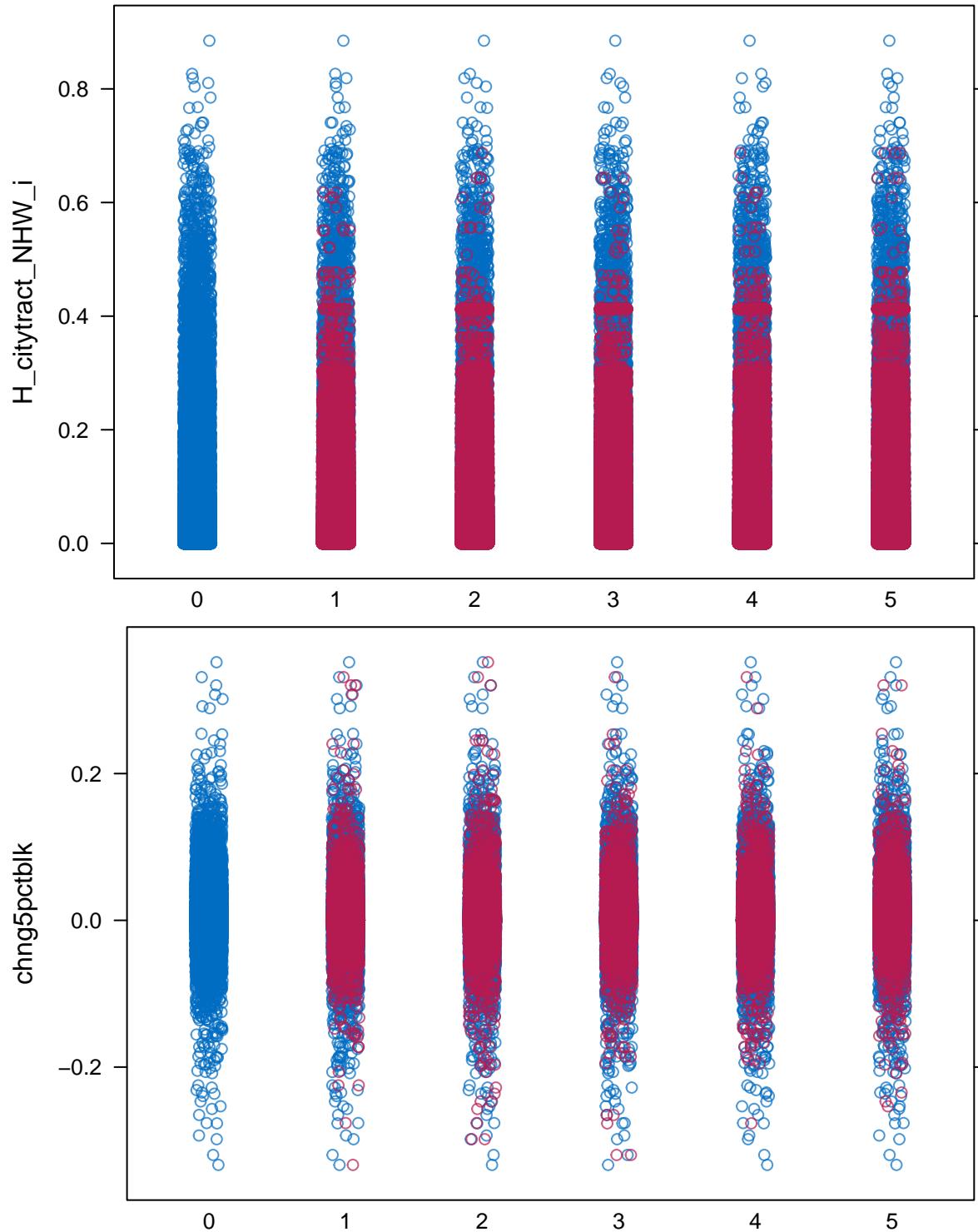
Finding no clear patterns in the missing data, I next performed multiple imputations (with 5 iterations) on the dataset. A non-stochastic imputation method, Classification and Regression Trees (CART), was used instead of the default because of an error with matrix inversion caused by the data. Before examining the results of Trounstein's model using the imputed data, I first run some diagnostic tests of the imputation results to make sure that everything is running as expected.

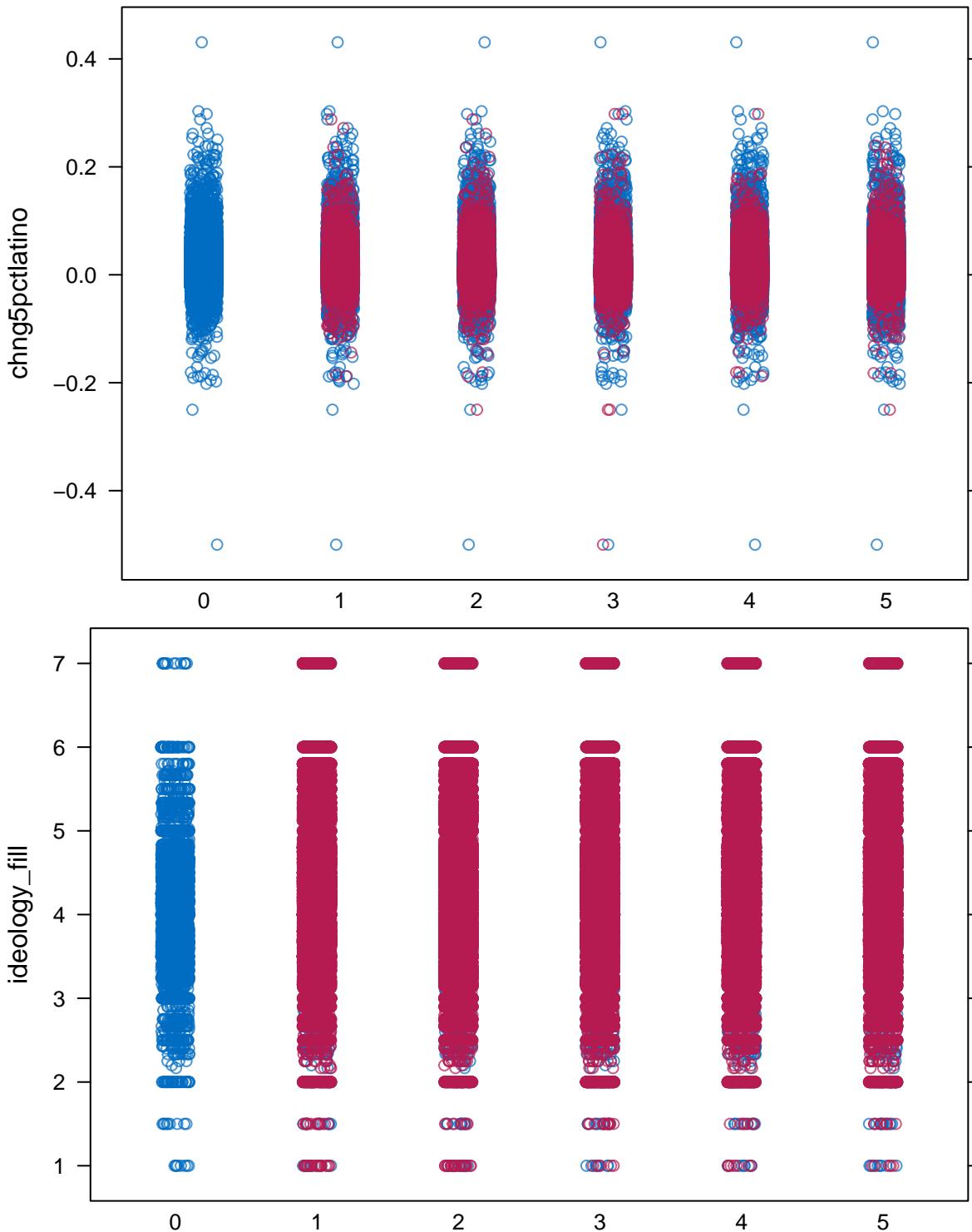
First, I check the convergence of the algorithm used within mice() for each of the variables. For the most part, the fits intertwine and do not exhibit any trends at later iterations, as desired.



We can also check the imputed values against the original values using `stripplot()`. Each column in each subplot represents a separate iteration. The magenta points represent the imputed data. The values of the variable in question are along the y axis. We expect the spread of the data to be similar if the imputations were done well. If the data were missing completely at random, then the imputed data should have the same distribution as the original data. In particular, we want to be sure that the imputations are within a plausible

range of the data. This is the case for all of our imputed variables, and there does not seem to be cause for alarm from these results.





Finally, we can also look at the density plots for each variable's actual data and for their imputed data from each of the iterations, which are represented in magenta. Overall, the density plots align quite well for the variables.

```
Error in density.default(x = c(NA_real_, NA_real_, NA_real_, NA_real_, : need at least 2 points to sele
```



```

6      chng5pctlatino -0.77989334 0.767938892 3.123802e-01
7      chng5pctasian  2.20234074 2.041358302 2.809437e-01
8      medinc_cpi    0.00589425 0.006690293 3.783217e-01
9  pctlocalgovworker_100 0.03058461 0.035468413 3.885312e-01
10     pctrentersinterp 0.14369964 0.333230292 6.663049e-01
11     pctover65       0.53679998 0.774926779 4.885004e-01
12  pctcollegegradinterp 3.41945066 0.819127369 3.000849e-05
13      logpop        -0.58243971 0.146434834 6.993432e-05

              term      estimate   std.error   p.value
1  H_citytract_NHW_i  0.143024472 0.399318809 7.249954e-01
2  diversityinterp   0.751188208 0.292782206 1.030571e-02
3  pctblkpopinterp   0.488306559 0.348789993 1.615307e-01
4  pctasianpopinterp 0.978740038 1.259502848 4.371195e-01
5  pctlatinopopinterp 1.579793555 0.437055242 3.016405e-04
6      medinc_cpi    0.005646432 0.006682476 3.981442e-01
7  pctlocalgovworker_100 0.031321702 0.035553248 3.783407e-01
8  pctrentersinterp   0.084033205 0.337804997 8.035511e-01
9  pctover65          0.603482776 0.776052959 4.367972e-01
10  pctcollegegradinterp 3.133684615 0.784339933 6.487400e-05
11      logpop        -0.588359060 0.146897573 6.221281e-05
12  ideology_fill    -0.012129920 0.025925864 6.533201e-01

```