# Synthetic Data Description

*Robert Arbon*

*08/04/2019*

## Introduction

In order to synthesize the data we are using an R package called *SynthPop*. The vignette is extremely informative and gives many good references. The documentation can be found here.

The method is very similar to Multiple Imputation by Chained Equations (MICE) except that missing values are treated as distinct and valid value (rather than something to be imputed). Please refer to the vignette section 2 for a description of the method.

## Our synthesis method

The method we have employed is as follows:

1. We use a Classification and Regression Tree (i.e. using `syn(method=cart)`) to build the conditional distributions for each variable, with the following parameters:
2. we used proper synthesis, i.e. `proper=TRUE`
3. smoothing was applied to the continuous variables: `iq`, `height_16` and `weight_16`.
4. only one replicant was produced, `m=1`.
5. The random seed was `seed = 8308313`.
6. The visit sequence was just the order of the dataframe.
7. All ther parameters were set to their defaults.
8. Once this was done we removed any observations that were in the original dataset using `anti_join` from `dplyr`.

For transparency we'll go through this method here:

### Synthesize data:

The following script was run to synthesize the data on a local compute cluster. A random seed was supplied as an argument:

```
library(synthpop)
library(tidyverse)

# Take a random seed from cluster job ID.  Only the first part is an integer.
args <- commandArgs(trailingOnly=TRUE)
my_seed <- strsplit(args[1], "\\.")[[1]][1]
my_seed <- as.integer(my_seed)
print(my_seed)

# Load data
df <- read_rds('data/maps-original.rds')

# Synthesize
df_syn <- syn(df, method="cart", m = 1,proper=T, smoothing=list("iq"="density", "height_16"="density",
```

```
# Write out
write_rds(x=df_syn, path=paste('output/cart/maps-synds-', my_seed, '.rds', sep=""))
```

## Load data

At this point the above script has been run and the output downloaded.

```
library(tidyverse)
```

```
## -- Attaching packages ----------------------------------------------------- tidyverse 1.2.1 --
```

```
## v ggplot2 3.0.0     v purrr   0.2.5
## v tibble  1.4.2     v dplyr   0.7.6
## v tidyr   0.8.1     v stringr 1.3.1
## v readr   1.1.1     v forcats 0.3.0
```

```
## -- Conflicts ------------------------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
library(synthpop)
```

```
## Loading required package: lattice
```

```
## Loading required package: MASS
```

```
##
## Attaching package: 'MASS'
```

```
## The following object is masked from 'package:dplyr':
##
##     select
```

```
## Loading required package: nnet
```

```
sds <- read_rds('output/cart_synth_ds/maps-synds-8308694.rds')
ods <- read_rds('data/maps-original.rds')
```

We can check the random seed and the call signature to make sure we've got the correct data.

```
sds$seed
```

```
## [1] 8308694
```

```
sds$call
```

```
## syn(data = df, method = "cart", m = 1, proper = T, smoothing = list(iq = "density",
##     height_16 = "density", weight_16 = "density"), seed = my_seed)
```

## Drop duplicates

The `sds` object is a SynthPop object which contains the synthetic dataframe amongst all the other associated (meta)data. Here we use an `anti_join` to select those rows of the synthetic dataframe that are an exact match to those in the original dataframe.

As this is perhaps not very intuitive, let's have an simple example:

```
table1<-data.frame(id=seq(1,6), vals=rnorm(6))
table2<-table1[c(1,3,4),]
table1
```

```
##   id        vals
## 1  1  0.01991039
## 2  2  1.78028122
## 3  3 -0.45165838
## 4  4  0.24391266
## 5  5  0.67606484
## 6  6  0.46407175
```

table2

```
##   id        vals
## 1  1  0.01991039
## 3  3 -0.45165838
## 4  4  0.24391266
```

anti_join(table1, table2)

```
## Joining, by = c("id", "vals")
```

```
##   id      vals
## 1  2 1.7802812
## 2  5 0.6760648
## 3  6 0.4640718
```

So we can use this to overwrite the synthetic dataset with one which has no rows in common with the original dataset:

```
# check the original dimensions
dim(sds$syn)
```

```
## [1] 14619    82
```

```
# Overwrite
sds$syn <- anti_join(sds$syn, ods, by=colnames(ods))
# Check dimensions of new dataset
dims <- dim(sds$syn)
# k is the number of observations
sds$k <- dims[1]
dims
```

```
## [1] 13734    82
```

# Comparisons

In this section we'll compare the dataset by running checking:

1. The number of complete cases on a selection of the variables (the outcomes and exposures)
2. Some logistic/multinomial regressions on some of the binary/ternary variables.
3. The marginal distributions of each variable (these are presented as charts and placed at the end)

## Complete cases

```
# The important variables
important_vars <- c('has_dep_diag', 'prim_diag', 'secd_diag', 'dep_thoughts',
                    'dep_score', 'comp_week', 'comp_wend')
```

```r
# This calculates the number of complete cases
n_comp_cases <- function(df){
  dims <- df %>%
      dplyr::select(important_vars) %>%
      dplyr::filter(complete.cases(.)) %>%
      dim()
  dims[1]
}

print(sprintf("Real DS complete cases: %4d", n_comp_cases(ods)))
```

```
## [1] "Real DS complete cases: 1206"
```

```r
print(sprintf("Synthetic DS complete cases: %d",n_comp_cases(sds$syn)))
```

```
## [1] "Synthetic DS complete cases: 1268"
```

## Regressions

We'll use the SynthPop functions to run and compare the regressions. First we'll regress the ICD diagnosis indicator against `sex`, `comp_week` and `comp_wend`.

```r
# We'll need these variables for the other regressions:
set.seed(42)
df <- sds$syn
col_names <- colnames(df)
n_cols <- length(col_names)
```

```r
# Fit the regression
mod <- glm.synds(formula=has_dep_diag ~ sex + comp_week + comp_wend,
                 family="binomial", data=sds)

# Get the comparison:
comp<- synthpop::compare(mod, ods, print.coef=F)
print(comp$ci.plot)
```

Z values for fit to has_dep_diag

Now we'll do some random regressions: 1. We'll randomly select columns, the first column will be the dependent variable. 2. Check whether the dependent variable has 2 or 3 unique levels (excluding NAs) 3. Check whether the covariates have 2 or more levels (excluding NAs). 4. Fit the model. 5. Plot the Z scores and CIs for each covariate on the same plot. 6. Repeat 10 times.

Note: Some comparisons won't work because different number of levels in the complete cases in original and synthetic data set. We'll ignore these.

First, here are some helper functions:

```r
# This tests whether the number of cases is greater than 2 for all variables.
# And whether the first variable has less than 4 cases (allowing for ternary outcomes)
is_valid_selection <- function(df, col_nums){
  n_vals <- sapply(col_nums, function(x) length(unique(df[,x])))
  flag <- (min(n_vals) > 2) & (n_vals[1] < 4)
  flag
}

# Creates a formula for a single col name for the DV, and a vector of col names for the IVs.
create_formula <- function(dv, ivs){
    form <- paste(sapply(ivs[2:length(ivs)], function(x) sprintf(' + %s', x)), collapse="")
    form <- paste(dv, sprintf('~ %s', ivs[1]), form)
    form <- as.formula(form)
    form
}
```

Now let's loop!

```r
count <- 1

while(count <= 10){

  # Select some random columns
  cols <- sample(n_cols, 5)
  # Check variables have more than 3 unique values and the first column has less than 4
  flag <- is_valid_selection(df, cols)

  if (flag){
    tryCatch({

      # Get the column names
      c_names <- sapply(cols, function(x) col_names[x])

      # Generate the formula for the regression
      form <- create_formula(c_names[1], c_names[2:length(c_names)])

      # Fit the regression
      mod <- glm.synds(formula=form, family="binomial", data=sds)

      # Get the comparison:
      comp<- synthpop::compare(mod, ods, print.coef=F)
      print(comp$ci.plot)
      count <- count + 1

    }, warning=function(e){
      # Some comparisons won't work because different number of complete cases in
      # original and synthetic data set.Just silently fail and move on.
    })
  }
}
```

Z values for fit to mat_anx_8m



Z values for fit to mat_anx_8m

7

Z values for fit to play_week



Z values for fit to has_dep_diag

Z values for fit to pat_pres_10



Z values for fit to emot_cruel

Z values for fit to mat_anx_1

Z values for fit to has_dep_diag

Z values for fit to tv_bed_9



Z values for fit to pat_pres

# Marginal distributions

Here are the comparisons of marginal distributions.

```r
for (colname in colnames(ods)){
  comparison <- synthpop::compare(sds,ods,colname)
  plt <- comparison$plots
  tbl <- comparison$tables
  plot(plt)
}
```
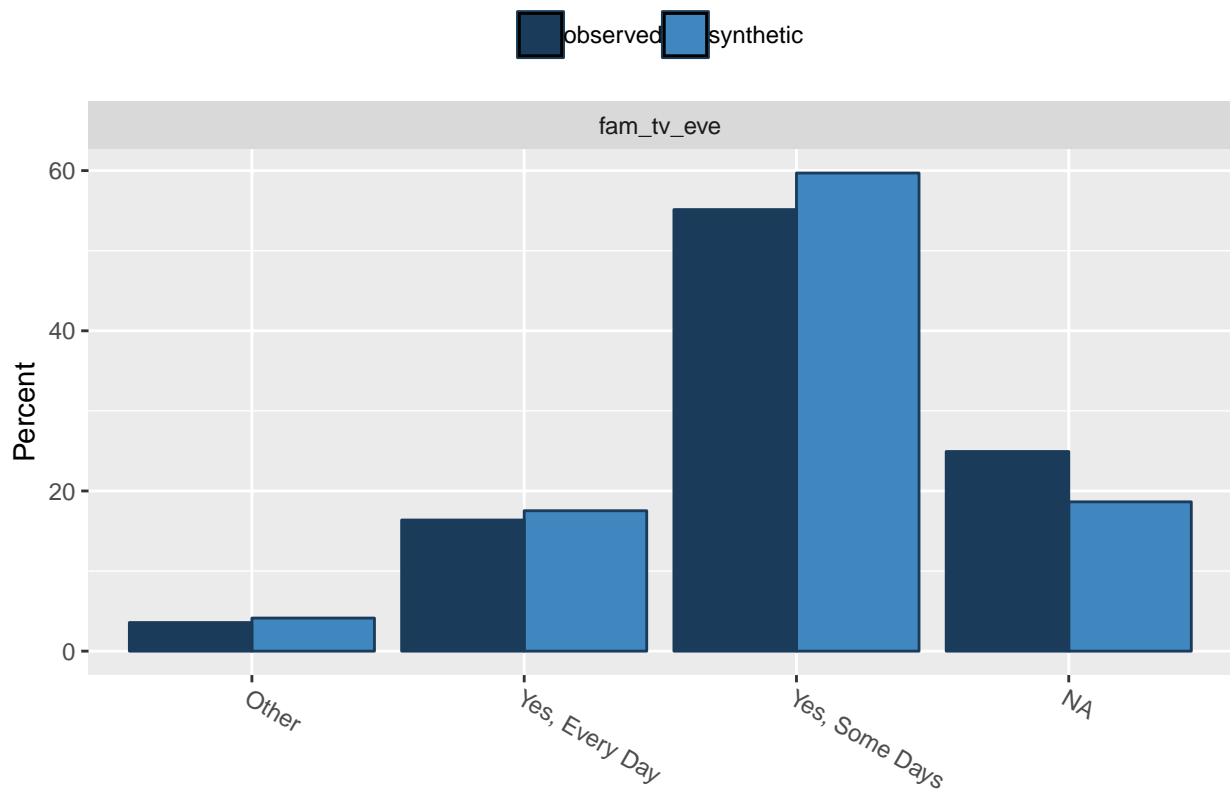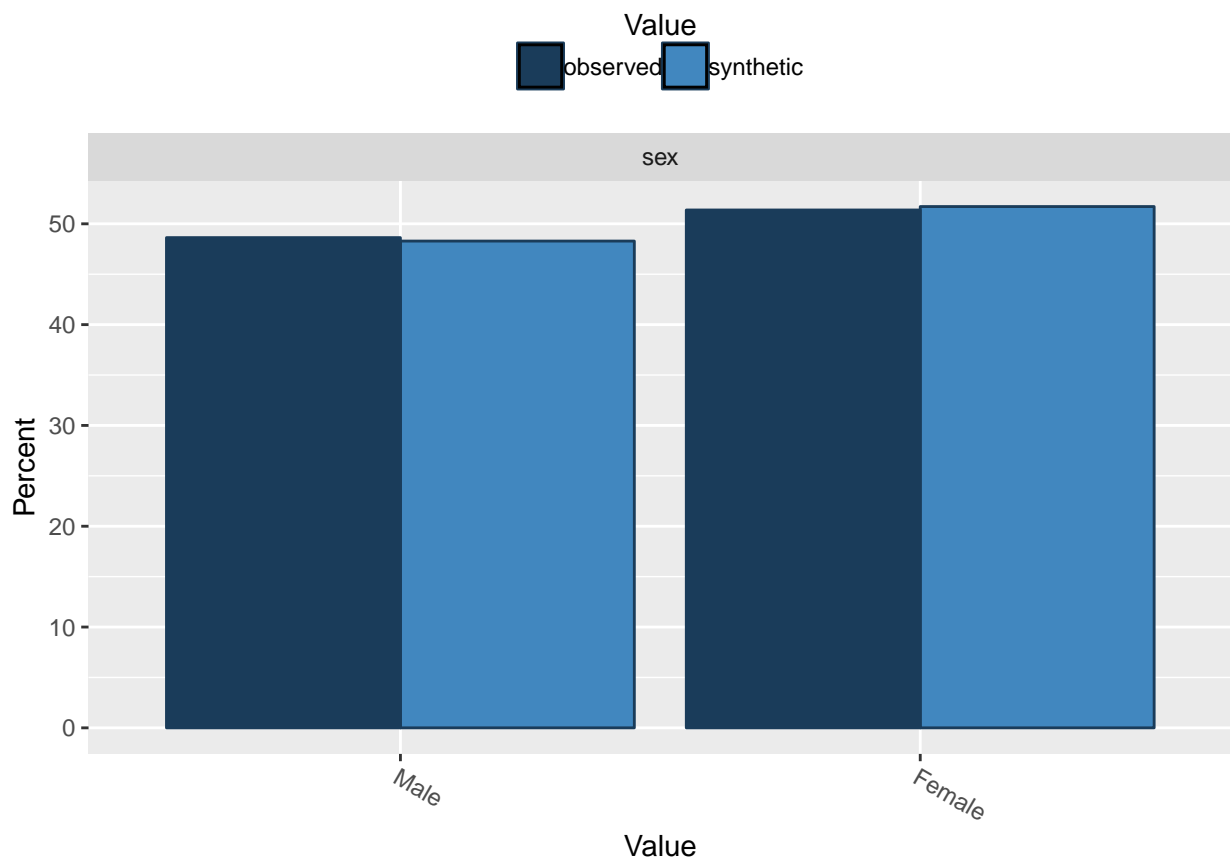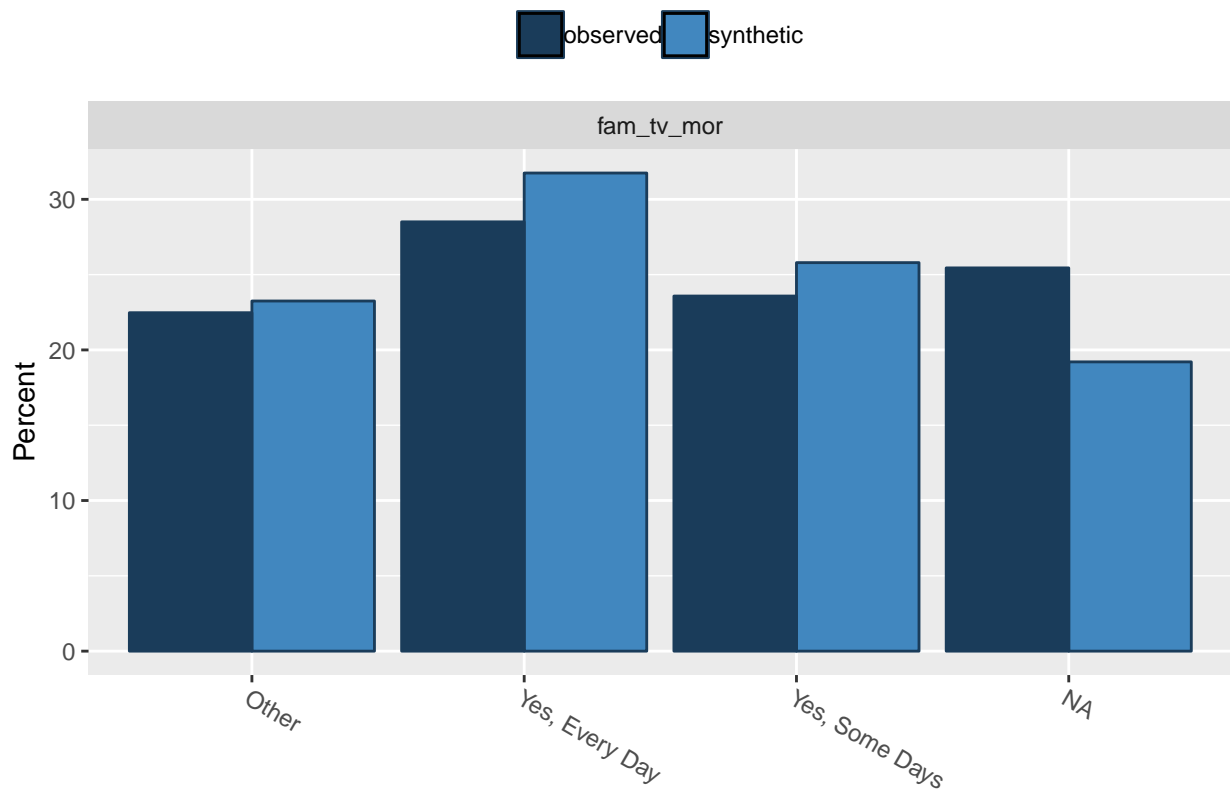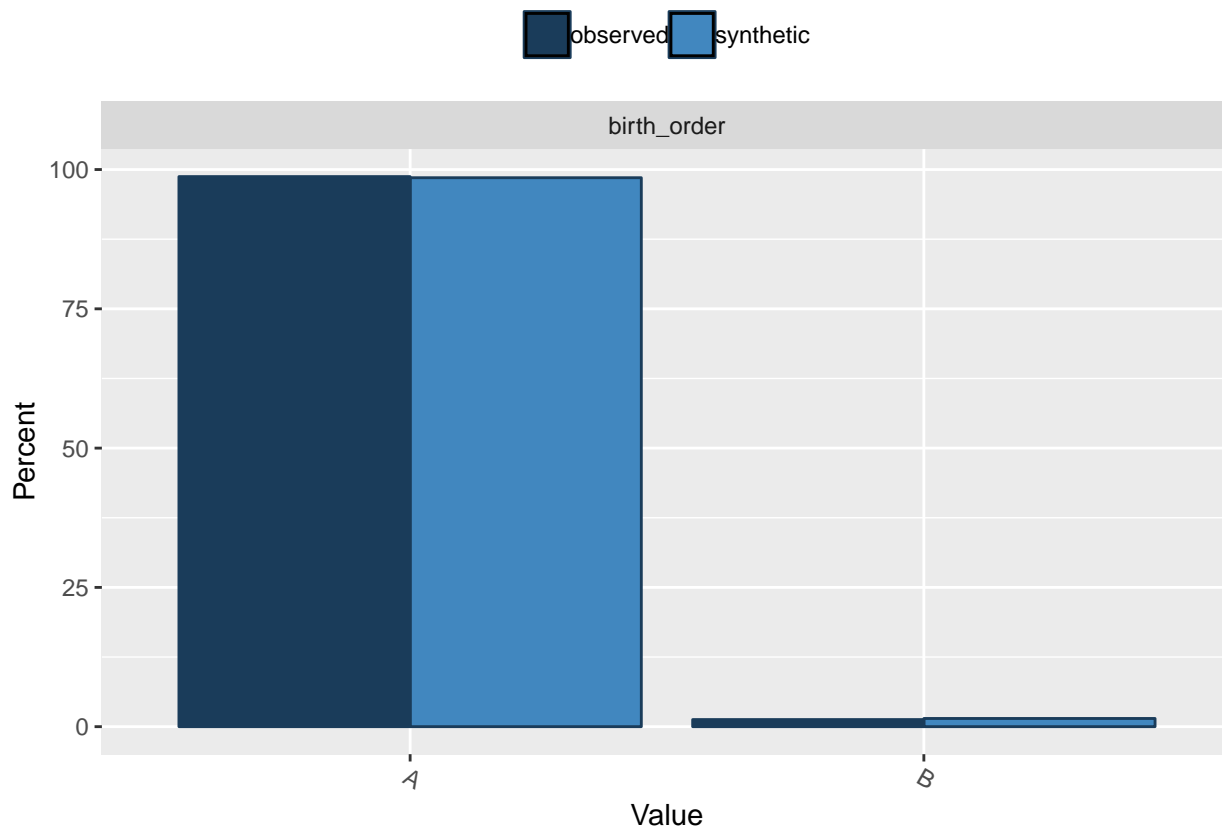
## Write synthetic data:

Finally we'll double check the data doesn't contain any unique rows from the original dataset, add a column indicating the data is synthetic and then write out as a csv.

To double check we can use the SynthPop function to remove any unique duplicates.

```
sds <- sdc(sds, ods, label='synthetic',  rm.replicated.uniques = T)
```

```
## no. of replicated uniques: 0
```

```
dim(sds$syn)
```

```
## [1] 13734    83
```

```
write.csv( sds$syn, 'output/maps-synthetic-data.csv')
```