

# Testing different imputation methods on PUMS (MCAR) - RandomForest

---

```
# load dataset: df
load('../Datasets/ordinalPUMS.Rdata')

# take 10,000 samples: df
n = 10000
sample <- sample(nrow(df), size = 10000)
df <- df[sample,]

# create MCAR scenario with 30% chance of missing: df_observed
set.seed(0)
missing_prob = 0.3
df_observed <- df
missing_col = colnames(df)[c(1,3,5,7,9,11)]
for (col in missing_col) {
  missing_ind <- rbernoulli(n,p = missing_prob)
  df_observed[missing_ind, col] <- NA
}
```

## MICE-CART

```
library(mice)

##
## Attaching package: 'mice'
## The following objects are masked from 'package:base':
##
##   cbind, rbind

df.imp <- missForest(df_observed, verbose = TRUE)

## missForest iteration 1 in progress...done!
## estimated error(s): 0.3092895
## difference(s): 0.1044364
## time: 6.057 seconds
##
## missForest iteration 2 in progress...done!
## estimated error(s): 0.262678
## difference(s): 0.03841818
## time: 5.538 seconds
##
## missForest iteration 3 in progress...done!
## estimated error(s): 0.2576005
## difference(s): 0.02488182
## time: 5.55 seconds
##
## missForest iteration 4 in progress...done!
## estimated error(s): 0.2556163
```

```
##      difference(s): 0.02058182
##      time: 5.347 seconds
##
##      missForest iteration 5 in progress...done!
##      estimated error(s): 0.2546922
##      difference(s): 0.01843636
##      time: 5.335 seconds
##
##      missForest iteration 6 in progress...done!
##      estimated error(s): 0.2536701
##      difference(s): 0.01782727
##      time: 5.547 seconds
##
##      missForest iteration 7 in progress...done!
##      estimated error(s): 0.2524536
##      difference(s): 0.01752727
##      time: 5.348 seconds
##
##      missForest iteration 8 in progress...done!
##      estimated error(s): 0.2526706
##      difference(s): 0.01741818
##      time: 5.35 seconds
##
##      missForest iteration 9 in progress...done!
##      estimated error(s): 0.2538219
##      difference(s): 0.01768182
##      time: 5.522 seconds
```

```
imputed_df <- df.imp$ximp
```

Diagnostics

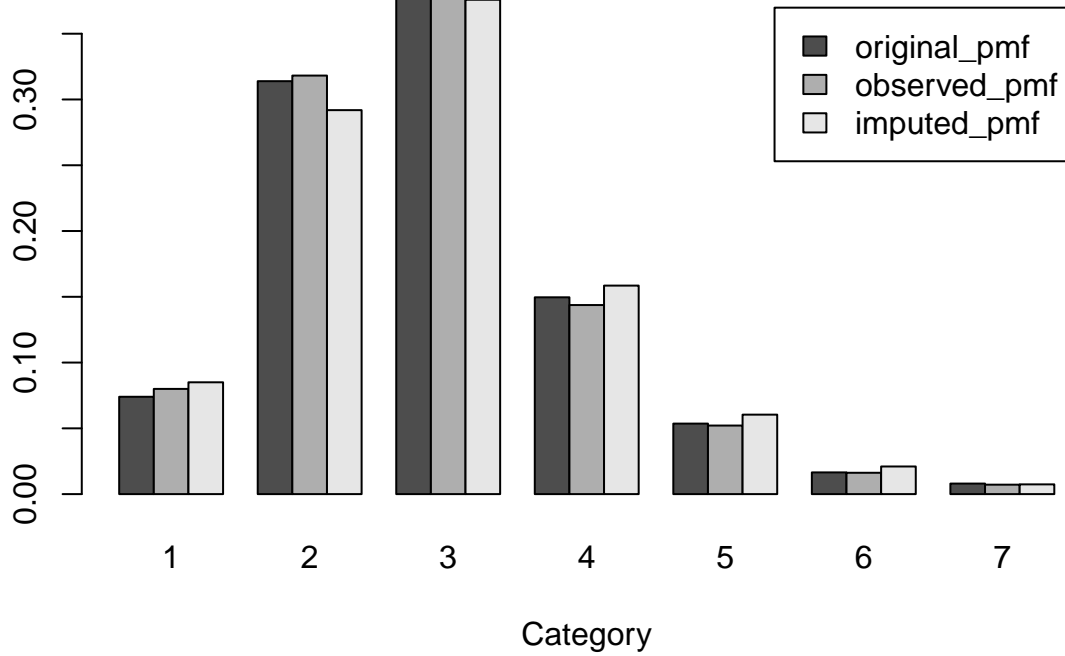
```
for (var_index in c(1,3,5,7,9,11)) {
  y_original = df[,var_index]
  original_pmf = table(y_original)/length(y_original)

  # Observed distribution
  missing_indicator = is.na(df_observed)[,var_index]
  y_observed = y_original[!missing_indicator]
  observed_pmf = table(y_observed)/length(y_observed)

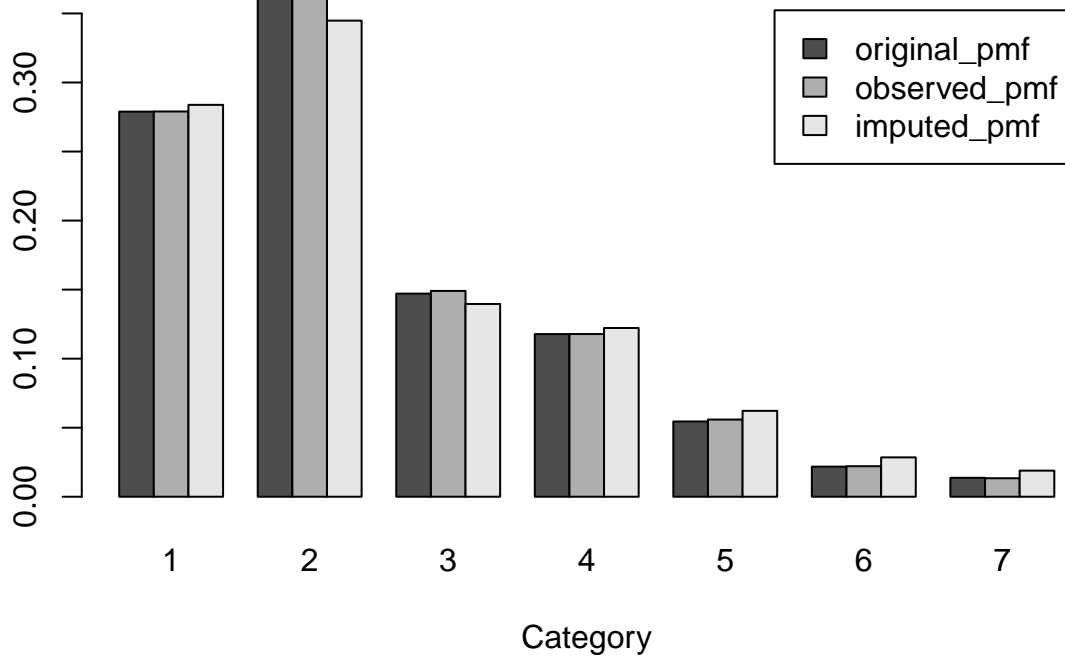
  # Marginal distribution after imputation
  imputed_pmf = table(imputed_df[, var_index])/sum(table(imputed_df[, var_index]))

  results = rbind(original_pmf,observed_pmf,imputed_pmf)
  colnames(results)<- 1:dim(imputed_pmf)
  barplot(results, xlab = 'Category', beside = TRUE,
          legend = TRUE,
          main = paste('Blocked Gibbs Sampling Assessment:', colnames(df)[var_index]))
}
```

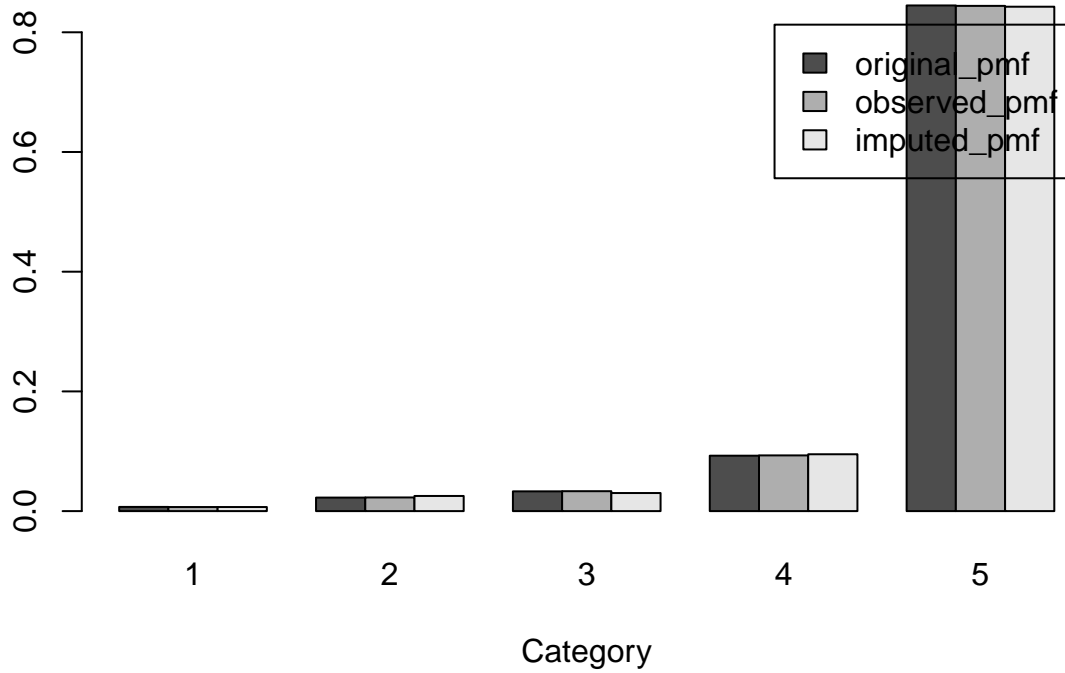
### Blocked Gibbs Sampling Assessment: VEH



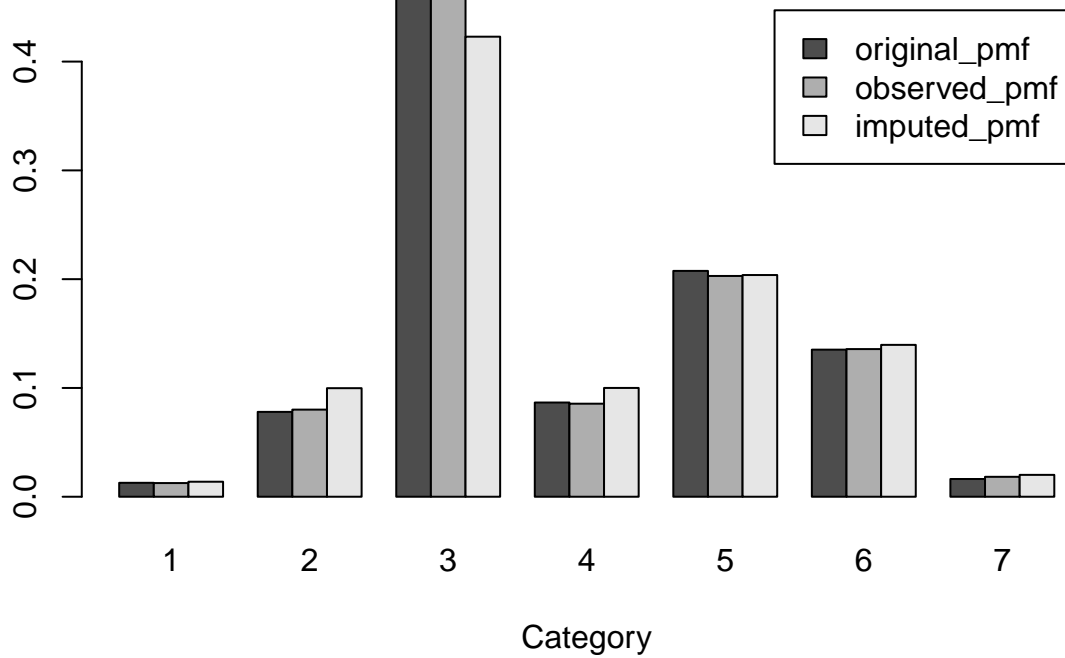
### Blocked Gibbs Sampling Assessment: NP



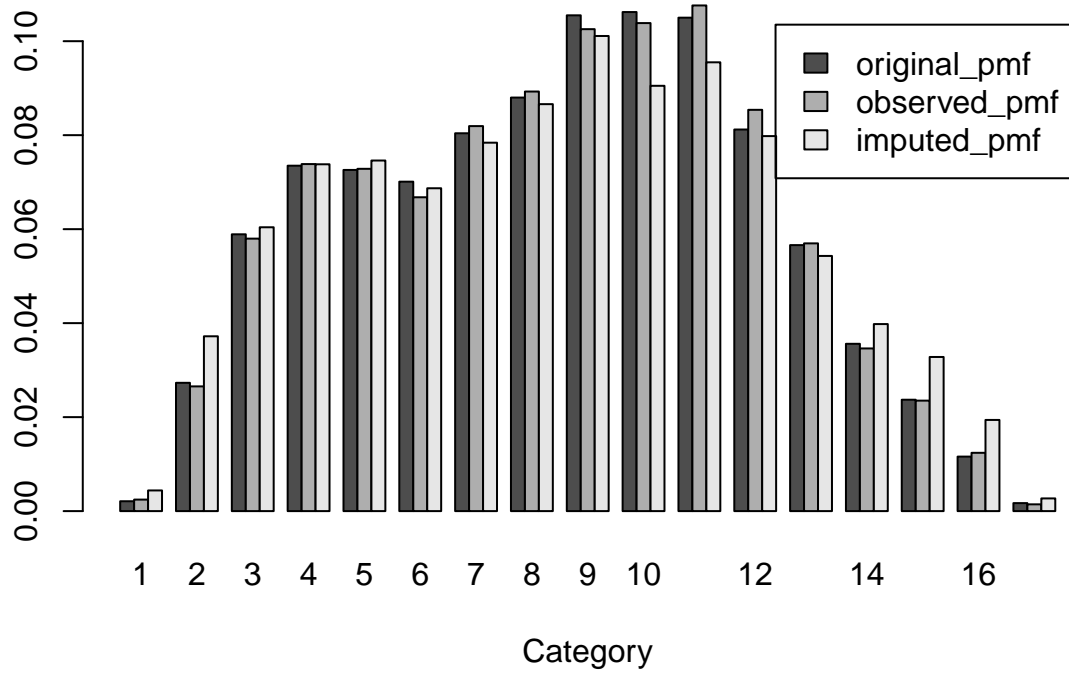
### Blocked Gibbs Sampling Assessment: ENG



### Blocked Gibbs Sampling Assessment: SCHL



### Blocked Gibbs Sampling Assessment: AGEF



### Blocked Gibbs Sampling Assessment: PINCP

