

Simple Missing Data Treatments

Utrecht University Winter School: Missing Data in R



**Utrecht
University**

Kyle M. Lang

Department of Methodology & Statistics
Utrecht University

Outline

Bad Methods

- Deletion-Based Methods

- Deterministic Imputation Methods

OK Methods



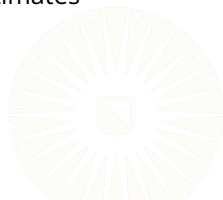
Bad Methods (These almost never work)

Listwise Deletion (Complete Case Analysis)

- Use only complete observations for the analysis
 - Very wasteful (can throw out lots of useful data)
 - Loss of statistical power

Pairwise Deletion (Available Case Analysis)

- Use only complete pairs of observations for analysis
 - Different samples sizes for different parameter estimates
 - Can cause computational issues



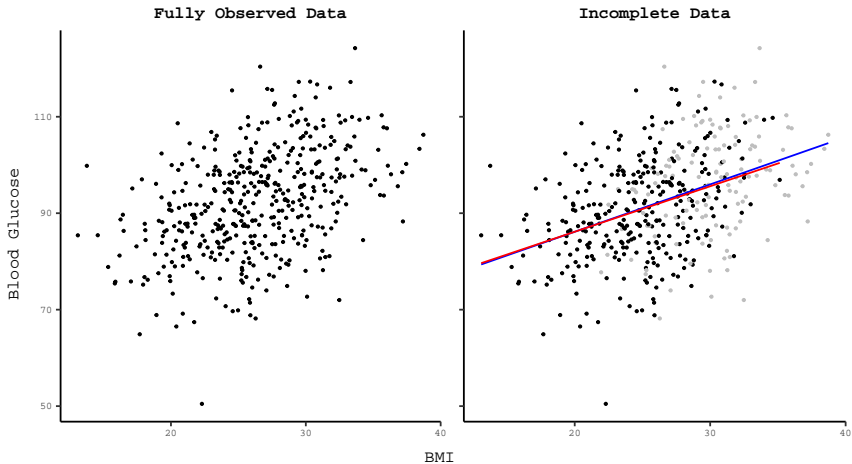
Example

```
## Read in some example data:
dat0 <- dat1 <- readRDS(paste0(dataDir, "diabetes_norm.rds"))

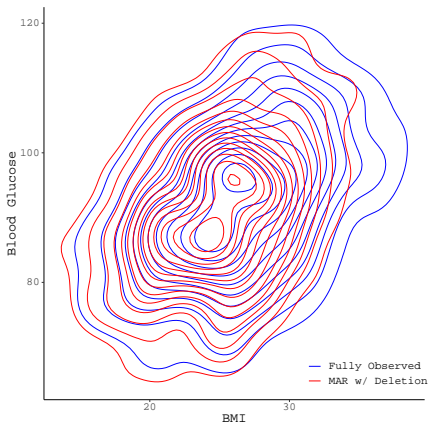
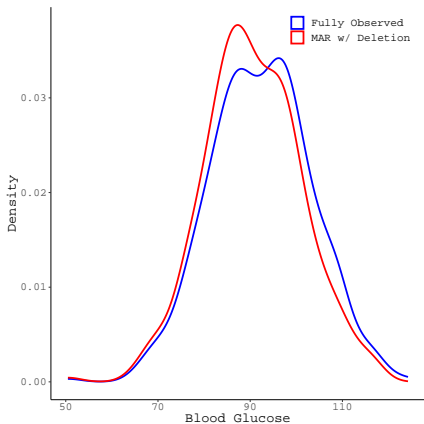
## Simulated missingness based on 'bmi':
m <- simLinearMissingness(data = dat1,
                           pm    = 0.30,
                           preds = "bmi",
                           auc   = 0.85)$r

## Impose missing on 'glu' according to the missingness above:
dat1[m, "glu"] <- NA
```

Example



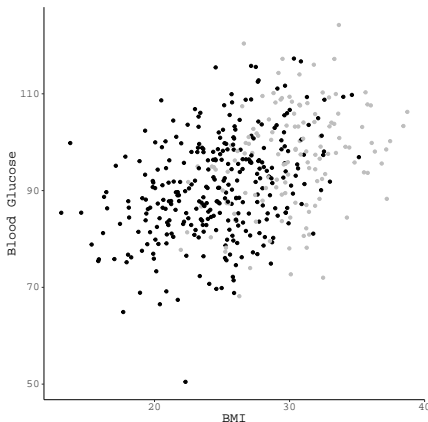
Example



Bad Methods (These almost never work)

(Unconditional) Mean Substitution

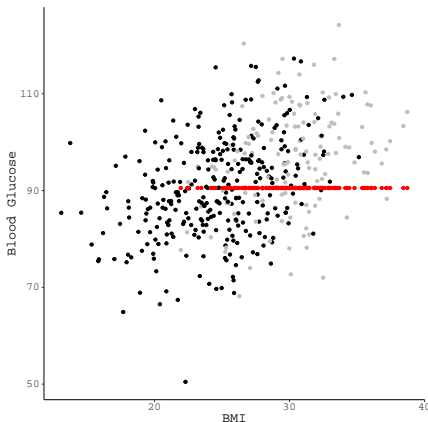
- Replace Y_{mis} with \bar{Y}_{obs}
 - Negatively biases regression slopes and correlations
 - Attenuates measures of linear association



Bad Methods (These almost never work)

(Unconditional) Mean Substitution

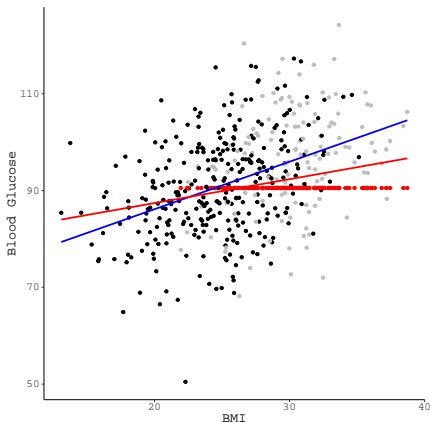
- Replace Y_{mis} with \bar{Y}_{obs}
 - Negatively biases regression slopes and correlations
 - Attenuates measures of linear association



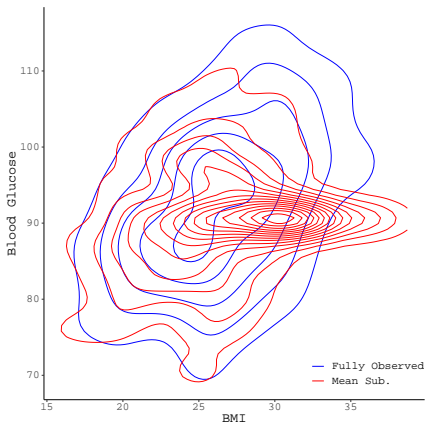
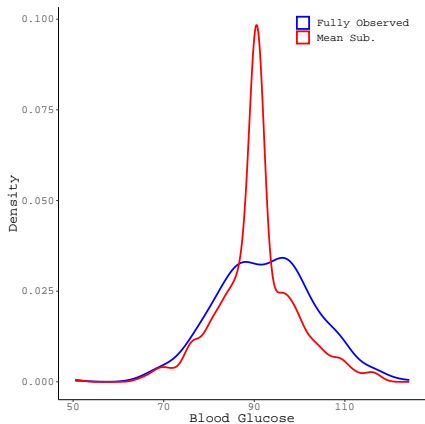
Bad Methods (These almost never work)

(Unconditional) Mean Substitution

- Replace Y_{mis} with \bar{Y}_{obs}
 - Negatively biases regression slopes and correlations
 - Attenuates measures of linear association



Example



Implementation

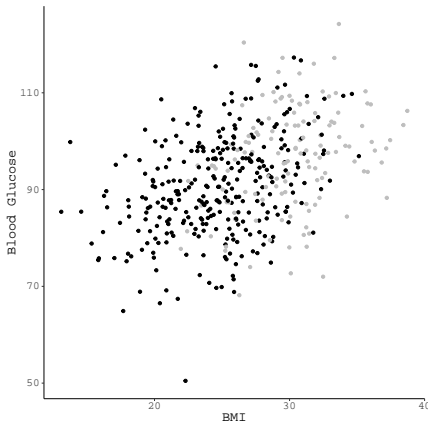
```
dat1[m, "glu"] <- mean(dat1$glu, na.rm = TRUE)

miceOut <- mice(data = dat1, m = 1, maxit = 1, method = "mean")
impData <- complete(miceOut, 1)
```

Bad Methods ('These almost never work')

Deterministic Regression Imputation (Conditional Mean Substitution)

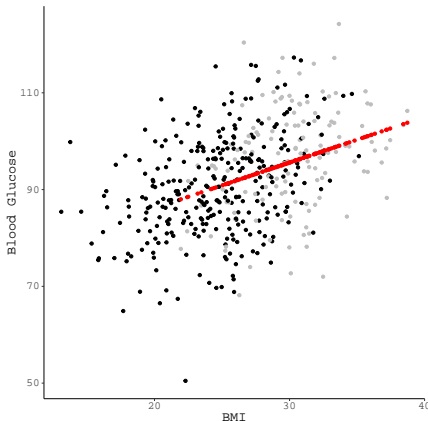
- Replace Y_{mis} with \hat{Y}_{mis} from some regression equation
 - Positively biases regression slopes and correlations
 - Inflates measures of linear association



Bad Methods (These almost never work)

Deterministic Regression Imputation (Conditional Mean Substitution)

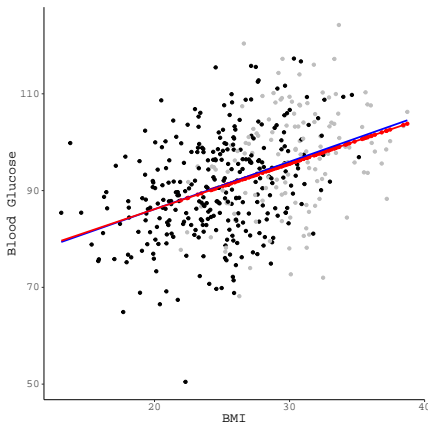
- Replace Y_{mis} with \hat{Y}_{mis} from some regression equation
 - Positively biases regression slopes and correlations
 - Inflates measures of linear association



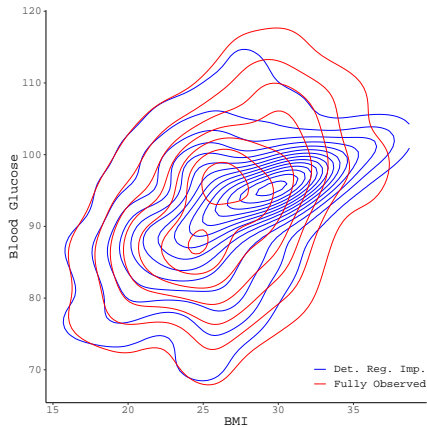
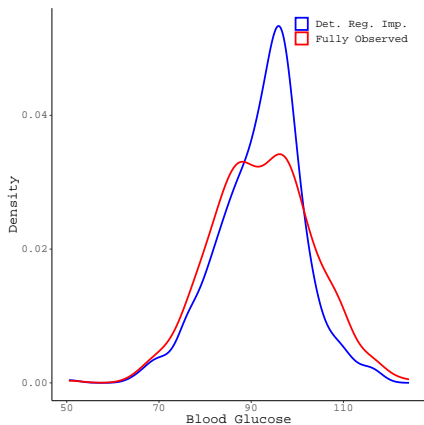
Bad Methods (These almost never work)

Deterministic Regression Imputation (Conditional Mean Substitution)

- Replace Y_{mis} with \hat{Y}_{mis} from some regression equation
 - Positively biases regression slopes and correlations
 - Inflates measures of linear association



Example



Implementation

```
miceOut <- mice(data = dat1, m = 1, method = "norm.predict")  
impData <- complete(miceOut, 1)
```



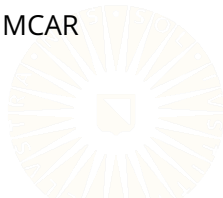
Bad Methods (These almost never work)

General Issues with Deletion-Based Methods

- Biased parameter estimates unless data are MCAR
- Generalizability issues

General Issues with Simple Single Imputation Methods

- Biased parameter estimates even when data are MCAR
- Attenuates variability in any treated variables



Bad Methods (These almost never work)

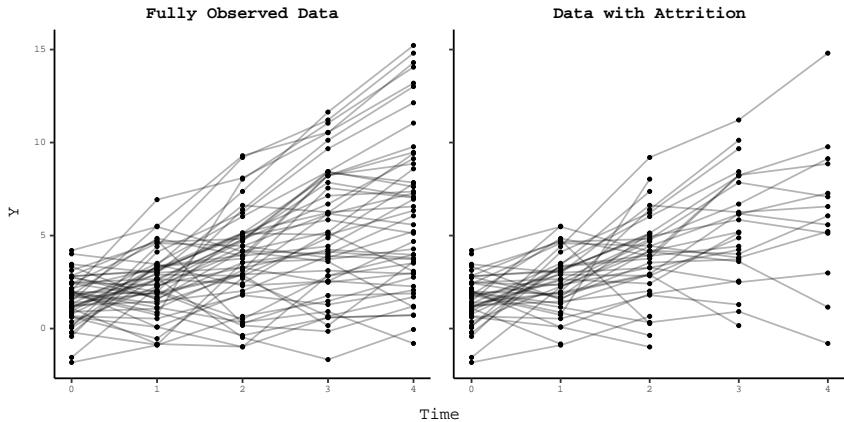
Averaging Available Items (Person-Mean Imputation)

- Compute aggregate scores using only available values
 - Missing data must be MCAR
 - Each item must contribute equally to the aggregate score

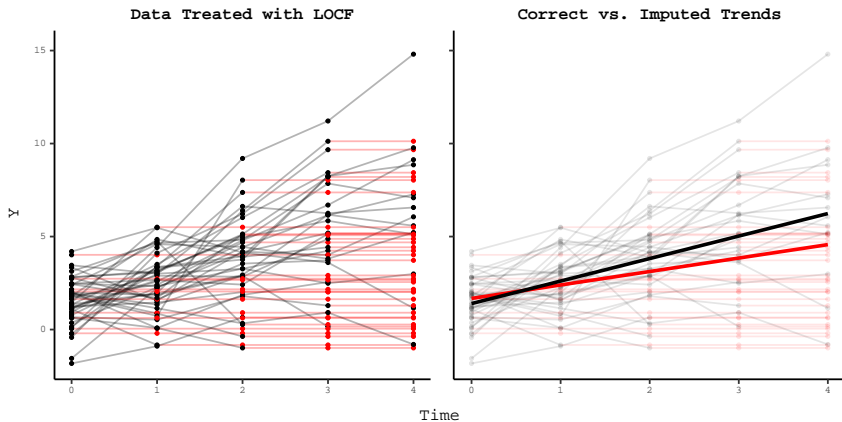
Last Observation Carried Forward (LOCF)

- Replace post-dropout values with the most recent observed value
 - Assume that dropouts would maintain their last known values
 - Attenuates estimates of growth/development

LOCF

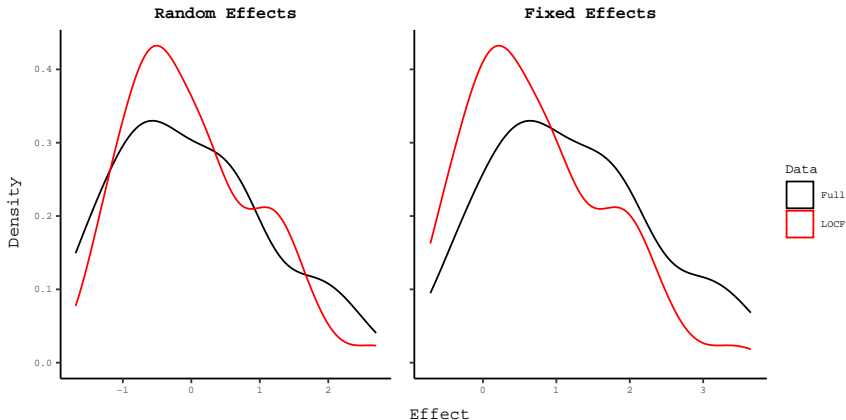


LOCF



Example

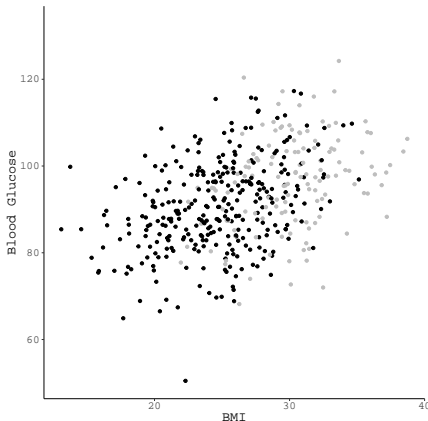
```
## Fit some multilevel regression models  
fit1 <- lmer(y ~ t + (t | id), data = fullData)  
fit2 <- lmer(y ~ t + (t | id), data = locfData)
```



OK Methods (These sometimes work)

Stochastic Regression Imputation

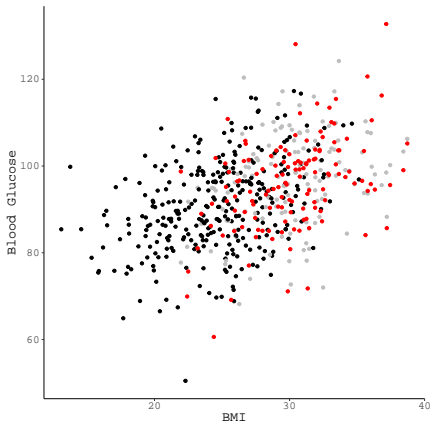
- Fill Y_{mis} with \hat{Y}_{mis} plus some random noise.
 - Produces unbiased parameter estimates and predictions
 - Computationally efficient
 - Attenuates standard errors
 - Makes CIs and prediction intervals too narrow



OK Methods (These sometimes work)

Stochastic Regression Imputation

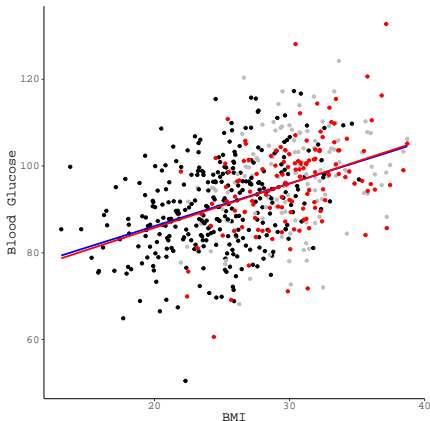
- Fill Y_{mis} with \hat{Y}_{mis} plus some random noise.
 - Produces unbiased parameter estimates and predictions
 - Computationally efficient
 - Attenuates standard errors
 - Makes CIs and prediction intervals too narrow



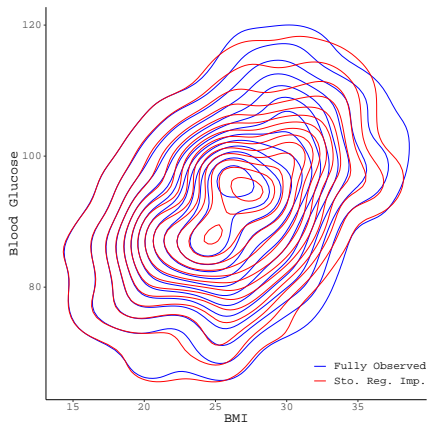
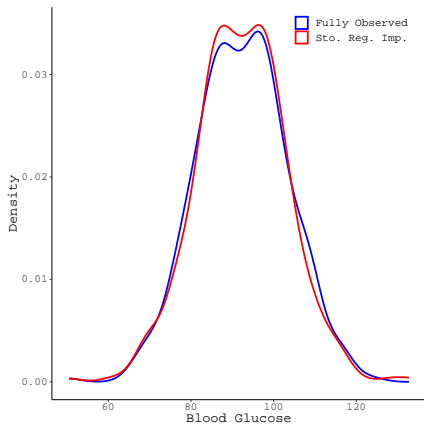
OK Methods (These sometimes work)

Stochastic Regression Imputation

- Fill Y_{mis} with \hat{Y}_{mis} plus some random noise.
 - Produces unbiased parameter estimates and predictions
 - Computationally efficient
 - Attenuates standard errors
 - Makes CIs and prediction intervals too narrow



Example



Implementation

```
miceOut <- mice(data = dat1, m = 1, seed = 42, method = "norm.nob")  
impData <- complete(1)
```



OK Methods (These sometimes work)

Nonresponse Weighting

- Weight the observed cases to correct for nonresponse bias
 - Popular in survey research and official statistics
 - Only worth considering with *Unit Nonresponse*
 - Doesn't make any sense with *Item Nonresponse*

