# 431 Class 25

thomaselove.github.io/431

2021-11-30

## Today's Agenda

1. What exactly is R doing if you ignore missing values when fitting models?
   - What does `type.convert()` do?
   - `na.omit` vs. `na.exclude` vs. `na.delete`
2. Use multiple imputation to deal with missing data in fitting a linear regression with `lm` using the `mice` package.

(MICE = Multiple Imputation through Chained Equations)

## My Setup

My R project for Class 25 is located in the following folder. . .

"unique_stuff_to_me/2021-431/431-classes/class25"

```
library(here)
```

```
here() starts at C:/Users/Thomas/Dropbox/2021-431/431-classes/
```

```
library(magrittr); library(knitr)
library(janitor); library(naniar)

library(mice)
  # mice = multiple imputation through chained equations

library(broom)
library(tidyverse)

theme_set(theme_bw())
```

**What happens if you fit a regression model without doing anything at all about missing data?**

## What happens if you ignore NAs?

Let's open a small, simulated data set with 100 subjects and some missing values.

```
sim1 <- read_csv(here("data", "c12_sim1.csv")) %>%
    type.convert(as.is = FALSE, na.strings = "NA")

head(sim1)
```

```
# A tibble: 6 x 6
  subject out_q out_b pred1 pred2 pred3
  <fct>   <dbl> <fct> <dbl> <dbl> <fct>
1 S001     81.1 Yes     8.8  20.5 Middle
2 S002    105.  No      7.1  24.9 High
3 S003      NA  <NA>    9.9  17.4 Middle
4 S004      NA  No      8.9  31.8 <NA>
5 S005     75.9 <NA>    NA   22    High
6 S006     79.8 No      9.7  NA   <NA>
```

# What does `type.convert()` actually do?

Tries to convert each column (individually) to either logical, integer, numeric, complex or (if a character vector) to factor.

- The first type (from that list) that can accept all non-missing values is chosen.
- If all of the values are missing, the column is converted to logical.
- Columns containing just F, T, FALSE, TRUE or NA values are made into logical.
- Use the `na.strings` parameter to add missing strings (default = "NA")
- `as.is = FALSE` converts characters to factors. `as.is = TRUE` will become the default soon, so you should specify.

# Our `sim1` data

| Variable | Description |
|----------|-------------|
| subject | Subject identifier |
| out_q | Quantitative outcome |
| out_b | Binary outcome with levels Yes, No |
| pred1 | Predictor 1 (quantitative) |
| pred2 | Predictor 2 (also quantitative) |
| pred3 | Predictor 3 (categories are Low, Middle, High) |

- Clean up the factors?
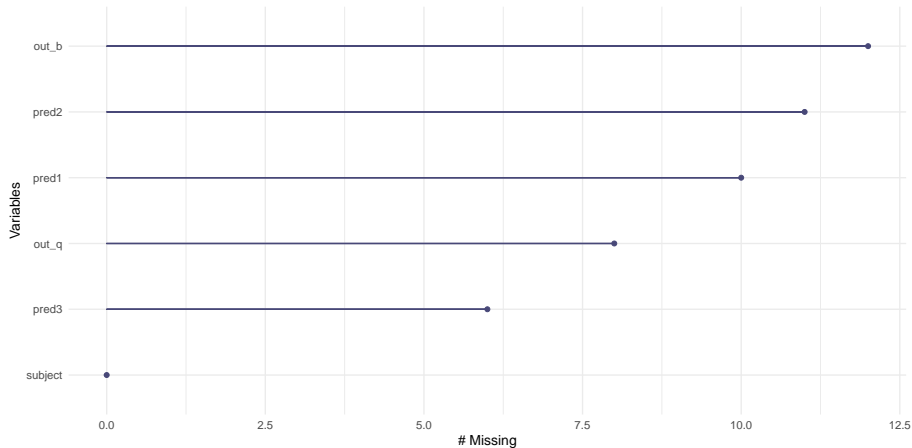
# Cleaning up `subject` and pred3

```
sim1 <- sim1 %>%
    mutate(subject = as.character(subject),
           pred3 = fct_relevel(pred3, "Low", "Middle"))

sim1 %>% tabyl(pred3, out_b)
```

```
  pred3 No Yes NA_
    Low 10  12   4
 Middle 12  17   4
   High 16  15   4
   <NA>  4   2   0
```

# How much missingness do we have?

```
gg_miss_var(sim1)
```

# How much missingness do we have?

```
miss_var_summary(sim1)

# A tibble: 6 x 3
  variable n_miss pct_miss
  <chr>     <int>    <dbl>
1 out_b        12       12
2 pred2        11       11
3 pred1        10       10
4 out_q         8        8
5 pred3         6        6
6 subject       0        0

n_miss(sim1)

[1] 47
```

# How much missingness do we have?

```
prop_complete_case(sim1)
```

```
[1] 0.65
```

```
miss_case_table(sim1)
```

```
# A tibble: 4 x 3
  n_miss_in_case n_cases pct_cases
           <int>   <int>     <dbl>
1              0      65        65
2              1      25        25
3              2       8         8
4              3       2         2
```

# Suppose we run a linear regression

without dealing with the missing data, so that we run:

```
mod1 <- lm(out_q ~ pred1 + pred2 + pred3, data = sim1)
```

How can we tell how many observations will be used?

# What happens when we run a regression model?

```
mod1 <- lm(out_q ~ pred1 + pred2 + pred3, data = sim1)

anova(mod1)

Analysis of Variance Table

Response: out_q
          Df   Sum Sq  Mean Sq F value Pr(>F)
pred1      1    209.8   209.81  0.5976 0.4423
pred2      1    132.1   132.14  0.3763 0.5417
pred3      2     86.5    43.24  0.1231 0.8843
Residuals 65 22821.9   351.11
```

- How many observations were used to fit this model?

# Summary of our linear model

```
> summary(mod1)

Call:
lm(formula = out_q ~ pred1 + pred2 + pred3, data = sim1)

Residuals:
    Min      1Q  Median      3Q     Max
-39.164 -13.900   2.419  15.541  34.156

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 105.2070    18.6185   5.651 3.82e-07 ***
pred1        -0.8361     1.3010  -0.643    0.523
pred2         0.2611     0.4614   0.566    0.573
pred3Middle  -1.3498     5.6802  -0.238    0.813
pred3High    -2.7443     5.5427  -0.495    0.622
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.74 on 65 degrees of freedom
  (30 observations deleted due to missingness)
Multiple R-squared:  0.01843,   Adjusted R-squared:  -0.04198
F-statistic: 0.3051 on 4 and 65 DF,  p-value: 0.8736
```

## Another way to see this

```
glance(mod1) %>% select(1:6)

# A tibble: 1 x 6
  r.squared adj.r.squared sigma statistic p.value    df
      <dbl>         <dbl> <dbl>     <dbl>   <dbl> <dbl>
1    0.0184       -0.0420  18.7     0.305   0.874     4

glance(mod1) %>% select(7:12)

# A tibble: 1 x 6
  logLik   AIC   BIC deviance df.residual  nobs
   <dbl> <dbl> <dbl>    <dbl>       <int> <int>
1  -302.  616.  629.   22822.          65    70
```

# How could we have known this would be 70, in advance?

```
sim1 %>% select(out_q, pred1, pred2, pred3) %>%
    miss_case_table()

# A tibble: 3 x 3
  n_miss_in_case n_cases pct_cases
           <int>   <int>     <dbl>
1              0      70        70
2              1      25        25
3              2       5         5
```

## Which observations were not used?

```
summary(mod1)$na.action
```

```
 3  4  5  6 13 16 19 26 27 29 30 34 39 48 51 56 62 66 67 68
 3  4  5  6 13 16 19 26 27 29 30 34 39 48 51 56 62 66 67 68
72 75 81 83 86 89 93 94 96 97
72 75 81 83 86 89 93 94 96 97
attr(,"class")
[1] "omit"
```

- A potentially more useful na.action setting in lm is na.exclude which pads out predicted values and residuals with NAs instead of omitting the 30 observations listed above.

```
lm(out_q ~ pred1 + pred2 + pred3,
      data = sim1, na.action = na.exclude)
```

# **Predictions from `mod1` with `na.omit` and `na.exclude`**

```
mod1 <- lm(out_q ~ pred1 + pred2 + pred3, data = sim1)
              ## note: by default na.action = na.omit here
head(predict(mod1))
```

```
        1         2         7         8         9        10
101.85279 103.02874  98.14391  96.57037 101.49208 101.01744
```

```
mod1_e <- lm(out_q ~ pred1 + pred2 + pred3, data = sim1,
            na.action = na.exclude)
head(predict(mod1_e))
```

```
       1        2        3        4        5        6
101.8528 103.0287       NA       NA       NA       NA
```

# Multiple Imputation: Potential and Pitfalls

# Sterne et al. 2009 *BMJ*

Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls

> *In this article, we review the reasons why missing data may lead to bias and loss of information in epidemiological and clinical research. We discuss the circumstances in which multiple imputation may help by reducing bias or increasing precision, as well as describing potential pitfalls in its application. Finally, we describe the recent use and reporting of analyses using multiple imputation in general medical journals, and suggest guidelines for the conduct and reporting of such analyses.*

- https://www.bmj.com/content/338/bmj.b2393

**Note**: The next 7 slides are derived from Sterne et al.

# An Example from Sterne et al.

Consider, for example, a study investigating the association of systolic blood pressure with the risk of subsequent coronary heart disease, in which data on systolic blood pressure are missing for some people.

The probability that systolic blood pressure is missing is likely to:

- decrease with age (doctors are more likely to measure it in older people),
- decrease with increasing body mass index, and
- decrease with history of smoking (doctors are more likely to measure it in people with heart disease risk factors or comorbidities).

If we assume that data are missing at random and that we have systolic blood pressure data on a representative sample of individuals within strata of age, smoking, body mass index, and coronary heart disease, then we can use multiple imputation to estimate the overall association between systolic blood pressure and coronary heart disease.

# Missing Data Mechanisms

- **Missing completely at random** There are no systematic differences between the missing values and the observed values.
  - For example, blood pressure measurements may be missing because of breakdown of an automatic sphygmomanometer.
- **Missing at random** Any systematic difference between the missing and observed values can be explained by other observed data.
  - For example, missing BP measurements may be lower than measured BPs but only because younger people more often have a missing BP.
- **Missing not at random** Even after the observed data are taken into account, systematic differences remain between the missing values and the observed values.
  - For example, people with high BP may be more likely to have headaches that cause them to miss clinic appointments.

"Missing at random" is an **assumption** that justifies the analysis, and is not a property of the data.

# Trouble: Data missing not at random

Sometimes, it is impossible to account for systematic differences between missing and observed values using the available data.

- In such (MNAR) cases, multiple imputation may give misleading results.
    - Those results can be either more or less misleading than a complete case analysis.
- For example, consider a study investigating predictors of depression.
    - If individuals are more likely to miss appointments because they are depressed on the day of the appointment, then it may be impossible to make the MAR assumption plausible, even if a large number of variables is included in the imputation model.

Where complete cases and multiple imputation analyses give different results, the analyst should attempt to understand why, and this should be reported in publications.

## What if the data are MCAR?

If we assume data are MAR, then unbiased and statistically more powerful analyses (compared with analyses based on complete cases) can generally be done by including individuals with incomplete data.

There are circumstances in which analyses of **complete cases** will not lead to bias.

- Missing data in predictor variables do not cause bias in analyses of complete cases if the reasons for the missing data are unrelated to the outcome.
  - In such cases, imputing missing data may lessen the loss of precision and power resulting from exclusion of individuals with incomplete predictor variables but are not required in order to avoid bias.

# Stages of Multiple Imputation (1 of 2)

*Multiple imputation ... aims to allow for the uncertainty about the missing data by creating several different plausible imputed data sets and appropriately combining results obtained from each of them.*

The first stage is to create multiple copies of the dataset, with the missing values replaced by imputed values.

- The imputation procedure must fully account for all uncertainty in predicting the missing values by injecting appropriate variability into the multiple imputed values; we can never know the true values of the missing data.

Note that single Imputation of missing values usually causes standard errors to be too small, since it fails to account for the fact that we are uncertain about the missing values.

# Stages of Multiple Imputation (2 of 2)

The second stage is to use standard statistical methods to fit the model of interest to each of the imputed datasets.

- Estimated associations in each of the imputed datasets will differ because of the variation introduced in the imputation of the missing values, and they are only useful when averaged together to give overall estimated associations.
- Standard errors are calculated using Rubin's rules, which take account of the variability in results between the imputed datasets, reflecting the uncertainty associated with the missing values.
- Valid inferences are obtained because we are averaging over the distribution of the missing data given the observed data.

# Comparing Two Linear Models including Multiple Imputation

# Framingham data

```
fram_raw <- read_csv(here("data/framingham.csv")) %>%
    clean_names()

dim(fram_raw)

[1] 4238    17

n_miss(fram_raw)

[1] 645
```

- See https://www.framinghamheartstudy.org/ for more details.

# fram_sub Tibble for Today

```
fram_sub <- fram_raw %>%
    mutate(educ = fct_recode(factor(education),
                           "Some HS" = "1",
                           "HS grad" = "2",
                           "Some Coll" = "3",
                           "Coll grad" = "4")) %>%
    mutate(obese = as.numeric(bmi >= 30)) %>%
    rename(smoker = "current_smoker",
           sbp = "sys_bp") %>%
    mutate(subj_id = as.character(subj_id)) %>%
    select(sbp, educ, smoker, obese, glucose, subj_id)
```

# Data Descriptions (variables we'll use today)

The variables describe n = 4238 adult subjects who were examined at baseline and then followed for ten years to see if they developed incident coronary heart disease during that time.
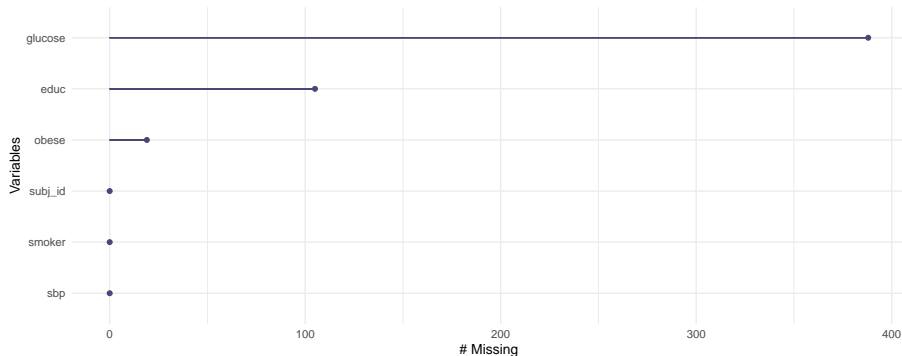
| Variable | Description |
|---------:|-------------|
| educ | four-level factor: educational attainment |
| smoker | 1 = current smoker at time of examination, else 0 |
| sbp | systolic blood pressure (mm Hg) |
| obese | 1 if subject's bmi is 30 or higher, else 0 |
| glucose | blood glucose level in mg/dl |

## Today's Goal

Use linear regression to predict sbp using two different models, in each case accounting for missingness via multiple imputation, where the predictors of interest are glucose, obese, educ, and smoker.

# Which variables are missing data?

gg_miss_var(fram_sub)

# Track missingness with shadow

```
fram_sub_sh <- bind_shadow(fram_sub)

head(fram_sub_sh)

# A tibble: 6 x 12
     sbp educ      smoker obese glucose subj_id sbp_NA educ_NA
   <dbl> <fct>      <dbl> <dbl>   <dbl> <chr>   <fct>  <fct>
1  106  Coll gr~       0     0      77 1       !NA    !NA
2  121  HS grad        0     0      76 2       !NA    !NA
3  128. Some HS        1     0      70 3       !NA    !NA
4  150  Some Co~       1     0     103 4       !NA    !NA
5  130  Some Co~       1     0      85 5       !NA    !NA
6  180  HS grad        0     1      99 6       !NA    !NA
# ... with 4 more variables: smoker_NA <fct>,
#   obese_NA <fct>, glucose_NA <fct>, subj_id_NA <fct>
```

# Our Two Models

Model 2: predict sbp using `glucose` and `obese`.

Model 4: predict sbp using `glucose`, `obese`, `educ`, and `smoker`.

## Model 2 (CC): Two-predictor model for `sbp`

Suppose we ignore the missingness and just run the model on the data with complete information on sbp, glucose and obese.

```
m2_cc <- fram_sub_sh %$% lm(sbp ~ glucose + obese)

tidy(m2_cc, conf.int = TRUE) %>% select(-statistic) %>%
    kable(digits = 3)
```

| term | estimate | std.error | p.value | conf.low | conf.high |
|------|---------:|----------:|--------:|---------:|----------:|
| (Intercept) | 121.671 | 1.244 | 0 | 119.232 | 124.110 |
| glucose | 0.111 | 0.015 | 0 | 0.082 | 0.139 |
| obese | 13.532 | 1.045 | 0 | 11.484 | 15.580 |

# Edited Summary of Model 2 (CC)

```
summary(m2_cc)    ## we'll just look at the bottom
```

```
Residual standard error: 21.42 on 3833 degrees of freedom
  (402 observations deleted due to missingness)
Multiple R-squared:  0.05857,   Adjusted R-squared:  0.05808
F-statistic: 119.2 on 2 and 3833 DF,  p-value: < 2.2e-16
```

```
glance(m2_cc) %>%
    select(nobs, r.squared, adj.r.squared, AIC, BIC) %>%
    kable(digits = c(0, 4, 4, 0, 0))
```

| nobs | r.squared | adj.r.squared | AIC | BIC |
|------|-----------|---------------|------|------|
| 3836 | 0.0586 | 0.0581 | 34401 | 34426 |

## Model 4 (CC): **Four-predictor model for** `sbp`

```
m4_cc <- fram_sub_sh %$%
    lm(sbp ~ glucose + obese + smoker + educ)

tidy(m4_cc, conf.int = TRUE) %>% select(-statistic) %>%
    kable(digits = 3)
```

| term | estimate | std.error | p.value | conf.low | conf.high |
|------|----------|-----------|---------|----------|-----------|
| (Intercept) | 127.107 | 1.388 | 0 | 124.385 | 129.829 |
| glucose | 0.106 | 0.015 | 0 | 0.078 | 0.135 |
| obese | 12.304 | 1.066 | 0 | 10.213 | 14.395 |
| smoker | -4.704 | 0.699 | 0 | -6.075 | -3.332 |
| educHS grad | -3.698 | 0.833 | 0 | -5.332 | -2.065 |
| educSome Coll | -4.724 | 1.010 | 0 | -6.704 | -2.744 |
| educColl grad | -5.954 | 1.158 | 0 | -8.225 | -3.683 |

# Edited Summary of Model 4 (CC)

```
summary(m4_cc)          ## we'll just look at the bottom
```

```
Residual standard error: 21.2 on 3733 degrees of freedom
  (498 observations deleted due to missingness)
Multiple R-squared: 0.08257,   Adjusted R-squared: 0.0811
F-statistic:    56 on 6 and 3733 DF,  p-value: < 2.2e-16
```

```
glance(m4_cc) %>%
    select(nobs, r.squared, adj.r.squared, AIC, BIC) %>%
    kable(digits = c(0, 4, 4, 0, 0))
```

| nobs | r.squared | adj.r.squared | AIC | BIC |
|------|-----------|---------------|-------|-------|
| 3740 | 0.0826 | 0.0811 | 33466 | 33516 |

## Variables used in our models 2 and 4

```
miss_var_summary(fram_sub)
```

```
# A tibble: 6 x 3
  variable n_miss pct_miss
  <chr>     <int>    <dbl>
1 glucose     388    9.16
2 educ        105    2.48
3 obese        19    0.448
4 sbp           0    0
5 smoker        0    0
6 subj_id       0    0
```

- Are we missing data on our outcome for these models?

# Create multiple imputations for this subset

How many subjects have complete / missing data that affect this model?

```
pct_complete_case(fram_sub)
```

```
[1] 88.24917
```

```
pct_miss_case(fram_sub)
```

```
[1] 11.75083
```

Let's create 15 imputed data sets.

```
set.seed(431431)
fram_mice24 <- mice(fram_sub, m = 15, printFlag = FALSE)
```

```
Warning: Number of logged events: 1
```

- Using printFlag = FALSE eliminates a lot of unnecessary (and not particularly informative) output here.

# Summary Information about Imputation Process

```
summary(fram_mice24)

Class: mids
Number of multiple imputations:  15
Imputation methods:
      sbp      educ    smoker     obese   glucose   subj_id
       "" "polyreg"        ""     "pmm"     "pmm"        ""
PredictorMatrix:
        sbp educ smoker obese glucose subj_id
sbp       0    1      1     1       1       0
educ      1    0      1     1       1       0
smoker    1    1      0     1       1       0
obese     1    1      1     0       1       0
glucose   1    1      1     1       0       0
subj_id   1    1      1     1       1       0
Number of logged events:  1
  it im dep     meth      out
```

# Options within `mice` for imputation approaches

Default methods include:

- `pmm` predictive mean matching (default choice for quantitative variables)
- `logreg` logistic regression (default for binary categorical variables)
- `polyreg` polytomous logistic regression (for nominal multi-categorical variables)
- `polr` proportional odds logistic regression (for ordinal categories)

but there are `cart` methods and many others available, too.

# What should we include in an imputation model?

1. If things you are imputing are not Normally distributed, this can pose special challenges, and either a transformation or choosing an imputation method which is robust to these concerns is helpful.
2. Include the outcome when imputing predictors. It causes you to conclude the relationship is weaker than it actually is, if you don't.
3. The MAR assumption may only be reasonable when a certain variable is included in the model.
   - As a result, it's usually a good idea to include as wide a range of variables in imputation models as possible. The concerns we'd have about parsimony in outcome models don't apply here.

# Store one (or more) of the imputed data sets

This will store the fifth imputed data set in `imp_5`.

```
imp_5 <- complete(fram_mice24, 5) %>% tibble()

dim(imp_5)
```

```
[1] 4238    6
```

```
n_miss(imp_5)
```

```
[1] 0
```

# Run Model 2 on each imputed data frame

```
m2_mods <- with(fram_mice24, lm(sbp ~ glucose + obese))
```

```
> summary(m2_mods)
# A tibble: 45 x 6
   term        estimate std.error statistic  p.value  nobs
   <chr>          <dbl>     <dbl>     <dbl>    <dbl> <int>
 1 (Intercept)  122.      1.15      106.     0        4238
 2 glucose        0.104   0.0134      7.74  1.27e-14   4238
 3 obese         13.3     0.984      13.5   5.93e-41   4238
 4 (Intercept)  121.      1.19      102.     0         4238
 5 glucose        0.112   0.0139      8.03  1.23e-15   4238
 6 obese         13.4     0.986      13.6   4.48e-41   4238
 7 (Intercept)  121.      1.15      105.     0         4238
 8 glucose        0.120   0.0135      8.93  6.49e-19   4238
 9 obese         13.1     0.986      13.3   1.89e-39   4238
10 (Intercept)  121.      1.20      101.     0         4238
# ... with 35 more rows
```

- 3 coefficients in each model, times 15 imputations = 45 rows.

## More detailed regression results?

Consider working with the analysis done on the 4th imputed data set (of the 15 created)...

```
m2_a4 <- m2_mods$analyses[[4]]
tidy(m2_a4) %>% kable(digits = 3)
```

| term | estimate | std.error | statistic | p.value |
|------|----------|-----------|-----------|---------|
| (Intercept) | 121.052 | 1.200 | 100.861 | 0 |
| glucose | 0.117 | 0.014 | 8.322 | 0 |
| obese | 13.243 | 0.985 | 13.441 | 0 |

## Pool Results across the 15 imputations

```
m2_pool <- pool(m2_mods)
summary(m2_pool, conf.int = TRUE, conf.level = 0.95)
```

```
          term     estimate   std.error statistic        df
1 (Intercept) 121.1962463 1.28910584 94.015745  408.3743
2     glucose   0.1153255 0.01529575  7.539707  332.3981
3       obese  13.2902350 0.99510223 13.355648 3704.3661
       p.value        2.5 %       97.5 %
1 0.000000e+00 118.66213491 123.7303577
2 4.527489e-13   0.08523679   0.1454141
3 0.000000e+00  11.33923297  15.2412370
```

## Model 2 (Complete Cases vs. Multiple Imputation)

```
tidy(m2_cc, conf.int = TRUE) %>% kable(digits = 3)
```

| term | estimate | std.error | statistic | p.value | conf.low | conf.high |
|------|----------|-----------|-----------|---------|----------|-----------|
| (Intercept) | 121.671 | 1.244 | 97.792 | 0 | 119.232 | 124.110 |
| glucose | 0.111 | 0.015 | 7.577 | 0 | 0.082 | 0.139 |
| obese | 13.532 | 1.045 | 12.954 | 0 | 11.484 | 15.580 |

```
summary(m2_pool, conf.int = TRUE, conf.level = 0.95) %>%
    select(-df) %>% kable(digits = 2)
```

| term | estimate | std.error | statistic | p.value | 2.5 % | 97.5 % |
|------|----------|-----------|-----------|---------|-------|--------|
| (Intercept) | 121.20 | 1.29 | 94.02 | 0 | 118.66 | 123.73 |
| glucose | 0.12 | 0.02 | 7.54 | 0 | 0.09 | 0.15 |
| obese | 13.29 | 1.00 | 13.36 | 0 | 11.34 | 15.24 |

# More Details on Multiple Imputation Modeling

m2_pool

```
> m2_pool
Class: mipo    m = 15
         term  m    estimate        ubar              b          t dfcom       df
1 (Intercept) 15 121.1962463 1.3726312457 2.710900e-01 1.66179387  4235 408.3743
2     glucose 15   0.1153255 0.0001883469 4.276229e-05 0.00023396  4235 332.3981
3       obese 15  13.2902350 0.9703219834 1.866232e-02 0.99022845  4235 3704.3661
        riv      lambda         fmi
1 0.21066301 0.17400631 0.17802209
2 0.24217612 0.19496118 0.19976168
3 0.02051532 0.02010291 0.02063153
```

Definitions of all of these terms are available in the mipo help file.

- riv = relative increase in variance attributable to nonresponse
- fmi = fraction of missing information due to nonresponse

# Model 4 run on each imputed data frame

```
m4_mods <- with(fram_mice24, lm(sbp ~ glucose +
                                obese + smoker + educ))
```

```
> summary(m4_mods)
# A tibble: 105 x 6
   term          estimate std.error statistic  p.value  nobs
   <chr>            <dbl>     <dbl>     <dbl>    <dbl> <int>
 1 (Intercept)    128.       1.27     100.    0          4238
 2 glucose          0.0969   0.0133     7.28  3.93e-13   4238
 3 obese           11.8      0.985     12.0   1.57e-32   4238
 4 smoker          -4.46     0.656     -6.80  1.23e-11   4238
 5 educHS grad     -3.68     0.780     -4.72  2.40e- 6   4238
 6 educSome Coll   -5.31     0.947     -5.61  2.13e- 8   4238
 7 educColl grad   -6.16     1.09      -5.65  1.71e- 8   4238
 8 (Intercept)    127.       1.31      97.2   0          4238
 9 glucose          0.104    0.0138     7.57  4.70e-14   4238
10 obese           11.8      0.986     12.0   1.12e-32   4238
# ... with 95 more rows
```

# Pool Results across the five imputations

```
m4_pool <- pool(m4_mods)

summary(m4_pool, conf.int = TRUE, conf.level = 0.95) %>%
    select(-df) %>% kable(digits = 2)
```

| term | estimate | std.error | statistic | p.value | 2.5 % | 97.5 % |
|---|---:|---:|---:|---:|---:|---:|
| (Intercept) | 126.91 | 1.41 | 90.25 | 0 | 124.14 | 129.67 |
| glucose | 0.11 | 0.02 | 7.16 | 0 | 0.08 | 0.14 |
| obese | 11.77 | 1.00 | 11.82 | 0 | 9.82 | 13.73 |
| smoker | -4.43 | 0.66 | -6.75 | 0 | -5.72 | -3.15 |
| educHS grad | -3.66 | 0.79 | -4.62 | 0 | -5.21 | -2.11 |
| educSome Coll | -5.25 | 0.96 | -5.47 | 0 | -7.13 | -3.37 |
| educColl grad | -6.13 | 1.10 | -5.58 | 0 | -8.28 | -3.97 |

## Complete Cases Result (Model 4)

```
tidy(m4_cc, conf.int = TRUE) %>% select(-statistic) %>%
    kable(digits = 3)
```

| term | estimate | std.error | p.value | conf.low | conf.high |
|---|---|---|---|---|---|
| (Intercept) | 127.107 | 1.388 | 0 | 124.385 | 129.829 |
| glucose | 0.106 | 0.015 | 0 | 0.078 | 0.135 |
| obese | 12.304 | 1.066 | 0 | 10.213 | 14.395 |
| smoker | -4.704 | 0.699 | 0 | -6.075 | -3.332 |
| educHS grad | -3.698 | 0.833 | 0 | -5.332 | -2.065 |
| educSome Coll | -4.724 | 1.010 | 0 | -6.704 | -2.744 |
| educColl grad | -5.954 | 1.158 | 0 | -8.225 | -3.683 |

# Multiple Imputation Result (Model 4)

```
summary(m4_pool, conf.int = TRUE) %>%
  select(-statistic, -df) %>% kable(digits = 3)
```

| term | estimate | std.error | p.value | 2.5 % | 97.5 % |
|---|---:|---:|---:|---:|---:|
| (Intercept) | 126.905 | 1.406 | 0 | 124.143 | 129.668 |
| glucose | 0.107 | 0.015 | 0 | 0.078 | 0.137 |
| obese | 11.774 | 0.996 | 0 | 9.820 | 13.728 |
| smoker | -4.435 | 0.657 | 0 | -5.723 | -3.147 |
| educHS grad | -3.660 | 0.792 | 0 | -5.214 | -2.106 |
| educSome Coll | -5.251 | 0.960 | 0 | -7.133 | -3.369 |
| educColl grad | -6.126 | 1.097 | 0 | -8.276 | -3.975 |

# More Details on Multiple Imputation Modeling

`m4_pool`

```
> m4_pool
Class: mipo    m = 15
          term  m    estimate         ubar            b          t dfcom       df
1  (Intercept) 15 126.9054997 1.6777883057 2.806862e-01 1.9771868946  4231 521.7802
2      glucose 15   0.1074957 0.0001848288 3.778876e-05 0.0002251368  4231 387.9514
3        obese 15  11.7738693 0.9713228490 2.028081e-02 0.9929557100  4231 3628.0260
4       smoker 15  -4.4347370 0.4295295322 1.829583e-03 0.4314810871  4231 4184.1354
5   educHS grad 15  -3.6597613 0.6073110889 1.939744e-02 0.6280016922  4231 3105.0759
6 educSome Coll 15 -5.2509672 0.8952687033 2.432866e-02 0.9212192783  4231 3333.3562
7 educColl grad 15 -6.1257072 1.1823783257 1.974515e-02 1.2034398227  4231 3808.7663
         riv      lambda          fmi
1 0.178448370 0.151426549 0.154660564
2 0.218083002 0.179037883 0.183237700
3 0.022271546 0.021786330 0.022325139
4 0.004543471 0.004522921 0.004998414
5 0.034069200 0.032946732 0.033569016
6 0.028986353 0.028169813 0.028752383
7 0.017812824 0.017501080 0.018016589
>
```

## Estimate $R^2$

```
pool.r.squared(m2_mods)
```

```
           est       lo 95      hi 95 fmi
R^2 0.05944041 0.04607775 0.07426549 NaN
```

```
pool.r.squared(m4_mods)
```

```
           est       lo 95      hi 95 fmi
R^2 0.08146963 0.06617182 0.09805258 NaN
```

- Can also calculated adjusted $R^2$ by using pool.r.squared(m2_mods, adjusted = TRUE). See next slide.

# Estimates of adjusted $R^2$

```
pool.r.squared(m2_mods, adjusted = TRUE)
```

```
              est       lo 95      hi 95 fmi
adj R^2 0.05899617 0.04567921 0.07377905 NaN
```

```
pool.r.squared(m4_mods, adjusted = TRUE)
```

```
           est      lo 95      hi 95 fmi
adj R^2 0.080167 0.06497656 0.09665148 NaN
```

# Tests of Nested Fits after imputation

The models must be nested (same outcome, one set of predictors is a subset of the other) for this to be appropriate.

```
fit4 <- with(fram_mice24,
          expr = lm(sbp ~ glucose + obese + smoker + educ))
fit2 <- with(fram_mice24,
          expr = lm(sbp ~ glucose + obese))
```

# Comparing Model 4 to Model 2 fits

We'll use the Wald test after a linear regression fit.

```
D1(fit4, fit2)

   test statistic df1      df2 dfcom      p.value
 1 ~~ 2  24.85276   4 4093.675  4231 2.322448e-20
        riv
 0.01982057
```

Could also use a likelihood ratio test.

```
D3(fit4, fit2)

   test statistic df1      df2 dfcom p.value          riv
 1 ~~ 2  24.76194   4 276830.7  4231       0 0.01339975
```
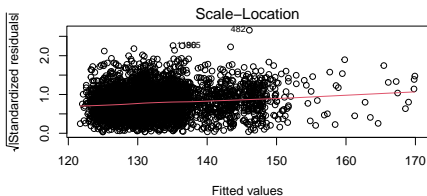
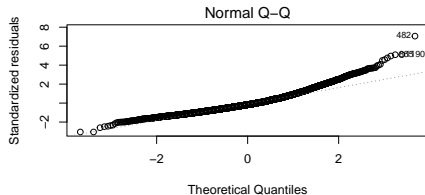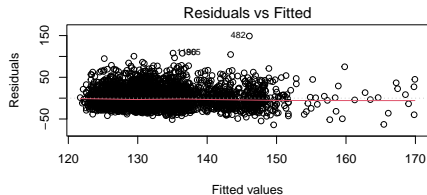# Residual Plots for `mod4` with 6th imputed data set

```
par(mfrow = c(2,2))
plot(m4_mods$analyses[[6]])
par(mfrow = c(1,1))
```

- See the next slide for the results.

# Residual Plots for `mod4` using imputation 6

# Guidelines for Reporting

# Guidelines for reporting, I (Sterne et al.)

How should we report on analyses potentially affected by missing data?

- Report the number of missing values for each variable of interest, or the number of cases with complete data for each important component of the analysis. Give reasons for missing values if possible, and indicate how many individuals were excluded because of missing data when reporting the flow of participants through the study. If possible, describe reasons for missing data in terms of other variables (rather than just reporting a universal reason such as treatment failure.)
- Clarify whether there are important differences between individuals with complete and incomplete data, for example, by providing a table comparing the distributions of key exposure and outcome variables in these different groups
- Describe the type of analysis used to account for missing data (eg, multiple imputation), and the assumptions that were made (eg, missing at random)

# Guidelines for reporting, II (Sterne et al.)

How should we report on analyses that involve multiple imputation?

- Provide details of the imputation modeling (software used, key settings, number of imputed datasets, variables included in imputation procedure, etc.)
- If a large fraction of the data is imputed, compare observed and imputed values.
- Where possible, provide results from analyses restricted to complete cases, for comparison with results based on multiple imputation. If there are important differences between the results, suggest explanations.
- It is also desirable to investigate the robustness of key inferences to possible departures from the missing at random assumption, by assuming a range of missing not at random mechanisms in sensitivity analyses.