

Solutions - Missing data workshop

dr. Iris Eekhout

2022-06-18

Missing value analyses

Solution 1: amount of missing values

```
# load library mice to load boys data
library(mice)
```

```
> Warning: package 'mice' was built under R version 4.1.2
```

a. How many variables have missing data?

All variables, except for age, have missing values, so in total 8 out of 9 variables have missing data.

```
# summary to see which variables have missing data
summary(boys)
```

```
>      age      hgt      wgt      bmi
> Min.   : 0.035  Min.   : 50.00  Min.   :  3.14  Min.   :11.77
> 1st Qu.: 1.581  1st Qu.: 84.88  1st Qu.: 11.70  1st Qu.:15.90
> Median :10.505  Median :147.30  Median : 34.65  Median :17.45
> Mean   : 9.159  Mean   :132.15  Mean   : 37.15  Mean   :18.07
> 3rd Qu.:15.267  3rd Qu.:175.22  3rd Qu.: 59.58  3rd Qu.:19.53
> Max.   :21.177  Max.   :198.00  Max.   :117.40  Max.   :31.74
>      NA's :20      NA's :4      NA's :21
>
>      hc      gen      phb      tv      reg
> Min.   :33.70  G1   : 56  P1   : 63  Min.   : 1.00  north: 81
> 1st Qu.:48.12  G2   : 50  P2   : 40  1st Qu.: 4.00  east :161
> Median :53.00  G3   : 22  P3   : 19  Median :12.00  west :239
> Mean   :51.51  G4   : 42  P4   : 32  Mean   :11.89  south:191
> 3rd Qu.:56.00  G5   : 75  P5   : 50  3rd Qu.:20.00  city : 73
> Max.   :65.00  NA's :503  P6   : 41  Max.   :25.00  NA's :  3
> NA's    :46      NA's :503  NA's   :522
```

b. How many rows in the data contain missing values?

In total 525 rows in the data have missing values, this is ~70%.

```
nic(boys)
```

```
> [1] 525
```

```
nic(boys)/nrow(boys)
```

```
> [1] 0.7018717
```

c. How many overall matrix entries are missing? And how many observed?

1622 matrix entries are missing and 5110 are observed; ~25% of the matrix entries are missing.

```
sum(is.na(boys))
```

```
> [1] 1622
```

```
sum(!is.na(boys))
```

```
> [1] 5110
```

```
1622/(1622+5110)
```

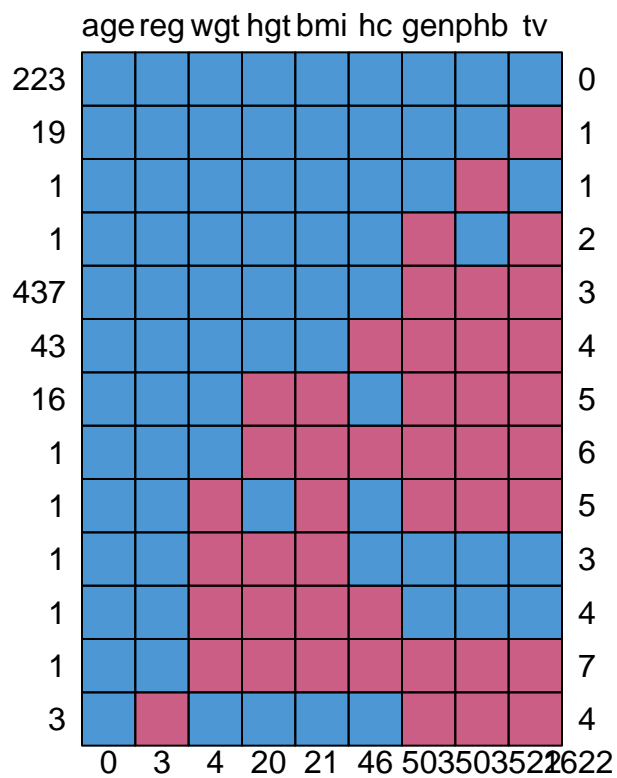
```
> [1] 0.2409388
```

Solution 2: missing data patterns

a. How many different missing data patterns occur in the data?

14 patterns

```
nrow(mice::md.pattern(boys, plot= T))
```

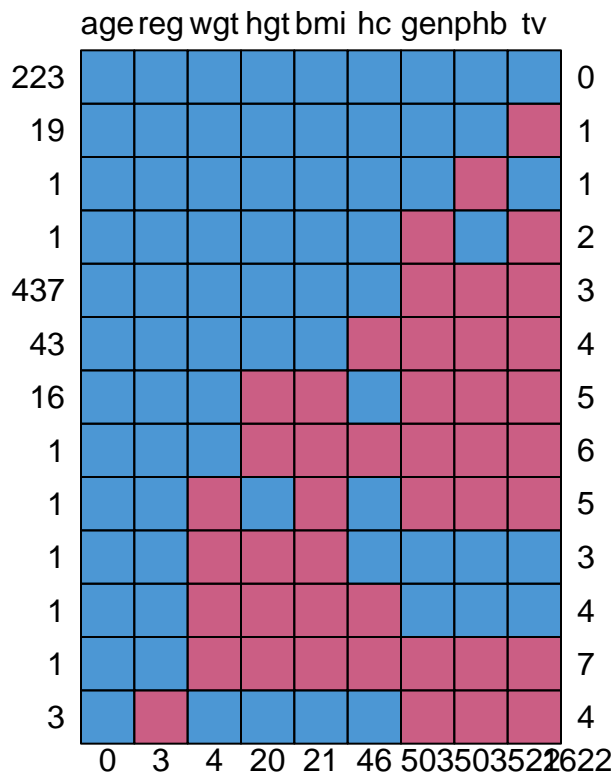


```
> [1] 14
```

b. What is the most frequently occurring pattern in the data?

The pattern with “gen”, “phb” and “tv” missing, occurs 437 times.

```
mice::md.pattern(boys, plot= T)
```



```

>      age reg wgt hgt bmi hc gen phb tv
> 223   1   1   1   1   1  1   1   1   1   0
> 19    1   1   1   1   1  1   1   1   0   1
> 1     1   1   1   1   1  1   1   0   1   1
> 1     1   1   1   1   1  1   1   0   1   0   2
> 437   1   1   1   1   1  1   1   0   0   0   3
> 43    1   1   1   1   1  1  0   0   0   0   4
> 16    1   1   1   0   0  0  1   0   0   0   5
> 1     1   1   1   0   0  0  0   0   0   0   6
> 1     1   1   0   1   0  1  0   0   0   0   5
> 1     1   1   0   0   0  0  1   1   1   1   3
> 1     1   1   0   0   0  0  0   1   1   1   4
> 1     1   1   0   0   0  0  0   0   0   0   7
> 3     1   0   1   1   1  1  0   0   0   0   4
>      0   3   4   20  21 46 503 503 522 1622

```

c. Looking at patterns that occur more than incidental (once or twice), which variables happen to be missing together often?

Variables that are most often missing at the same time are “gen”, “phb”, and “tv”. The patterns that occur more than once involve mostly all of these variables (pattern with “hc” and “gen”, “phb”, “tv” 43 times and the pattern with “hgt”, “bmi” and “gen”, “phb”, “tv” 16 times, pattern “reg” and “gen”, “phb”, “tv”, 3 times, pattern with “tv” missing 19 times).

d. Inspect the missing data pairs. With what other variable(s) is height observed together with in more than half of the cases?

The answer can be found looking at the first matrix `rr`. In the row of “hgt”, the column values that are higher than 374 indicate variables that are observed with hgt more than half of the time: “age”, “wgt”, “bmi”, “hc”, and “reg”.

```
mice::md.pairs(boys)
```

```
> $rr
>      age hgt wgt bmi  hc gen phb  tv reg
> age 748 728 744 727 702 245 245 226 745
> hgt 728 728 727 727 685 243 243 224 725
> wgt 744 727 744 727 700 243 243 224 741
> bmi 727 727 727 727 684 243 243 224 724
> hc  702 685 700 684 702 244 244 225 699
> gen 245 243 243 243 244 245 244 226 245
> phb 245 243 243 243 244 244 245 225 245
> tv  226 224 224 224 225 226 225 226 226
> reg 745 725 741 724 699 245 245 226 745
>
> $rm
>      age hgt wgt bmi hc gen phb  tv reg
> age   0  20  4  21 46 503 503 522  3
> hgt   0   0  1   1 43 485 485 504  3
> wgt   0  17  0  17 44 501 501 520  3
> bmi   0   0  0   0 43 484 484 503  3
> hc    0  17  2  18  0 458 458 477  3
> gen   0   2  2   2  1  0  1  19  0
> phb   0   2  2   2  1  1  0  20  0
> tv    0   2  2   2  1  0  1   0  0
> reg   0  20  4  21 46 500 500 519  0
>
> $mr
>      age hgt wgt bmi  hc gen phb tv reg
> age   0   0  0   0   0  0  0  0  0
> hgt  20   0 17   0  17  2  2  2 20
> wgt   4   1  0   0   2  2  2  2  4
> bmi  21   1 17   0  18  2  2  2 21
> hc   46  43 44  43   0  1  1  1 46
> gen 503 485 501 484 458  0  1  0 500
> phb 503 485 501 484 458  1  0  1 500
> tv  522 504 520 503 477 19 20  0 519
> reg   3   3  3   3   3  0  0  0  0
>
> $mm
>      age hgt wgt bmi hc gen phb  tv reg
> age   0   0  0   0  0  0  0  0  0
> hgt   0  20  3  20  3 18 18 18  0
> wgt   0   3  4   4  2  2  2  2  0
> bmi   0  20  4  21  3 19 19 19  0
> hc    0   3  2   3 46 45 45 45  0
> gen   0  18  2  19 45 503 502 503  3
> phb   0  18  2  19 45 502 503 502  3
> tv    0  18  2  19 45 503 502 522  3
> reg   0   0  0   0  0  3  3  3  3
```

Solution 3: understanding missing data mechanisms

a. What is the mean and standard deviation of knee pain score? And the association between BMI and knee pain (coefficient, standard error and p-value)?

Mean= 14.81, sd=3.21; The association is significant with coefficient=0.35; se=0.14.

b. What are the mean and standard deviation of the knee pain score? What is association between BMI and knee pain?

Mean= 14.70, sd=3.24; The association is not significant (coefficient=0.33; se=0.18).

c. How do these results compare to the complete data results?

The mean and standard deviation have not changed much, however the power for the association was decreased which caused the association between BMI and knee pain to not be significant anymore.

d. What happens to the association between BMI and knee pain? Explain differences with the previous answer (sample size 100).

At 0% missing: coefficient=0.57 and se=0.08 (significant); at 30% missing: coefficient=0.53 and se=0.10 (significant). The association does not change (much) and remains significant. The difference with answer c is explained by the change in sample size. Larger sample size, makes the analysis more robust against (MCAR) missing data.

e. What is the association between BMI and knee pain? How does this compare to the association when the data were MCAR?

There is a significant association with a coefficient of 0.49, se=0.11; the association is now less strong; the coefficient is lower (biased) (was 0.57 for 0% missing) and the standard error has increased slightly (was 0.08 for 0% missing).

f. When there are 30% MAR missing data at sample size 250, at what BMI values do missing data on knee pain occur (inspect the scatterplot and the boxplots).

More missing values at higher BMI values.

g. Comparing the histograms, what knee pain values are mostly missing?

In the MNAR situation, more missings are present in the higher values of knee pain.

h. What happens with the association between BMI and knee pain?

- MCAR: coefficient= 0.46; se=0.10
- MAR: coefficient=0.37; se=0.11
- MNAR: coefficient=0.27; se=0.11
- 0% missing: coefficient=0.54; se=0.07.

The association becomes less strong when you change from MCAR to MAR to MNAR, so coefficients get more biased. Also for MCAR the missings are nicely distributed over the BMI values. In the MAR mechanism there are more missings at higher values of BMI but also lower Knee Pain scores are missing. In the MNAR mechanism, mostly higher values of Knee Pain scores are missing.

i What happens with the mean and standard deviation of the knee pain score?

- MCAR: mean=14.79; sd=2.98
- MAR: mean=14.01; sd=2.85
- MNAR: mean= 13.34; sd=2.69
- 0%missing: mean=14.77; sd=3.11

Both mean and standard deviation decrease when changing from MCAR to MAR to MNAR.

Solution 4: evaluating the missing data mechanism

a. Evaluate the missing data mechanism for the airquality data with univariate tests. What are your conclusions?

Evaluation using T-tests for continuous variables and the Chi-square for categorical variables.

First create the missing data indicators for each variable with missing data.

```
library(dplyr)
summary(airquality)
```

```
>      Ozone      Solar.R      Wind      Temp
> Min.   : 1.00   Min.   : 7.0   Min.   : 1.700   Min.   :56.00
> 1st Qu.: 18.00   1st Qu.:115.8   1st Qu.: 7.400   1st Qu.:72.00
> Median : 31.50   Median :205.0   Median : 9.700   Median :79.00
> Mean   : 42.13   Mean   :185.9   Mean   : 9.958   Mean   :77.88
> 3rd Qu.: 63.25   3rd Qu.:258.8   3rd Qu.:11.500   3rd Qu.:85.00
> Max.   :168.00   Max.   :334.0   Max.   :20.700   Max.   :97.00
> NA's   :37      NA's    :7
>      Month      Day
> Min.   :5.000   Min.   : 1.0
> 1st Qu.:6.000   1st Qu.: 8.0
> Median :7.000   Median :16.0
> Mean   :6.993   Mean   :15.8
> 3rd Qu.:8.000   3rd Qu.:23.0
> Max.   :9.000   Max.   :31.0
>
```

```
airqualitym <- airquality %>%
  #create missing data indicators
  mutate(ROzone = is.na(Ozone),
         RSolar.R = is.na(Solar.R))
```

Do a T-test for each missing data indicator with the continuous variables and a chi square test for the categorical variables. We investigate Month as categorical, but it can also be used as continuous.

- Missing data in Ozone are related to Month, the probability for missing Ozone data is higher earlier in the year.
- The Chi-square test Solar.R and Month seems significant, but also throws a warning. So we cannot really be sure about the results. Probably also because there are only few missings in Solar.R. Tests for the other variables do not show a relation with the probability of missing data in Solar.R

Based on the univariate analyses we may conclude that the missing data are not Missing Completely at Random (not-MCAR). There are other measured variable related to the probability of missing data.

```
# Univariate tests for Ozone
t.test(Solar.R ~ ROzone, data = airqualitym)
```

```
>
> Welch Two Sample t-test
>
```

```

> data: Solar.R by ROzone
> t = -0.27457, df = 58.995, p-value = 0.7846
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -39.05621 29.63124
> sample estimates:
> mean in group FALSE mean in group TRUE
> 184.8018 189.5143

```

```
t.test(Wind ~ ROzone, data = airqualitym)
```

```

>
> Welch Two Sample t-test
>
> data: Wind by ROzone
> t = -0.60911, df = 63.646, p-value = 0.5446
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -1.6893132 0.8999377
> sample estimates:
> mean in group FALSE mean in group TRUE
> 9.862069 10.256757

```

```
t.test(Temp ~ ROzone, data = airqualitym)
```

```

>
> Welch Two Sample t-test
>
> data: Temp by ROzone
> t = -0.026831, df = 60.447, p-value = 0.9787
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -3.643306 3.546847
> sample estimates:
> mean in group FALSE mean in group TRUE
> 77.87069 77.91892

```

```
t.test(Month ~ ROzone, data = airqualitym)
```

```

>
> Welch Two Sample t-test
>
> data: Month by ROzone
> t = 4.0092, df = 92.075, p-value = 0.0001236
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> 0.4273815 1.2664675
> sample estimates:
> mean in group FALSE mean in group TRUE
> 7.198276 6.351351

```



```
t.test(Day ~ ROzone, data = airqualitym)
```

```
>
> Welch Two Sample t-test
>
> data: Day by ROzone
> t = -0.64426, df = 57.826, p-value = 0.522
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -4.576080 2.347749
> sample estimates:
> mean in group FALSE mean in group TRUE
> 15.53448 16.64865
```

```
chisq.test(airqualitym$ROzone, airqualitym$Month)
```

```
>
> Pearson's Chi-squared test
>
> data: airqualitym$ROzone and airqualitym$Month
> X-squared = 44.751, df = 4, p-value = 4.48e-09
```

```
# Univariate tests for Solar.R
```

```
t.test(Ozone ~ RSolar.R, data = airqualitym)
```

```
>
> Welch Two Sample t-test
>
> data: Ozone by RSolar.R
> t = -0.052696, df = 4.4917, p-value = 0.9602
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -36.0892 34.6874
> sample estimates:
> mean in group FALSE mean in group TRUE
> 42.0991 42.8000
```

```
t.test(Wind ~ RSolar.R, data = airqualitym)
```

```
>
> Welch Two Sample t-test
>
> data: Wind by RSolar.R
> t = 0.65629, df = 6.4571, p-value = 0.5343
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -2.674488 4.681338
> sample estimates:
> mean in group FALSE mean in group TRUE
> 10.00342 9.00000
```

```
t.test(Temp ~ RSolar.R, data = airqualitym)
```

```
>
> Welch Two Sample t-test
>
> data: Temp by RSolar.R
> t = 0.98706, df = 6.2689, p-value = 0.3602
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -7.436381 17.669258
> sample estimates:
> mean in group FALSE mean in group TRUE
> 78.11644 73.00000
```

```
t.test(Month ~ RSolar.R, data = airqualitym)
```

```
>
> Welch Two Sample t-test
>
> data: Month by RSolar.R
> t = 1.2018, df = 6.4489, p-value = 0.2717
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -0.7432664 2.2266323
> sample estimates:
> mean in group FALSE mean in group TRUE
> 7.027397 6.285714
```

```
t.test(Day ~ RSolar.R, data = airqualitym)
```

```
>
> Welch Two Sample t-test
>
> data: Day by RSolar.R
> t = 2.1941, df = 6.6803, p-value = 0.06612
> alternative hypothesis: true difference in means between group FALSE and group TRUE is not equal to 0
> 95 percent confidence interval:
> -0.6161166 14.5769776
> sample estimates:
> mean in group FALSE mean in group TRUE
> 16.123288 9.142857
```

```
chisq.test(airqualitym$RSolar.R, airqualitym$Month)
```

```
> Warning in chisq.test(airqualitym$RSolar.R, airqualitym$Month): Chi-squared
> approximation may be incorrect
```

```
>
> Pearson's Chi-squared test
>
> data: airqualitym$RSolar.R and airqualitym$Month
> X-squared = 11.136, df = 4, p-value = 0.02507
```

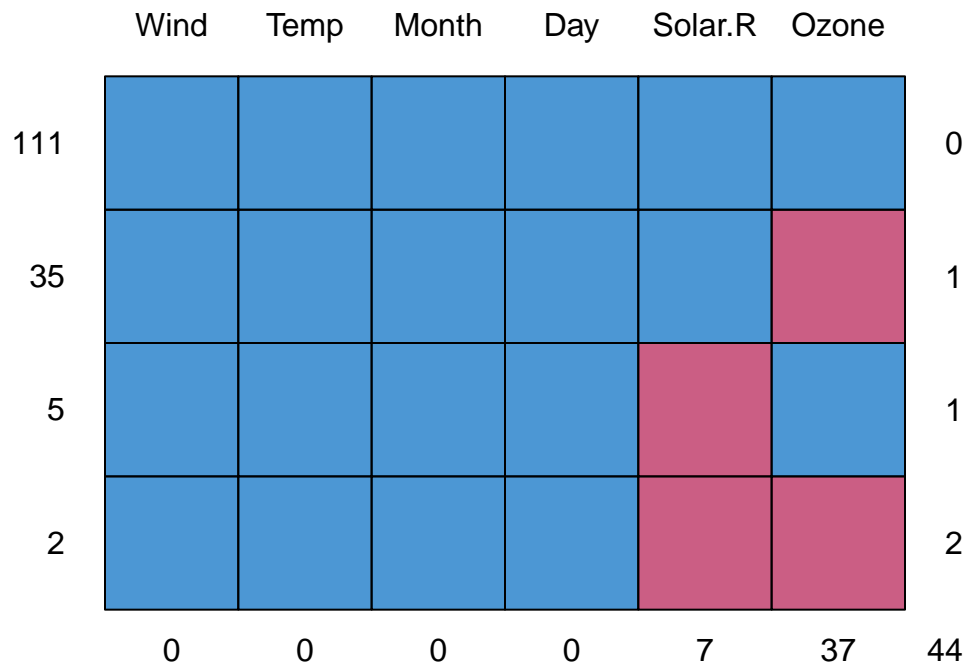
b. Evaluate the missing data mechanism for the airquality data with a multivariate test. What are your conclusions?

First, investigate which variables have missing data in the `airquality` data: Ozone and Solar.R

```
summary(airquality)
```

```
>      Ozone      Solar.R      Wind      Temp
> Min.   : 1.00   Min.   : 7.0   Min.   : 1.700   Min.   :56.00
> 1st Qu.: 18.00   1st Qu.:115.8   1st Qu.: 7.400   1st Qu.:72.00
> Median : 31.50   Median :205.0   Median : 9.700   Median :79.00
> Mean   : 42.13   Mean   :185.9   Mean   : 9.958   Mean   :77.88
> 3rd Qu.: 63.25   3rd Qu.:258.8   3rd Qu.:11.500   3rd Qu.:85.00
> Max.   :168.00   Max.   :334.0   Max.   :20.700   Max.   :97.00
> NA's   :37      NA's    :7
>      Month      Day
> Min.   :5.000   Min.   : 1.0
> 1st Qu.:6.000   1st Qu.: 8.0
> Median :7.000   Median :16.0
> Mean   :6.993   Mean   :15.8
> 3rd Qu.:8.000   3rd Qu.:23.0
> Max.   :9.000   Max.   :31.0
>
```

```
mice::md.pattern(airquality)
```



```

>      Wind Temp Month Day Solar.R Ozone
> 111    1    1    1    1        1    1  0
> 35     1    1    1    1        1    0  1
> 5      1    1    1    1        0    1  1
> 2      1    1    1    1        0    0  2
>      0    0    0    0        7   37 44

```

Create missing data indicators for these two variables. Additional, we can also make one missing indicator for any missing values. Note that this is only useful if we have at least some variables that have no missing data.

```

airqualitym <- airquality %>%
  mutate(ROzone = is.na(Ozone),
         RSolar.R = is.na(Solar.R),
         Rind = is.na(Ozone) | is.na(Solar.R))

```

Do a logistic regression analysis for each of the missing data indicators. Both Temp and Month seem to be related to the missing values in Ozone. There are no measured variables related to the missing values in Solar.R. Based on these results we can conclude that the missing values in the airquality dataset are not-MCAR.

```

glm(ROzone ~ Solar.R + Wind + Temp + Month + Day, data = airqualitym, family = "binomial") %>% summary

```

```

>
> Call:
> glm(formula = ROzone ~ Solar.R + Wind + Temp + Month + Day, family = "binomial",
>      data = airqualitym)
>
> Deviance Residuals:
>      Min       1Q   Median       3Q      Max
> -1.4547  -0.8277  -0.5250  -0.2253   2.2683
>
> Coefficients:
>              Estimate Std. Error z value Pr(>|z|)
> (Intercept) -3.028609    2.316740  -1.307  0.191120
> Solar.R      -0.002155    0.002496  -0.864  0.387769
> Wind         0.057605    0.063930   0.901  0.367552
> Temp         0.081839    0.031363   2.609  0.009069 **
> Month        -0.726536    0.218087  -3.331  0.000864 ***
> Day          0.012030    0.023231   0.518  0.604555
> ---
> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>
> (Dispersion parameter for binomial family taken to be 1)
>
>      Null deviance: 160.82  on 145  degrees of freedom
> Residual deviance: 144.39  on 140  degrees of freedom
>   (7 observations deleted due to missingness)
> AIC: 156.39
>
> Number of Fisher Scoring iterations: 5

```

```
glm(RSolar.R ~ Ozone + Wind + Temp + Month + Day, data = airqualitym, family = "binomial") %>% summary
```

```
>
> Call:
> glm(formula = RSolar.R ~ Ozone + Wind + Temp + Month + Day, family = "binomial",
>   data = airqualitym)
>
> Deviance Residuals:
>      Min       1Q   Median       3Q      Max
> -0.84141  -0.27298  -0.15226  -0.06878   2.54340
>
> Coefficients:
>              Estimate Std. Error z value Pr(>|z|)
> (Intercept)  0.08394     6.41576   0.013  0.9896
> Ozone        -0.02426     0.02503  -0.969  0.3324
> Wind         -0.22019     0.21582  -1.020  0.3076
> Temp          0.06002     0.10130   0.592  0.5536
> Month        -0.46001     0.50740  -0.907  0.3646
> Day          -0.16317     0.08894  -1.835  0.0666 .
> ---
> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>
> (Dispersion parameter for binomial family taken to be 1)
>
> Null deviance: 41.223  on 115  degrees of freedom
> Residual deviance: 32.229  on 110  degrees of freedom
> (37 observations deleted due to missingness)
> AIC: 44.229
>
> Number of Fisher Scoring iterations: 7
```

```
# analysis for the indicator of overall missing data
```

```
glm(Rind ~ Wind + Temp + Month + Day, data = airqualitym, family = "binomial") %>% summary
```

```
>
> Call:
> glm(formula = Rind ~ Wind + Temp + Month + Day, family = "binomial",
>   data = airqualitym)
>
> Deviance Residuals:
>      Min       1Q   Median       3Q      Max
> -1.3113  -0.8743  -0.5758   1.2091   2.2435
>
> Coefficients:
>              Estimate Std. Error z value Pr(>|z|)
> (Intercept) -0.665206     2.083304  -0.319  0.749496
> Wind          0.012608     0.059386   0.212  0.831873
> Temp          0.050067     0.025975   1.927  0.053922 .
> Month        -0.631889     0.186734  -3.384  0.000715 ***
> Day          -0.003769     0.021176  -0.178  0.858728
> ---
> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

>
> (Dispersion parameter for binomial family taken to be 1)
>
> Null deviance: 179.83 on 152 degrees of freedom
> Residual deviance: 165.15 on 148 degrees of freedom
> AIC: 175.15
>
> Number of Fisher Scoring iterations: 4

```

Multiple imputation

Solution 5: multiple imputation in mice

a. How many imputed datasets are generated?

The output returned from the mice functions states: “Number of multiple imputations: 5”.

```
imp <- mice(nhanes2)
```

```

>
> iter imp variable
> 1 1 bmi hyp chl
> 1 2 bmi hyp chl
> 1 3 bmi hyp chl
> 1 4 bmi hyp chl
> 1 5 bmi hyp chl
> 2 1 bmi hyp chl
> 2 2 bmi hyp chl
> 2 3 bmi hyp chl
> 2 4 bmi hyp chl
> 2 5 bmi hyp chl
> 3 1 bmi hyp chl
> 3 2 bmi hyp chl
> 3 3 bmi hyp chl
> 3 4 bmi hyp chl
> 3 5 bmi hyp chl
> 4 1 bmi hyp chl
> 4 2 bmi hyp chl
> 4 3 bmi hyp chl
> 4 4 bmi hyp chl
> 4 5 bmi hyp chl
> 5 1 bmi hyp chl
> 5 2 bmi hyp chl
> 5 3 bmi hyp chl
> 5 4 bmi hyp chl
> 5 5 bmi hyp chl

```

```
imp
```

```

> Class: mids
> Number of multiple imputations: 5
> Imputation methods:

```

```

>      age      bmi      hyp      chl
>      ""      "pmm" "logreg" "pmm"
> PredictorMatrix:
>      age bmi hyp chl
> age    0    1    1    1
> bmi    1    0    1    1
> hyp    1    1    0    1
> chl    1    1    1    0

```

b. How many sets of results are generated

The `with` function can be used to automatically analyze all imputed datasets. The analysis results are stored as `fit`. There are 5 sets of results generated.

```

fit <- with(imp, lm(bmi ~ age + hyp + chl))
fit

```

```

> call :
> with.mids(data = imp, expr = lm(bmi ~ age + hyp + chl))
>
> call1 :
> mice(data = nhanes2)
>
> nmis :
> age bmi hyp chl
>  0   9   8  10
>
> analyses :
> [[1]]
>
> Call:
> lm(formula = bmi ~ age + hyp + chl)
>
> Coefficients:
> (Intercept)    age40-59    age60-99    hypyes      chl
>    21.1212    -4.5387    -6.6866     5.4800     0.0345
>
>
> [[2]]
>
> Call:
> lm(formula = bmi ~ age + hyp + chl)
>
> Coefficients:
> (Intercept)    age40-59    age60-99    hypyes      chl
>    22.89223    -2.21384    -1.81372     0.12090     0.02124
>
>
> [[3]]
>
> Call:
> lm(formula = bmi ~ age + hyp + chl)
>
> Coefficients:

```

```

> (Intercept)      age40-59      age60-99      hypyes      chl
>    17.85584      -4.24203      -5.81019      2.22769      0.05607
>
>
> [[4]]
>
> Call:
> lm(formula = bmi ~ age + hyp + chl)
>
> Coefficients:
> (Intercept)      age40-59      age60-99      hypyes      chl
>    18.22754      -3.71703      -8.74612      2.26970      0.05725
>
>
> [[5]]
>
> Call:
> lm(formula = bmi ~ age + hyp + chl)
>
> Coefficients:
> (Intercept)      age40-59      age60-99      hypyes      chl
>    18.25935      -3.80209      -5.80226      1.92283      0.05036

```

c. What are the most relevant predictors for bmi?

The most relevant predictors are age and cholesterol. Older respondents have a lower bmi, whereas cholesterol is positively related (higher cholesterol means a higher bmi).

```

combi <- pool(fit)
summary(combi)

```

```

>      term      estimate std.error statistic      df      p.value
> 1 (Intercept) 19.67123335 4.41489251  4.4556540  9.785551 0.001291374
> 2   age40-59  -3.70274054 2.23244647 -1.6586022 12.939466 0.121225723
> 3   age60-99 -5.77177779 3.58390221 -1.6104730  4.506041 0.174559336
> 4    hypyes   2.40421707 2.97801196  0.8073228  5.720634 0.451731863
> 5      chl    0.04388315 0.02608737  1.6821608  7.075610 0.135960320

```

Solution 6: multiple imputation model and convergence

a. Adjust the predictor matrix so that the variables that have more than 50% missing values are excluded as predictors for the imputation.

First create the predictor matrix for the boys dataset with the `make.predictorMatrix()` function in `mice`.

```

pred <- make.predictorMatrix(boys)

```

Inspect the boys dataset, and find out what variables have more than 50% missing values. These are “gen”, “phb”, and “tv”.

```

colMeans(is.na(boys))

```



```

>          age      hgt      wgt      bmi      hc      gen
> 0.000000000 0.026737968 0.005347594 0.028074866 0.061497326 0.672459893
>          phb      tv      reg
> 0.672459893 0.697860963 0.004010695

```

Now, exclude variables “gen”, “phb”, and “tv” as predictors from this matrix. Note that predictors are in the columns.

```

pred[,c("gen", "phb", "tv")] <- 0
pred

```

```

>      age hgt wgt bmi hc gen phb tv reg
> age   0   1   1   1   1   0   0   0   1
> hgt   1   0   1   1   1   0   0   0   1
> wgt   1   1   0   1   1   0   0   0   1
> bmi   1   1   1   0   1   0   0   0   1
> hc    1   1   1   1   0   0   0   0   1
> gen   1   1   1   1   1   0   0   0   1
> phb   1   1   1   1   1   0   0   0   1
> tv    1   1   1   1   1   0   0   0   1
> reg   1   1   1   1   1   0   0   0   0

```

Perform multiple imputation on the boys data with the predictor matrix designed at assignment 6a with 10 imputations and 10 iterations.

```

imp <- mice(boys, m = 10, maxit = 10, predictorMatrix = pred)

```

```

>
> iter imp variable
>  1  1  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  2  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  3  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  4  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  5  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  6  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  7  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  8  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1  9  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  1 10  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  1  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  2  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  3  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  4  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  5  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  6  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  7  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  8  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2  9  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  2 10  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  3  1  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  3  2  hgt  wgt  bmi  hc  gen  phb  tv  reg
>  3  3  hgt  wgt  bmi  hc  gen  phb  tv  reg

```

[illegible]

```

> 8 8 hgt wgt bmi hc gen phb tv reg
> 8 9 hgt wgt bmi hc gen phb tv reg
> 8 10 hgt wgt bmi hc gen phb tv reg
> 9 1 hgt wgt bmi hc gen phb tv reg
> 9 2 hgt wgt bmi hc gen phb tv reg
> 9 3 hgt wgt bmi hc gen phb tv reg
> 9 4 hgt wgt bmi hc gen phb tv reg
> 9 5 hgt wgt bmi hc gen phb tv reg
> 9 6 hgt wgt bmi hc gen phb tv reg
> 9 7 hgt wgt bmi hc gen phb tv reg
> 9 8 hgt wgt bmi hc gen phb tv reg
> 9 9 hgt wgt bmi hc gen phb tv reg
> 9 10 hgt wgt bmi hc gen phb tv reg
> 10 1 hgt wgt bmi hc gen phb tv reg
> 10 2 hgt wgt bmi hc gen phb tv reg
> 10 3 hgt wgt bmi hc gen phb tv reg
> 10 4 hgt wgt bmi hc gen phb tv reg
> 10 5 hgt wgt bmi hc gen phb tv reg
> 10 6 hgt wgt bmi hc gen phb tv reg
> 10 7 hgt wgt bmi hc gen phb tv reg
> 10 8 hgt wgt bmi hc gen phb tv reg
> 10 9 hgt wgt bmi hc gen phb tv reg
> 10 10 hgt wgt bmi hc gen phb tv reg

```

b. What methods are used for the imputation of each variable and explain why these are used.

- For “hgt”, “wgt”, “bmi”, “hc” and “tv” the “pmm” (predictive mean matching) method is used. This is the default for continuous variables and the variables indicated are all continuous.
- For “gen” and “phb” the “polr” method is used, which is the default for ordinal variables. The imputation function is a proportional odds model.
- For “reg” the “polyreg” method is used. This method is the default for unordered nominal variables and is the polytomous logistic regression.

```
imp$method
```

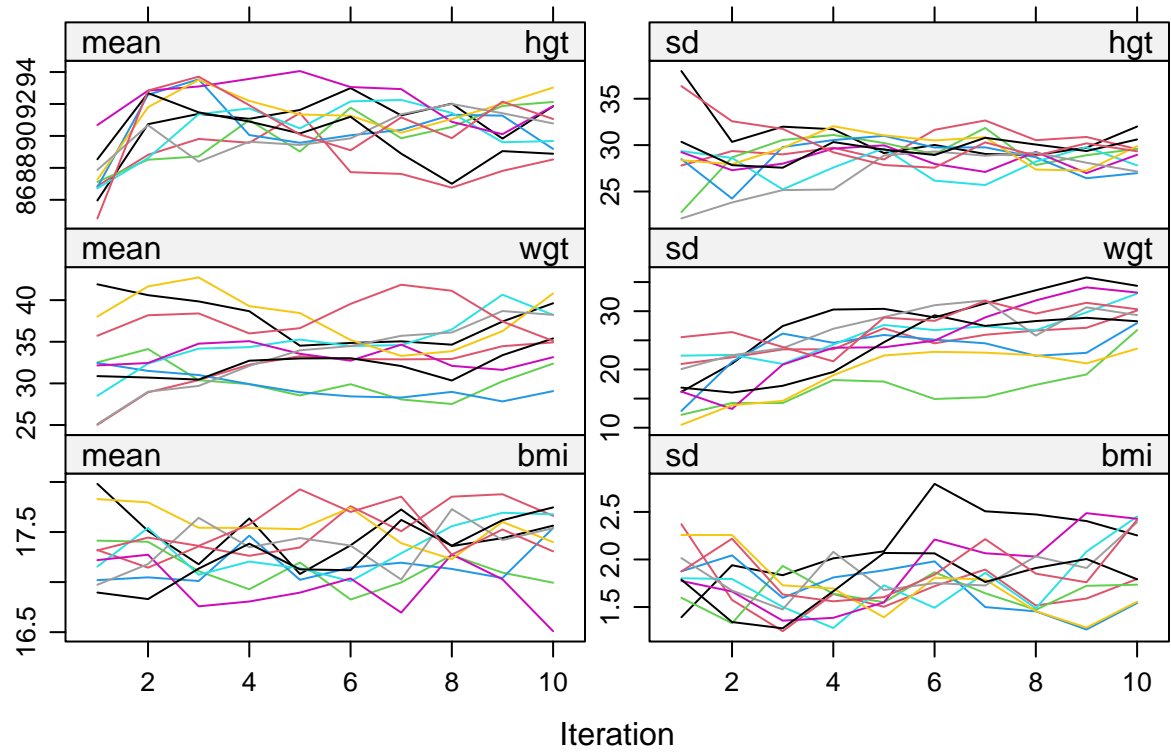
```

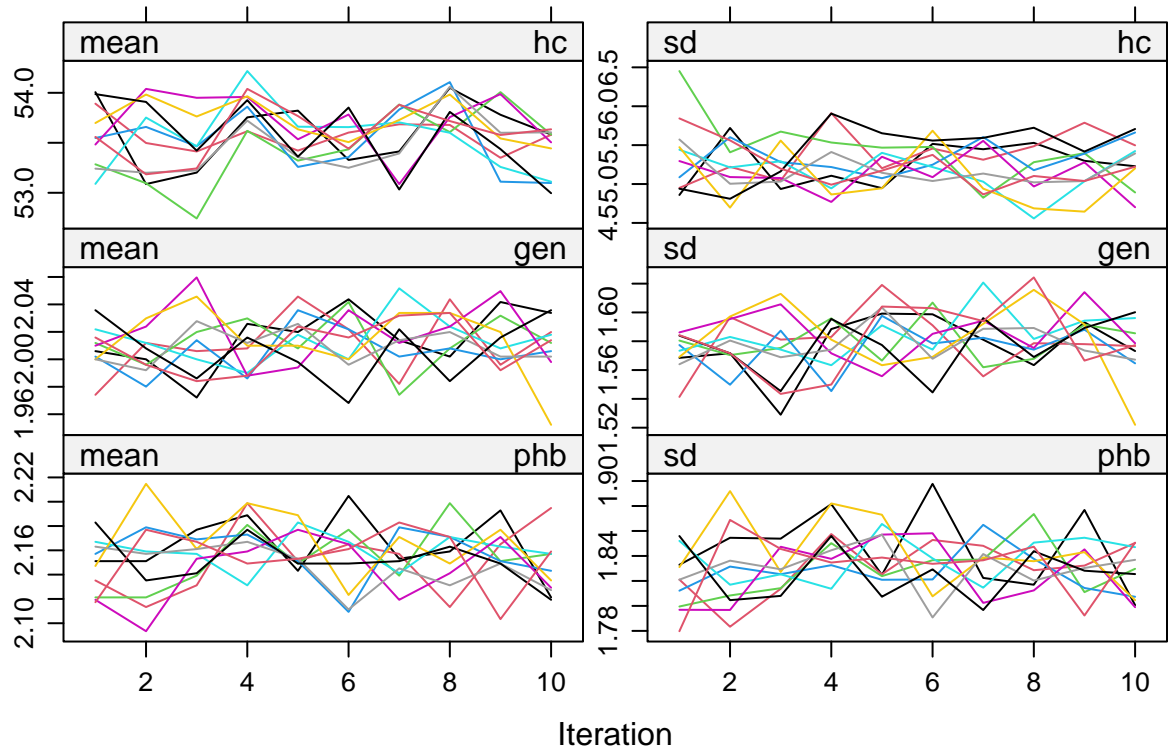
>      age      hgt      wgt      bmi      hc      gen      phb      tv
>      ""      "pmm"    "pmm"    "pmm"    "pmm"    "polr"    "polr"    "pmm"
>      reg
> "polyreg"

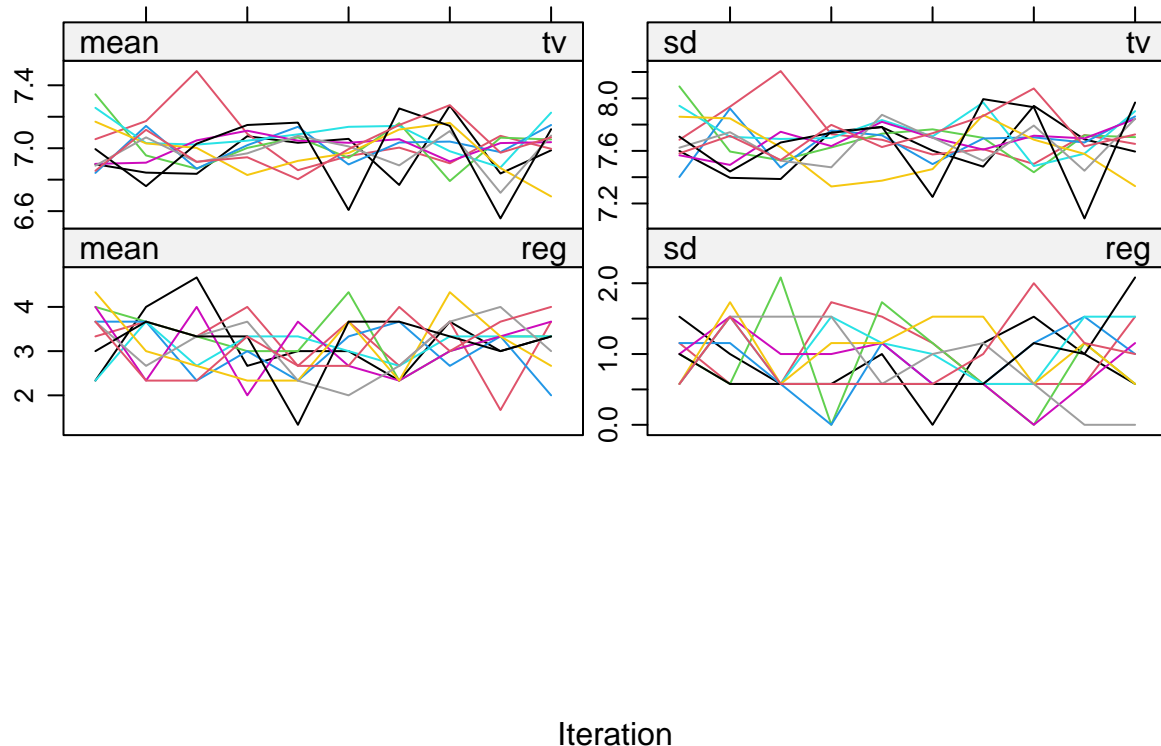
```

c. Inspect the iteration plots. What are your observations?

```
plot(imp)
```







The iteration plots for the variables “hc”, “gen”, “phb”, “tv”, and “reg” all show crossing interacting lines that are more or less centered around. For the variables “hgt”, “wgt” and “bmi” the lines do not cross as much and some lines are below the mean for all iterations, whereas other lines are above the mean for all iterations. So the iteration plots for “hgt”, “wgt” and “bmi” do not show a very good convergence.

d. Adjust the predictor matrix such that hgt and wgt are not imputed by bmi, or vice versa, and that hgt and wgt are not used as predictors together with bmi. Why do you think that these changes are needed?

Option 1:

First, we remove “hgt” and “wgt” as predictors for the imputation of “bmi”. Second, we remove “bmi” as a predictor for all variables, to ensure that “bmi” is not used as predictor in a model together with “wgt” and “hgt”.

```
pred1 <- pred
pred1["bmi", c("hgt", "wgt")] <- 0
pred1[, "bmi"] <- 0
```

Option 2:

We remove “hgt” and “wgt” as predictors for all variables, that way they won’t be together in the model with “bmi” and they will not be used as predictors for imputing “bmi”.

```
pred2 <- pred
pred2[, c("hgt", "wgt")] <- 0
pred2[c("hgt", "wgt"), "bmi"] <- 0
```

Maybe there are other options?

These changes are needed, because bmi is calculated directly from height and weight. So using these variables together results in multi-collinearity problems with the model.

e. Use the adjusted predictor matrix to impute the boys dataset again, with 10 imputations and 10 iterations. Inspect the iteration plots again, do you see improvements for hgt, wgt and bmi?

```
imp1 <- mice(boys, m = 10, maxit = 10, pred = pred1)
```

```
>
> iter imp variable
> 1 1 hgt wgt bmi hc gen phb tv reg
> 1 2 hgt wgt bmi hc gen phb tv reg
> 1 3 hgt wgt bmi hc gen phb tv reg
> 1 4 hgt wgt bmi hc gen phb tv reg
> 1 5 hgt wgt bmi hc gen phb tv reg
> 1 6 hgt wgt bmi hc gen phb tv reg
> 1 7 hgt wgt bmi hc gen phb tv reg
> 1 8 hgt wgt bmi hc gen phb tv reg
> 1 9 hgt wgt bmi hc gen phb tv reg
> 1 10 hgt wgt bmi hc gen phb tv reg
> 2 1 hgt wgt bmi hc gen phb tv reg
> 2 2 hgt wgt bmi hc gen phb tv reg
> 2 3 hgt wgt bmi hc gen phb tv reg
> 2 4 hgt wgt bmi hc gen phb tv reg
> 2 5 hgt wgt bmi hc gen phb tv reg
> 2 6 hgt wgt bmi hc gen phb tv reg
> 2 7 hgt wgt bmi hc gen phb tv reg
> 2 8 hgt wgt bmi hc gen phb tv reg
> 2 9 hgt wgt bmi hc gen phb tv reg
> 2 10 hgt wgt bmi hc gen phb tv reg
> 3 1 hgt wgt bmi hc gen phb tv reg
> 3 2 hgt wgt bmi hc gen phb tv reg
> 3 3 hgt wgt bmi hc gen phb tv reg
> 3 4 hgt wgt bmi hc gen phb tv reg
> 3 5 hgt wgt bmi hc gen phb tv reg
> 3 6 hgt wgt bmi hc gen phb tv reg
> 3 7 hgt wgt bmi hc gen phb tv reg
> 3 8 hgt wgt bmi hc gen phb tv reg
> 3 9 hgt wgt bmi hc gen phb tv reg
> 3 10 hgt wgt bmi hc gen phb tv reg
> 4 1 hgt wgt bmi hc gen phb tv reg
> 4 2 hgt wgt bmi hc gen phb tv reg
> 4 3 hgt wgt bmi hc gen phb tv reg
> 4 4 hgt wgt bmi hc gen phb tv reg
> 4 5 hgt wgt bmi hc gen phb tv reg
> 4 6 hgt wgt bmi hc gen phb tv reg
> 4 7 hgt wgt bmi hc gen phb tv reg
> 4 8 hgt wgt bmi hc gen phb tv reg
> 4 9 hgt wgt bmi hc gen phb tv reg
> 4 10 hgt wgt bmi hc gen phb tv reg
> 5 1 hgt wgt bmi hc gen phb tv reg
> 5 2 hgt wgt bmi hc gen phb tv reg
```

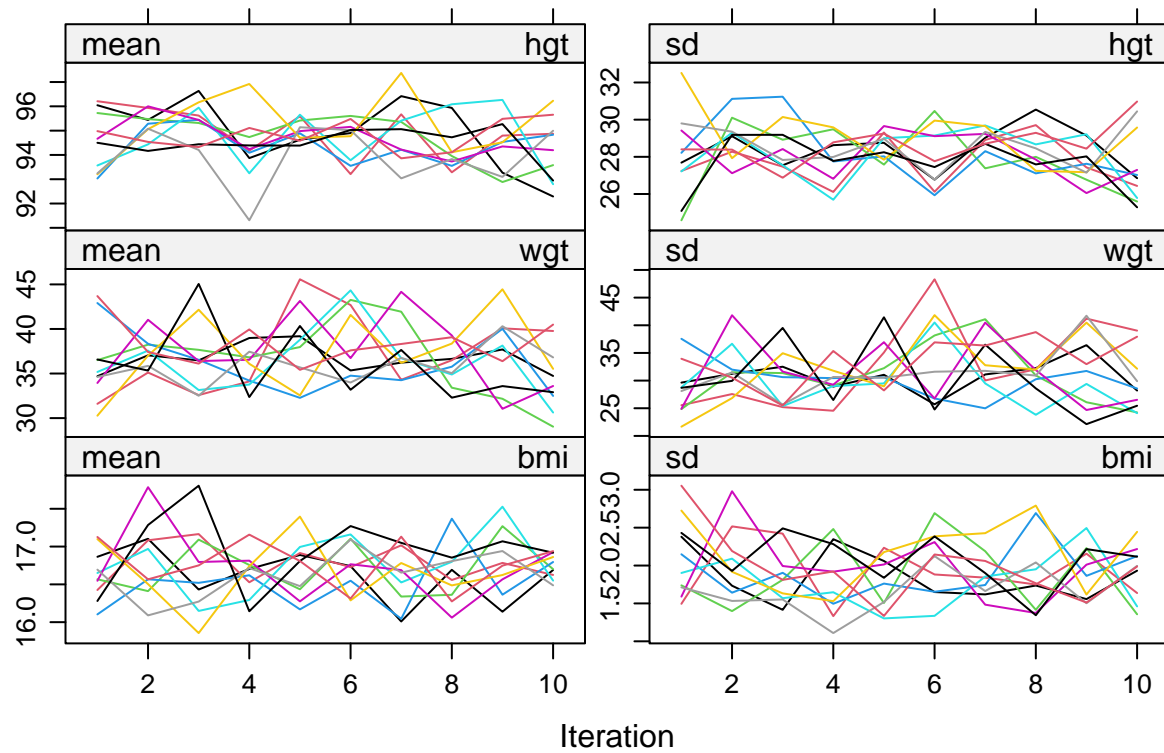
[illegible]

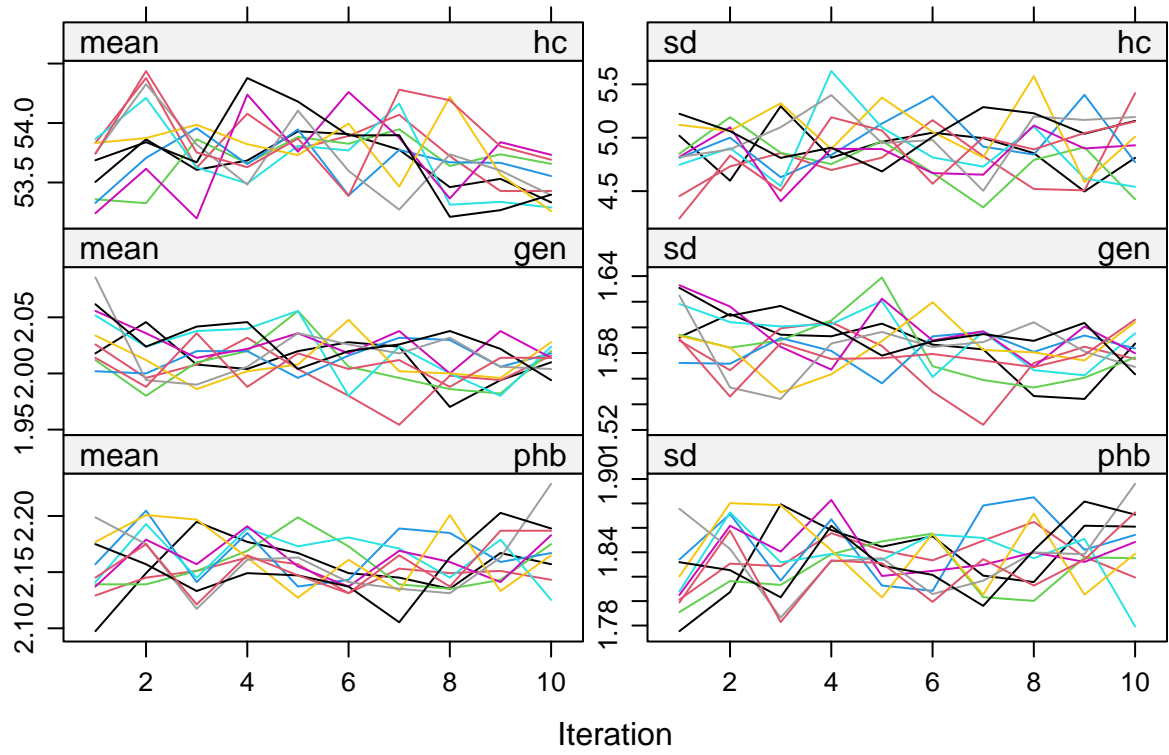

```

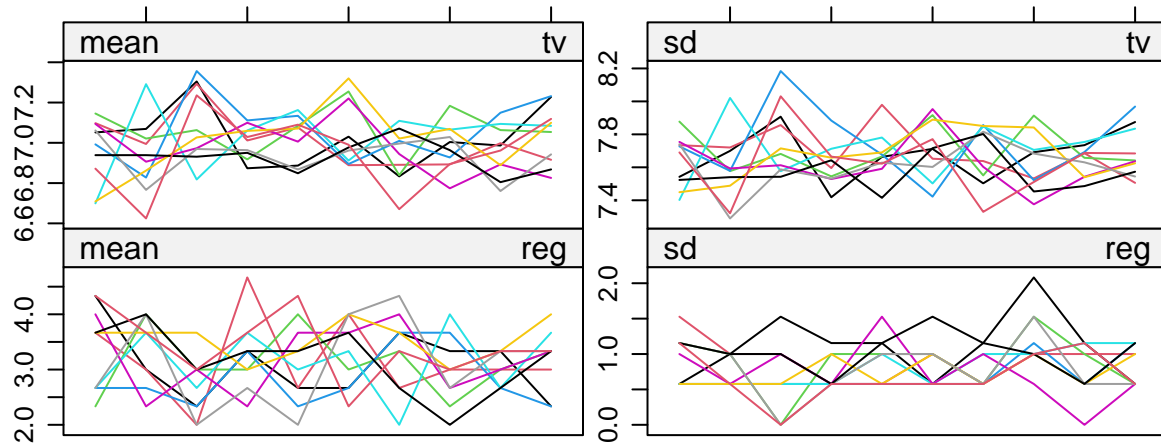
> 10 7 hgt wgt bmi hc gen phb tv reg
> 10 8 hgt wgt bmi hc gen phb tv reg
> 10 9 hgt wgt bmi hc gen phb tv reg
> 10 10 hgt wgt bmi hc gen phb tv reg

```

```
plot(imp1)
```







Iteration

```
imp2 <- mice(boys, m = 10, maxit = 10, pred = pred2)
```

```
>
> iter imp variable
> 1 1 hgt wgt bmi hc gen phb tv reg
> 1 2 hgt wgt bmi hc gen phb tv reg
> 1 3 hgt wgt bmi hc gen phb tv reg
> 1 4 hgt wgt bmi hc gen phb tv reg
> 1 5 hgt wgt bmi hc gen phb tv reg
> 1 6 hgt wgt bmi hc gen phb tv reg
> 1 7 hgt wgt bmi hc gen phb tv reg
> 1 8 hgt wgt bmi hc gen phb tv reg
> 1 9 hgt wgt bmi hc gen phb tv reg
> 1 10 hgt wgt bmi hc gen phb tv reg
> 2 1 hgt wgt bmi hc gen phb tv reg
> 2 2 hgt wgt bmi hc gen phb tv reg
> 2 3 hgt wgt bmi hc gen phb tv reg
> 2 4 hgt wgt bmi hc gen phb tv reg
> 2 5 hgt wgt bmi hc gen phb tv reg
> 2 6 hgt wgt bmi hc gen phb tv reg
> 2 7 hgt wgt bmi hc gen phb tv reg
> 2 8 hgt wgt bmi hc gen phb tv reg
> 2 9 hgt wgt bmi hc gen phb tv reg
> 2 10 hgt wgt bmi hc gen phb tv reg
> 3 1 hgt wgt bmi hc gen phb tv reg
```

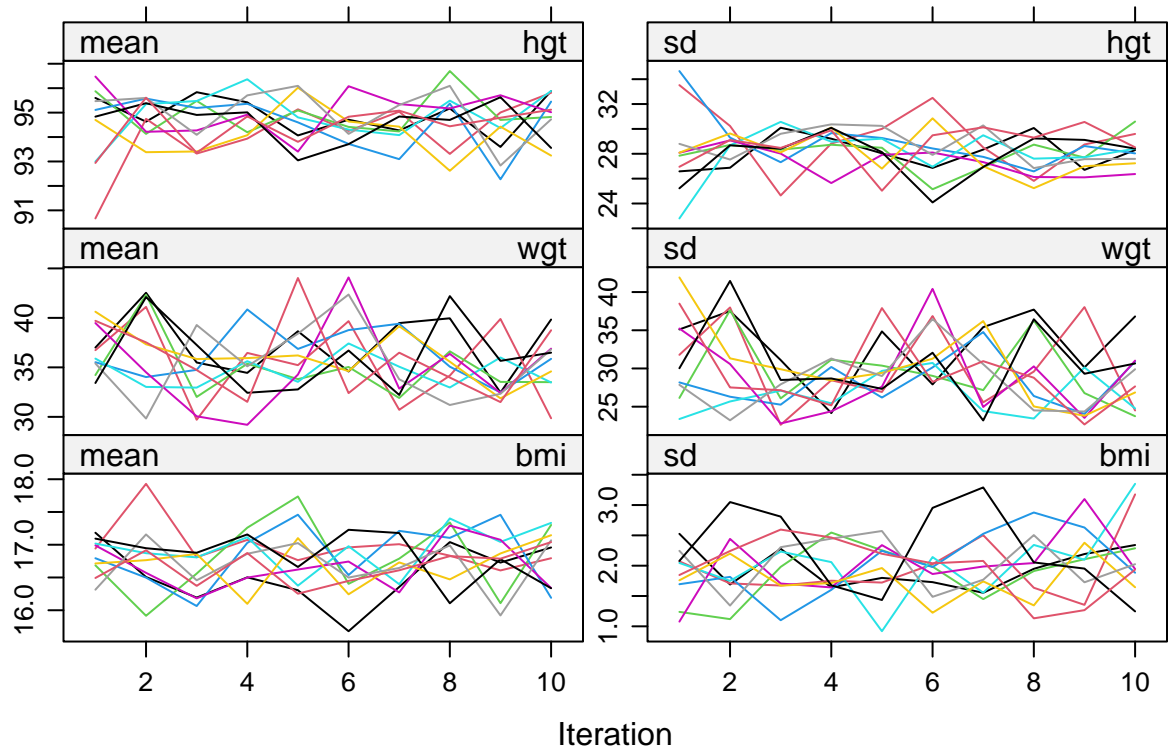
[illegible]

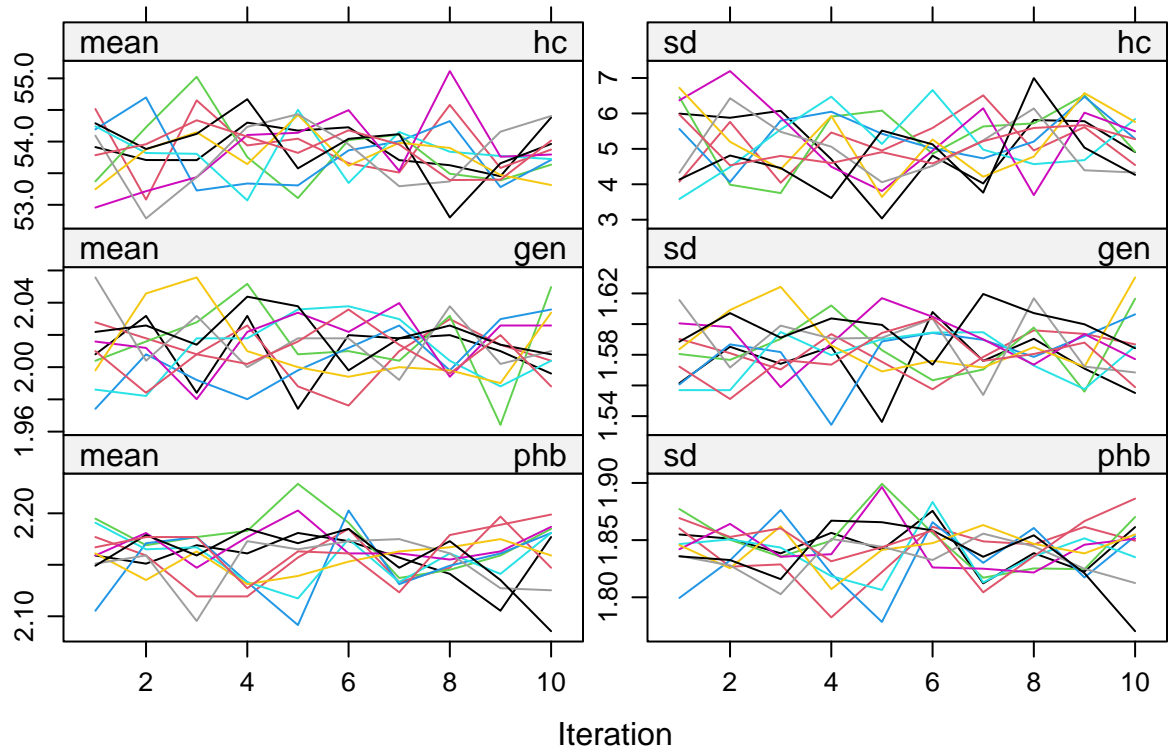
```

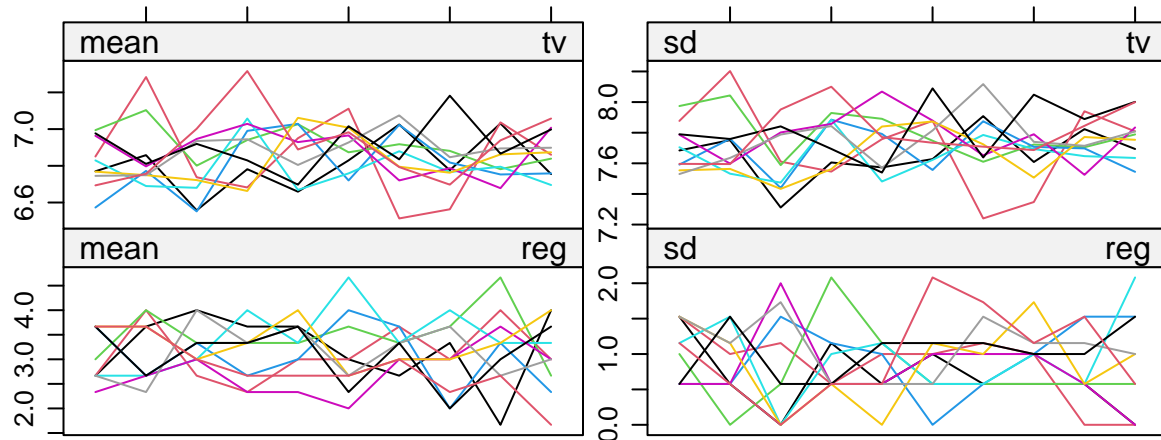
> 8 6 hgt wgt bmi hc gen phb tv reg
> 8 7 hgt wgt bmi hc gen phb tv reg
> 8 8 hgt wgt bmi hc gen phb tv reg
> 8 9 hgt wgt bmi hc gen phb tv reg
> 8 10 hgt wgt bmi hc gen phb tv reg
> 9 1 hgt wgt bmi hc gen phb tv reg
> 9 2 hgt wgt bmi hc gen phb tv reg
> 9 3 hgt wgt bmi hc gen phb tv reg
> 9 4 hgt wgt bmi hc gen phb tv reg
> 9 5 hgt wgt bmi hc gen phb tv reg
> 9 6 hgt wgt bmi hc gen phb tv reg
> 9 7 hgt wgt bmi hc gen phb tv reg
> 9 8 hgt wgt bmi hc gen phb tv reg
> 9 9 hgt wgt bmi hc gen phb tv reg
> 9 10 hgt wgt bmi hc gen phb tv reg
> 10 1 hgt wgt bmi hc gen phb tv reg
> 10 2 hgt wgt bmi hc gen phb tv reg
> 10 3 hgt wgt bmi hc gen phb tv reg
> 10 4 hgt wgt bmi hc gen phb tv reg
> 10 5 hgt wgt bmi hc gen phb tv reg
> 10 6 hgt wgt bmi hc gen phb tv reg
> 10 7 hgt wgt bmi hc gen phb tv reg
> 10 8 hgt wgt bmi hc gen phb tv reg
> 10 9 hgt wgt bmi hc gen phb tv reg
> 10 10 hgt wgt bmi hc gen phb tv reg

```

```
plot(imp2)
```







Iteration

f. Another way to deal with the relation between hgt and wgt with bmi is to use “passive imputation” to impute the bmi variable. Adjust the method and predictor matrix in such a way that the bmi variable is passively imputed by hgt and wgt. Inspect the iteration plots again, do you see improvements for hgt, wgt and bmi?

The formula for bmi is $bmi = \frac{weight}{height_{meters}^2}$. This formula is added as imputation method to compute missing bmi values, from the imputed weight and height values. We use the second option of the predictor matrix (option 1 is also possible), so that the updated bmi variable (from imputed hgt and wgt) is used as predictor for other variables. Option 1 might be preferred when bmi is used in the substantial analyses.

```
method <- imp$method
method["bmi"] <- "~I(wgt / (hgt/100)^2)"

imp <- mice(boys, m = 10, maxit = 10, pred = pred2, method = method)
```

```
>
> iter imp variable
> 1 1 hgt wgt bmi hc gen phb tv reg
> 1 2 hgt wgt bmi hc gen phb tv reg
> 1 3 hgt wgt bmi hc gen phb tv reg
> 1 4 hgt wgt bmi hc gen phb tv reg
> 1 5 hgt wgt bmi hc gen phb tv reg
> 1 6 hgt wgt bmi hc gen phb tv reg
> 1 7 hgt wgt bmi hc gen phb tv reg
> 1 8 hgt wgt bmi hc gen phb tv reg
> 1 9 hgt wgt bmi hc gen phb tv reg
```



```

> 1 10 hgt wgt bmi hc gen phb tv reg
> 2 1 hgt wgt bmi hc gen phb tv reg
> 2 2 hgt wgt bmi hc gen phb tv reg
> 2 3 hgt wgt bmi hc gen phb tv reg
> 2 4 hgt wgt bmi hc gen phb tv reg
> 2 5 hgt wgt bmi hc gen phb tv reg
> 2 6 hgt wgt bmi hc gen phb tv reg
> 2 7 hgt wgt bmi hc gen phb tv reg
> 2 8 hgt wgt bmi hc gen phb tv reg
> 2 9 hgt wgt bmi hc gen phb tv reg
> 2 10 hgt wgt bmi hc gen phb tv reg
> 3 1 hgt wgt bmi hc gen phb tv reg
> 3 2 hgt wgt bmi hc gen phb tv reg
> 3 3 hgt wgt bmi hc gen phb tv reg
> 3 4 hgt wgt bmi hc gen phb tv reg
> 3 5 hgt wgt bmi hc gen phb tv reg
> 3 6 hgt wgt bmi hc gen phb tv reg
> 3 7 hgt wgt bmi hc gen phb tv reg
> 3 8 hgt wgt bmi hc gen phb tv reg
> 3 9 hgt wgt bmi hc gen phb tv reg
> 3 10 hgt wgt bmi hc gen phb tv reg
> 4 1 hgt wgt bmi hc gen phb tv reg
> 4 2 hgt wgt bmi hc gen phb tv reg
> 4 3 hgt wgt bmi hc gen phb tv reg
> 4 4 hgt wgt bmi hc gen phb tv reg
> 4 5 hgt wgt bmi hc gen phb tv reg
> 4 6 hgt wgt bmi hc gen phb tv reg
> 4 7 hgt wgt bmi hc gen phb tv reg
> 4 8 hgt wgt bmi hc gen phb tv reg
> 4 9 hgt wgt bmi hc gen phb tv reg
> 4 10 hgt wgt bmi hc gen phb tv reg
> 5 1 hgt wgt bmi hc gen phb tv reg
> 5 2 hgt wgt bmi hc gen phb tv reg
> 5 3 hgt wgt bmi hc gen phb tv reg
> 5 4 hgt wgt bmi hc gen phb tv reg
> 5 5 hgt wgt bmi hc gen phb tv reg
> 5 6 hgt wgt bmi hc gen phb tv reg
> 5 7 hgt wgt bmi hc gen phb tv reg
> 5 8 hgt wgt bmi hc gen phb tv reg
> 5 9 hgt wgt bmi hc gen phb tv reg
> 5 10 hgt wgt bmi hc gen phb tv reg
> 6 1 hgt wgt bmi hc gen phb tv reg
> 6 2 hgt wgt bmi hc gen phb tv reg
> 6 3 hgt wgt bmi hc gen phb tv reg
> 6 4 hgt wgt bmi hc gen phb tv reg
> 6 5 hgt wgt bmi hc gen phb tv reg
> 6 6 hgt wgt bmi hc gen phb tv reg
> 6 7 hgt wgt bmi hc gen phb tv reg
> 6 8 hgt wgt bmi hc gen phb tv reg
> 6 9 hgt wgt bmi hc gen phb tv reg
> 6 10 hgt wgt bmi hc gen phb tv reg
> 7 1 hgt wgt bmi hc gen phb tv reg
> 7 2 hgt wgt bmi hc gen phb tv reg
> 7 3 hgt wgt bmi hc gen phb tv reg

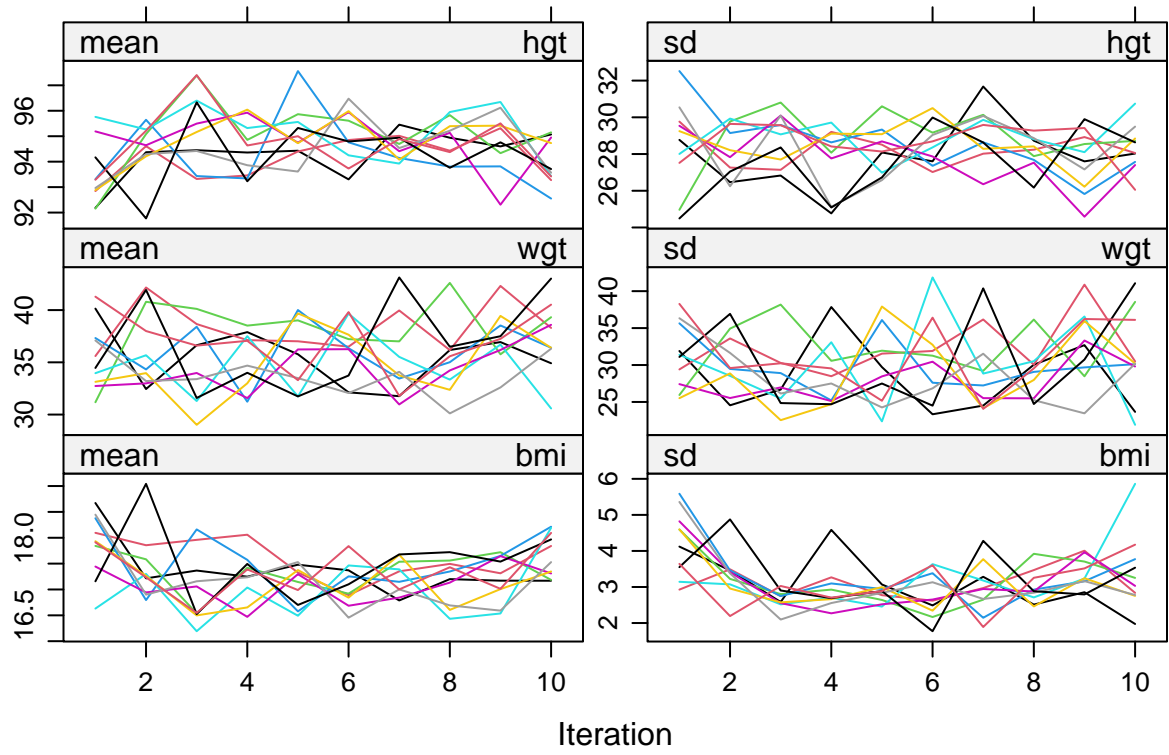
```

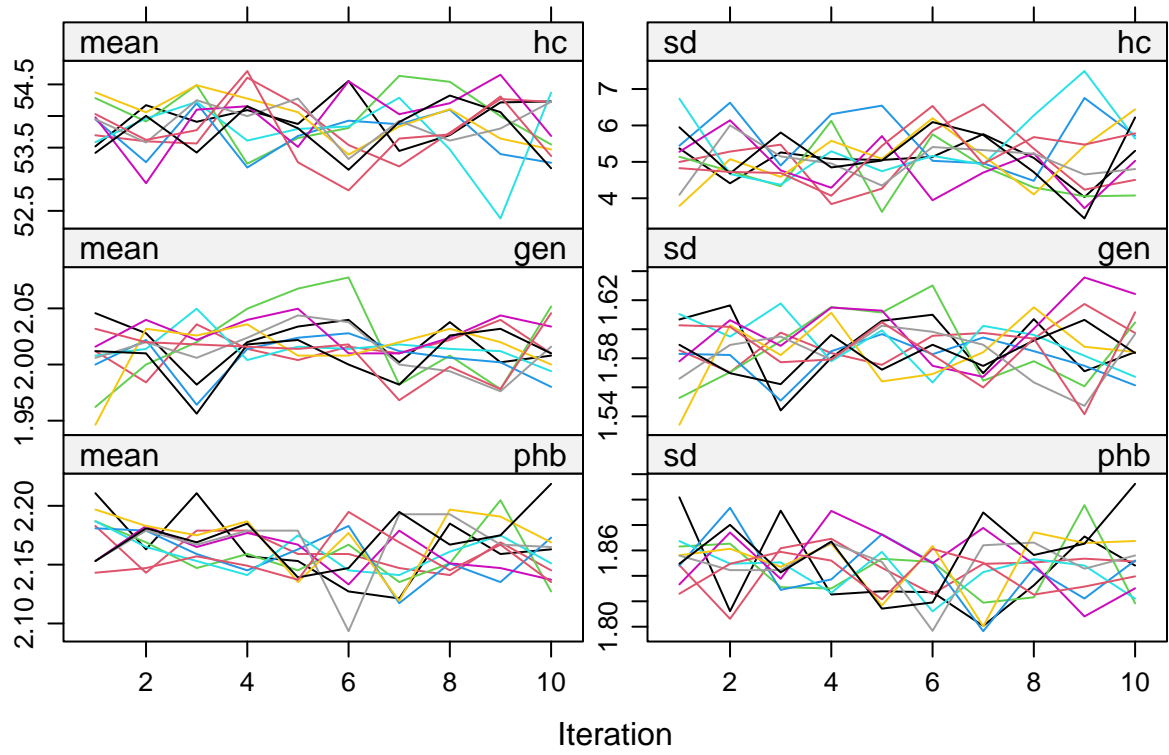
```

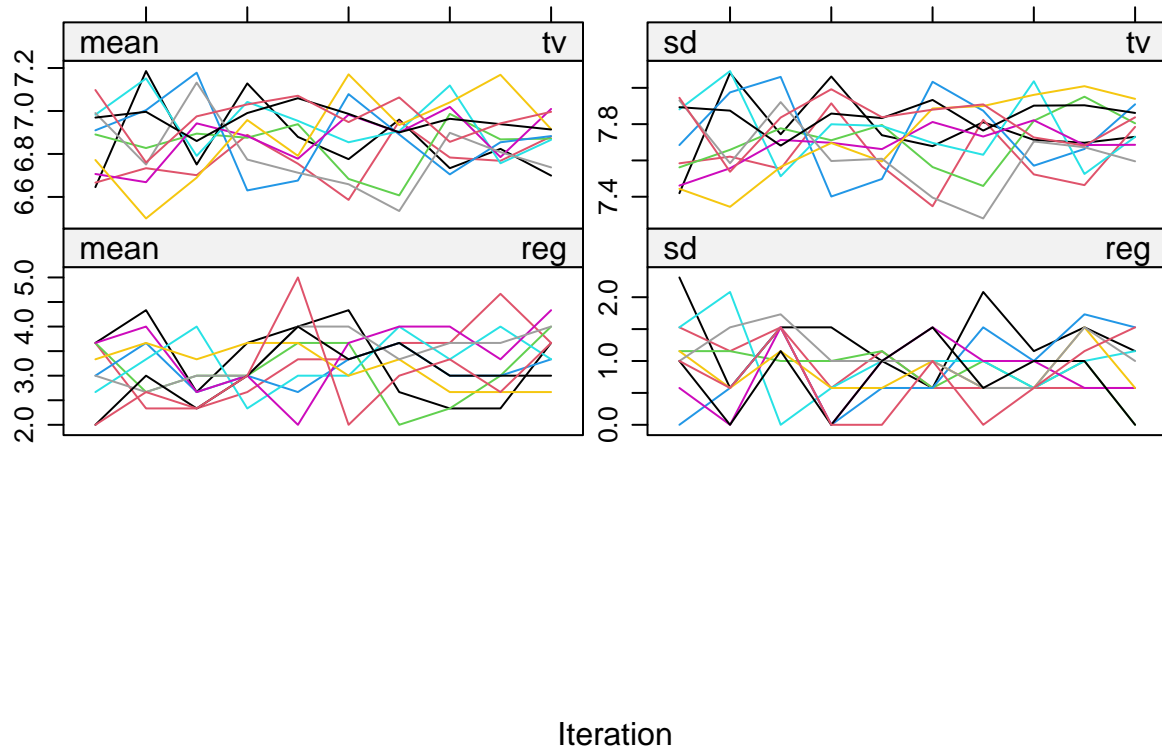
> 7 4 hgt wgt bmi hc gen phb tv reg
> 7 5 hgt wgt bmi hc gen phb tv reg
> 7 6 hgt wgt bmi hc gen phb tv reg
> 7 7 hgt wgt bmi hc gen phb tv reg
> 7 8 hgt wgt bmi hc gen phb tv reg
> 7 9 hgt wgt bmi hc gen phb tv reg
> 7 10 hgt wgt bmi hc gen phb tv reg
> 8 1 hgt wgt bmi hc gen phb tv reg
> 8 2 hgt wgt bmi hc gen phb tv reg
> 8 3 hgt wgt bmi hc gen phb tv reg
> 8 4 hgt wgt bmi hc gen phb tv reg
> 8 5 hgt wgt bmi hc gen phb tv reg
> 8 6 hgt wgt bmi hc gen phb tv reg
> 8 7 hgt wgt bmi hc gen phb tv reg
> 8 8 hgt wgt bmi hc gen phb tv reg
> 8 9 hgt wgt bmi hc gen phb tv reg
> 8 10 hgt wgt bmi hc gen phb tv reg
> 9 1 hgt wgt bmi hc gen phb tv reg
> 9 2 hgt wgt bmi hc gen phb tv reg
> 9 3 hgt wgt bmi hc gen phb tv reg
> 9 4 hgt wgt bmi hc gen phb tv reg
> 9 5 hgt wgt bmi hc gen phb tv reg
> 9 6 hgt wgt bmi hc gen phb tv reg
> 9 7 hgt wgt bmi hc gen phb tv reg
> 9 8 hgt wgt bmi hc gen phb tv reg
> 9 9 hgt wgt bmi hc gen phb tv reg
> 9 10 hgt wgt bmi hc gen phb tv reg
> 10 1 hgt wgt bmi hc gen phb tv reg
> 10 2 hgt wgt bmi hc gen phb tv reg
> 10 3 hgt wgt bmi hc gen phb tv reg
> 10 4 hgt wgt bmi hc gen phb tv reg
> 10 5 hgt wgt bmi hc gen phb tv reg
> 10 6 hgt wgt bmi hc gen phb tv reg
> 10 7 hgt wgt bmi hc gen phb tv reg
> 10 8 hgt wgt bmi hc gen phb tv reg
> 10 9 hgt wgt bmi hc gen phb tv reg
> 10 10 hgt wgt bmi hc gen phb tv reg

```

```
plot(imp)
```







Solution 7: multiple imputation iteration plots