

6. Model Diagnosis on Arsentic dataset

```
library(tidyverse)
library(ggplot2)
library(olsrr)
library(haven)
library(car)

# data preperation
arsenic <- read_dta("arsenic.dta")
arsenic$id <- c(1:21)

summary(arsenic)

##           age           sex           drinkuse           cookuse
## Min.      : 8.00   Length:21   Min.      :1.000   Min.      :2.000
## 1st Qu.:41.00   Class :character   1st Qu.:4.000   1st Qu.:5.000
## Median :45.00   Mode  :character   Median :5.000   Median :5.000
## Mean    :47.57                      Mean    :4.333   Mean    :4.857
## 3rd Qu.:53.00                      3rd Qu.:5.000   3rd Qu.:5.000
## Max.    :86.00                      Max.    :5.000   Max.    :5.000
##      arswater      arsnails      id
## Min.      :0.00000   Min.      :0.0730   Min.      : 1
## 1st Qu.:0.00000   1st Qu.:0.1180   1st Qu.: 6
## Median :0.00100   Median :0.1750   Median :11
## Mean    :0.01624   Mean    :0.3664   Mean    :11
## 3rd Qu.:0.01800   3rd Qu.:0.3580   3rd Qu.:16
## Max.    :0.13700   Max.    :2.2520   Max.    :21

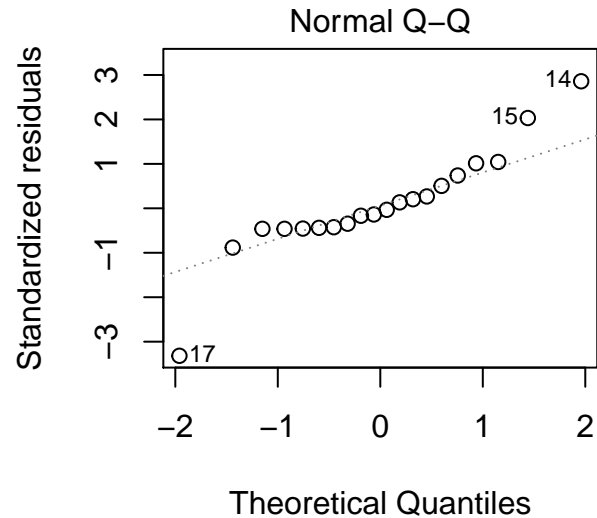
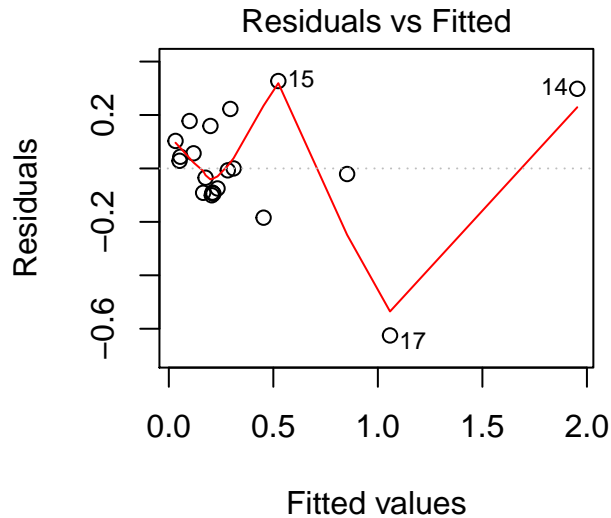
mod <- lm(arsnails ~ age + sex + drinkuse + cookuse + arswater + arswater, data = arsenic)
summary(mod)

##
## Call:
## lm(formula = arsnails ~ age + sex + drinkuse + cookuse + arswater +
##      arswater, data = arsenic)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.62510 -0.09117 -0.00714  0.10297  0.32719
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.452972   0.418132   1.083   0.296
## age         -0.001290   0.003444  -0.374   0.713
## sexMale     -0.145038   0.107448  -1.350   0.197
## drinkuse    -0.011719   0.047010  -0.249   0.807
## cookuse     -0.027471   0.082861  -0.332   0.745
## arswater    13.195586   1.639792   8.047 8.01e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2302 on 15 degrees of freedom
## Multiple R-squared:  0.8323, Adjusted R-squared:  0.7764
## F-statistic: 14.89 on 5 and 15 DF,  p-value: 2.339e-05
```

Diagnostic plots

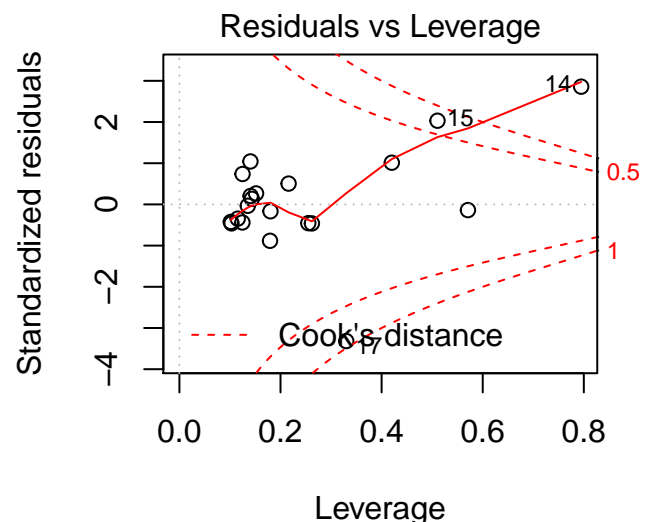
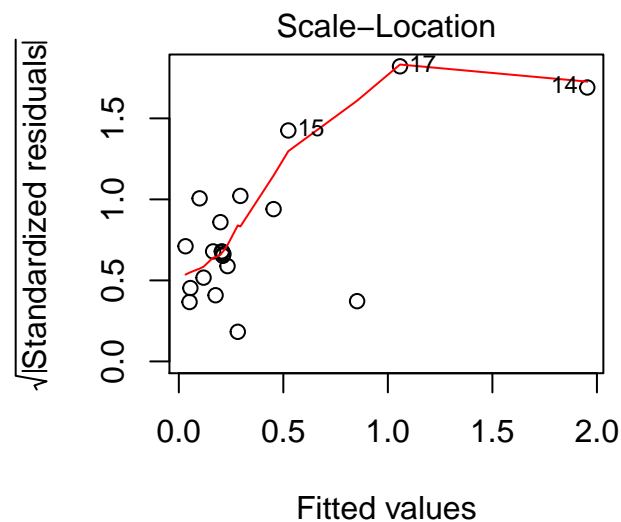
```
plot(mod)
```

```
## Warning: not plotting observations with leverage one:
##      9
```



```
ails ~ age + sex + drinkuse + cookuse + arswateails ~ age + sex + drinkuse + cookuse + arswate
```

```
## Warning: not plotting observations with leverage one:
##      9
```



```
ails ~ age + sex + drinkuse + cookuse + arswateails ~ age + sex + drinkuse + cookuse + arswate
```

From the residuals - fitted value plot, we can see a existence of heteroscedasticity, the residual variance tends to be larger when the fitted values are larger. Also, we can see three abnormal cases, 14, 15, and 17, with both larger studentized residuals and cook's distance, indicating them to be outlying and highly influential. Besides, there the 9th case has leverage 1, being an abnormal case as well.

```
arsenic[c(9,14,15,17),]
```

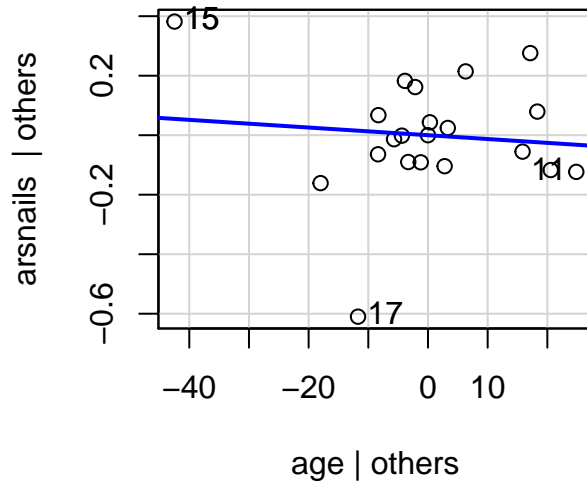
```
## # A tibble: 4 x 7
##   age sex  drinkuse cookuse arswater arsnails id
##   <dbl> <chr>    <dbl>    <dbl>    <dbl>    <dbl> <int>
```

## 1	41 Female	3	2	0	0.310	9
## 2	86 Female	5	5	0.137	2.25	14
## 3	8 Female	5	5	0.0210	0.851	15
## 4	44 Male	5	5	0.0760	0.433	17

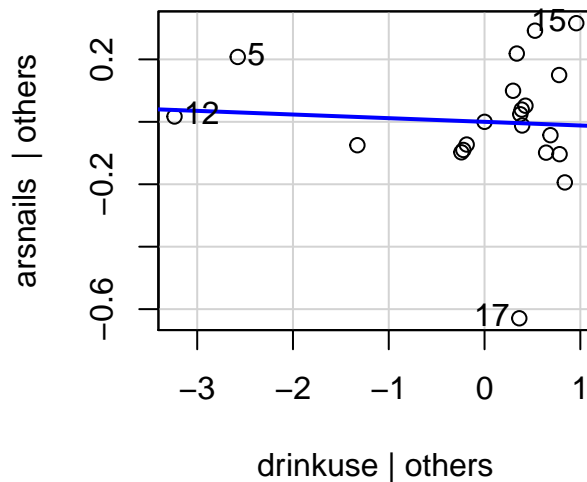
Partial Residual Plots

To analyze whether the predictors are useful, and if there is nonlinearity in each predictor, we then made partial residual plots for covariates. For the sex group, we also made a group wise boxplot to detect outliers.

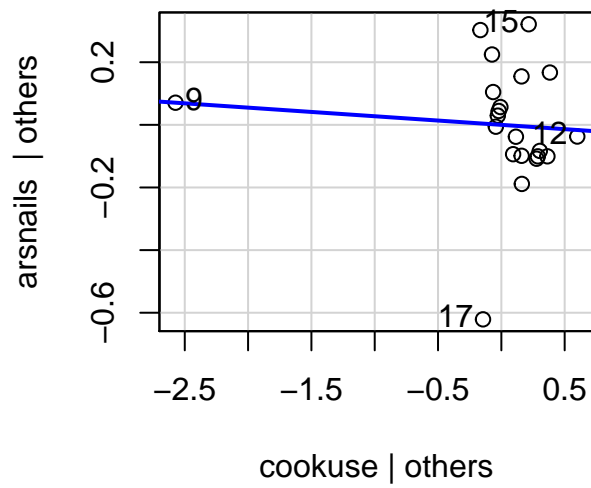
```
avPlots(mod, ~ age)
```



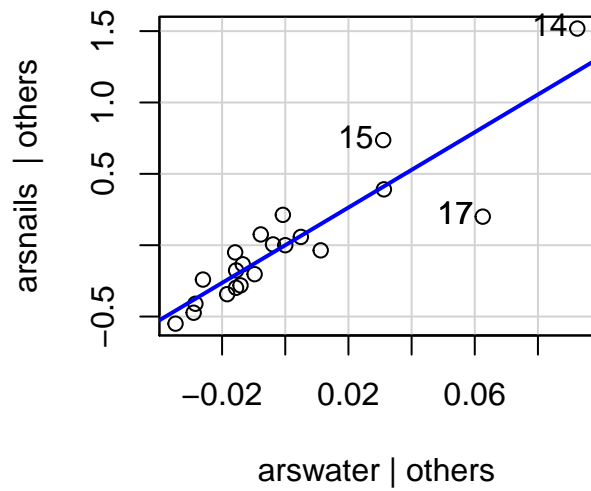
```
avPlots(mod, ~ drinkuse)
```



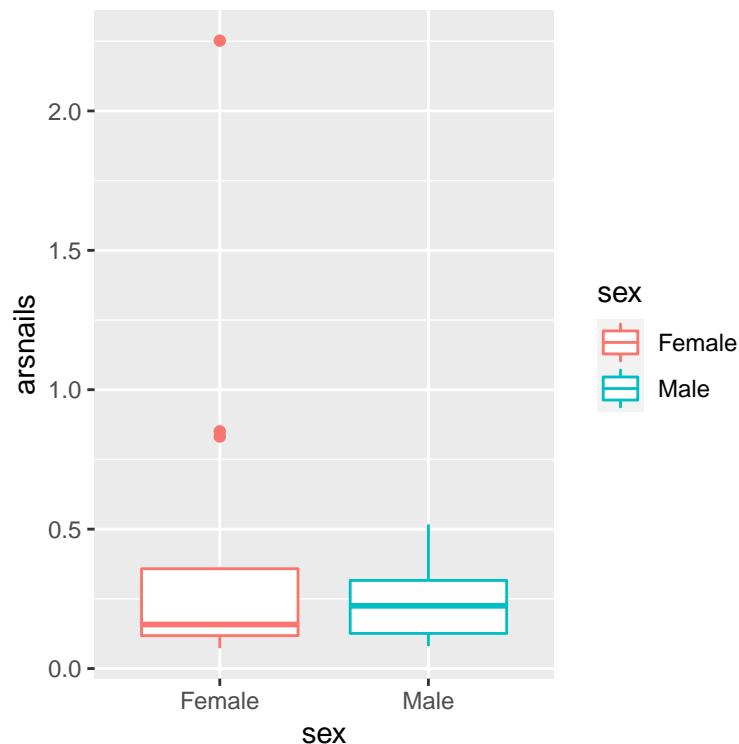
```
avPlots(mod, ~ cookuse)
```



```
avPlots(mod, ~ arswater)
```



```
ggplot(arsenic, aes(sex, arsnails, color = sex)) +  
  geom_boxplot()
```



- Case 15 has an abnormally low age, also a high toe nail arsenic concentration, being an outlier with also high influence on model.
- Case 9 has an unusual frequency of cookuse, as it is the only case in the dataset to have 2 for cookuse, and the rest of the cases all have value 5.
- Case 14, 17 have an abnormally high outcome, which is the main reason why they are influential. However, according to the partial residual plots, the positive linear relationship between toenail arsenic concentration with well water arsenic concentration is not completely driven only by these outliers.

Regression without outlier

```
mod1 <- lm(arsnails ~ age + sex + drinkuse + arswater + arswater, data = arsenic[-c(9,14,15,17),])
summary(mod1)
```

```
##
## Call:
## lm(formula = arsnails ~ age + sex + drinkuse + arswater + arswater,
##     data = arsenic[-c(9, 14, 15, 17), ])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.084450 -0.040869 -0.015782  0.009797  0.239076
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.2681601  0.1493702   1.795   0.0978 .
## age         -0.0002232  0.0021761  -0.103   0.9200
## sexMale      0.0183767  0.0485070   0.379   0.7114
## drinkuse    -0.0278381  0.0224720  -1.239   0.2391
## arswater     12.8679620  2.2049209   5.836 8.01e-05 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.09439 on 12 degrees of freedom
## Multiple R-squared:  0.8245, Adjusted R-squared:  0.766
## F-statistic: 14.1 on 4 and 12 DF,  p-value: 0.0001737
```

When we excluded all the high influential cases, the regression result stays the same, indicating a robust positive linear relationship between toenail arsenic concentration and the arsenic concentration in their well-water. And the other factors, age, sex, frequency of drink and cook use remain insignificant.

Summary

In the analysis above, we studied the ingestion of arsenic-containing water of 21 people and what factors have impact on the ingestion. We built up multiple linear regression model with toenail arsenic concentration as outcome, and had age, sex, frequency of private well water for drinking and cooking, and the arsenic concentration in their well water as predictors. The model revealed a significant positive linear relationship between the arsenic concentration in well water and toenail arsenic concentration.

Within the 21 participants, 13 of them are female and 8 are male. Most of the participants are aged from 40 - 60, while case 9 is very young of age 8, and case 11 is very old of age 86. Most participants except for 3 reported a high drink use of 4 or 5, and all participants except for 1 reported a very high cook use of 5. Arsenic concentration were detected in 15 of the 21 well-water samples from the participants, but we still include the 6 cases with no arsenic concentration detected, imputing a value of 0. 3 participants had abnormally high toenail arsenic concentration, and all of them are female.

The model results and direction of linear relationship remained valid when we excluded all influential outlying cases. Therefore, despite the extreme data points and the small sample size, our conclusion that toenail arsenic concentration has a significant positive linear relationship with the arsenic concentration of participants' well water arsenic concentration is very stable.