

BIFX 553 - Discussion 3

Randy Johnson

February 2, 2017

Assessing Model Fit and Assumptions

Regression Assumptions

We will primarily use visual inspection and the `car` package for checking regression assumptions, but there are many other resources in R to do this (e.g. the `rms` and `gvlma` packages).

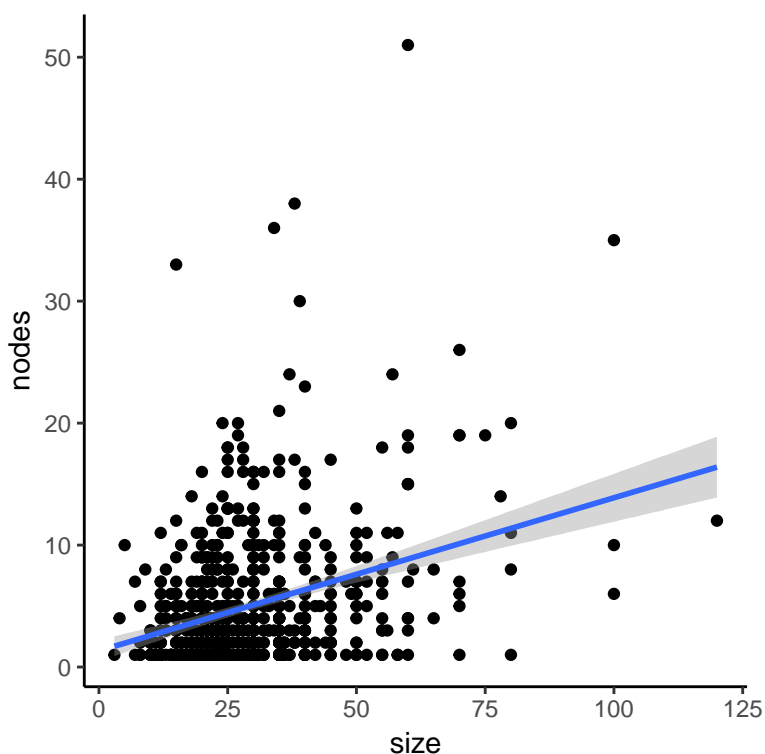
- Linear relationship
- Multivariate Normality
- No/little multicollinearity
- No autocorrelation
- Homoscedasticity

How does our model hold up?

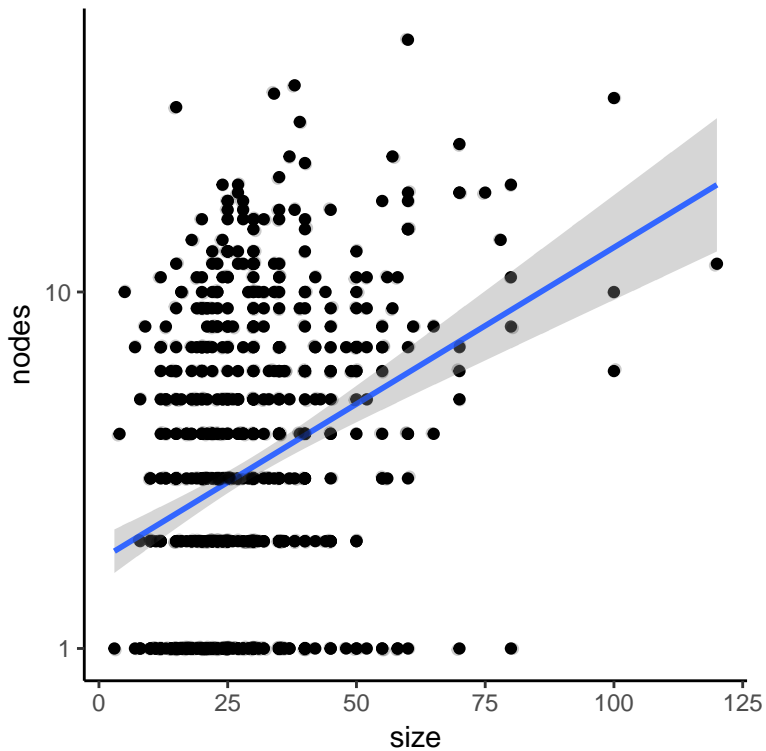
```
load('../1-26/gbsg.RData')

# model from Discussion2
gbsg.lm <- lm(nodes ~ age + meno + size + grade + pgr + er + hormon, data = gbsg)

ggplot(gbsg, aes(size, nodes)) +
  geom_point() +
  geom_smooth(method = 'lm')
```



```
# this is probably a better way to look at the data
ggplot(gbsg, aes(size, nodes)) +
  geom_point() +
  geom_smooth(method = 'lm') +
  geom_jitter(alpha = .2) +
  scale_y_log10()
```

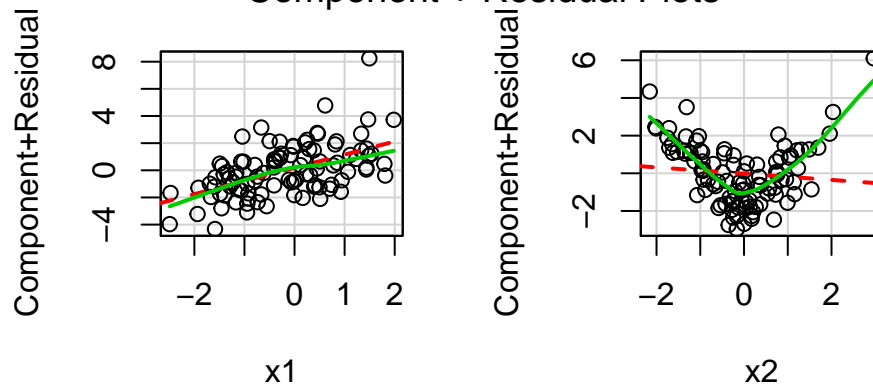


Linear Relationship

Component residual plots are a way to see if the predictors have a linear relationship with the outcome variable. The red, dashed line in the figures below represents the best fit of each predictor and the residuals, and the solid, green line is a running, smoothed average along the x-axis. In this example,

```
tmp <- data_frame(x1 = rnorm(100),
                  x2 = rnorm(100),
                  y = x1 + x2^2 + rnorm(100))
lm(y ~ x1 + x2, data = tmp) %>%
  crPlots()
```

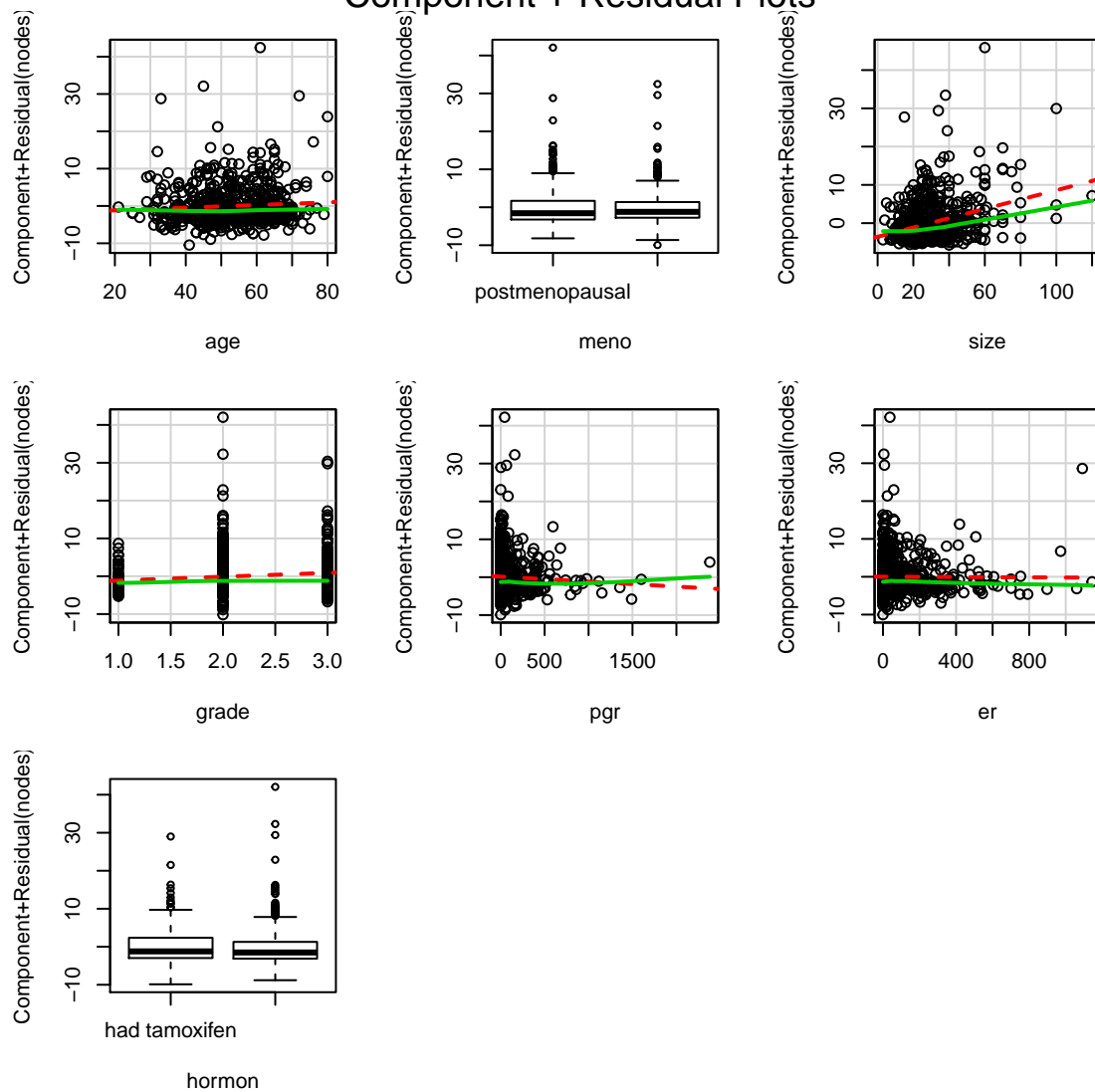
Component + Residual Plots



How does our model look?

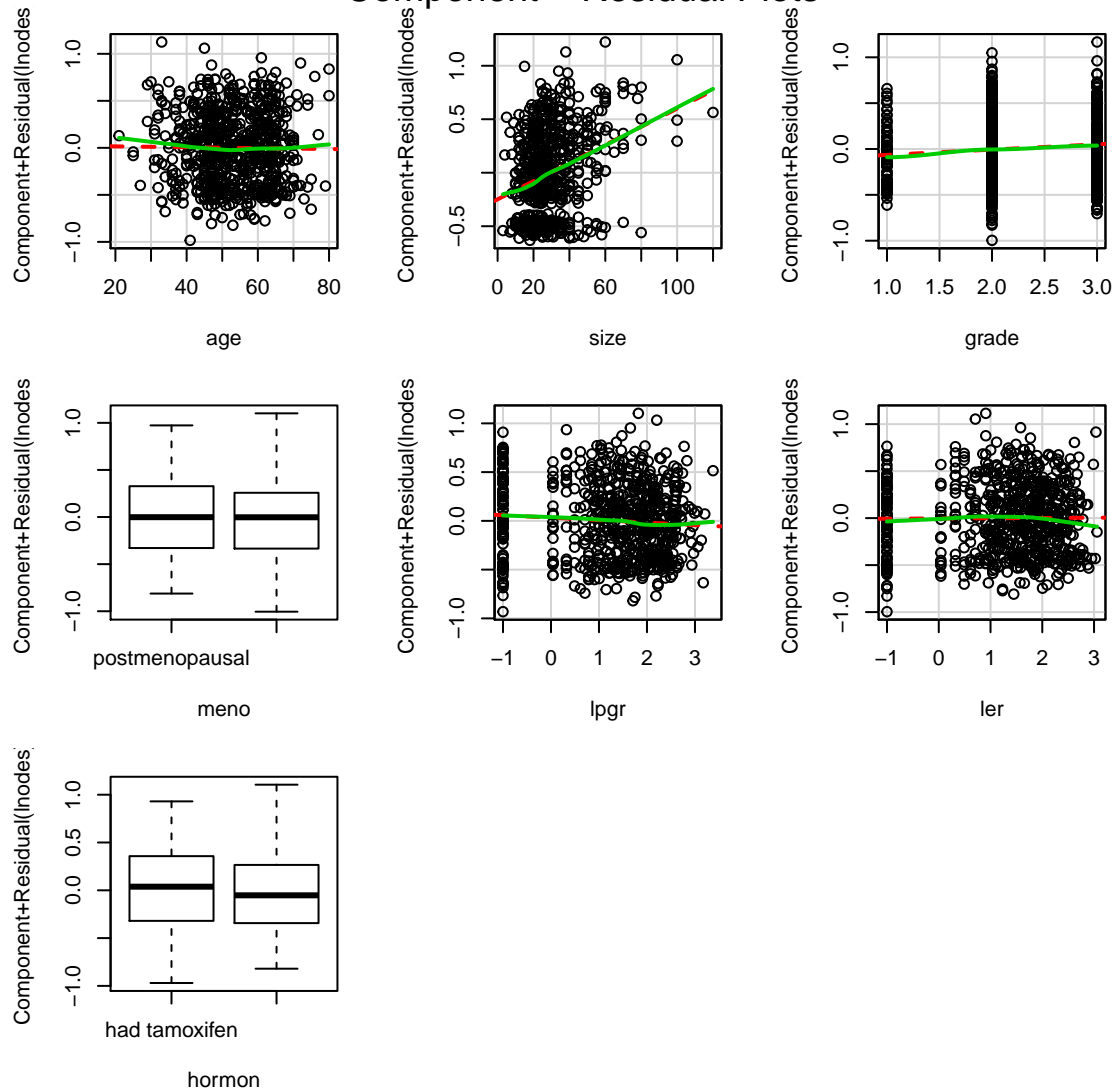
```
# old model
crPlots(gbsg.lm)
```

Component + Residual Plots



```
# new model
crPlots(gbsg.lm2)
```

Component + Residual Plots



Multivariate Normality

There are a couple of ways we can look at normality. The Shapiro-Wilk test for normality will give you a quantitative measure of whether your residuals are normally distributed, and the QQ plot will give you a graphical way to see what is going on with your residuals.

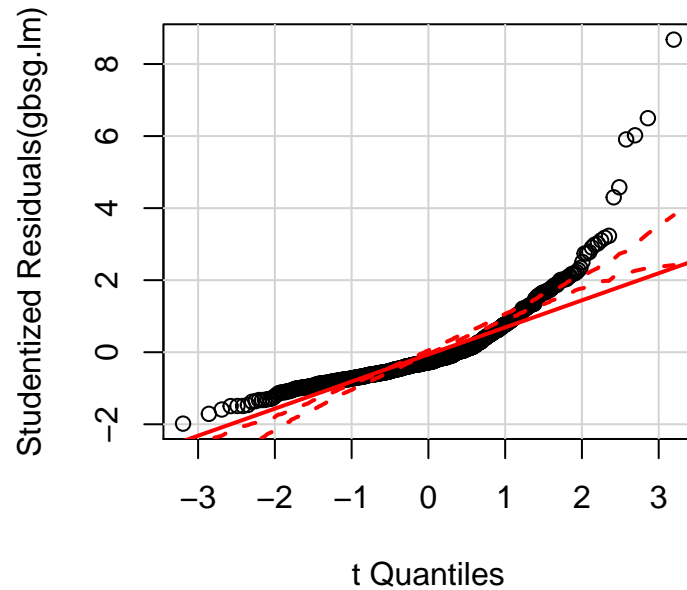
```
# Shapiro-Wilk test for normality
with(augment(gbsg.lm),
      shapiro.test(.std.resid))
```

```
##
##  Shapiro-Wilk normality test
##
## data:  .std.resid
## W = 0.79776, p-value < 2.2e-16
```

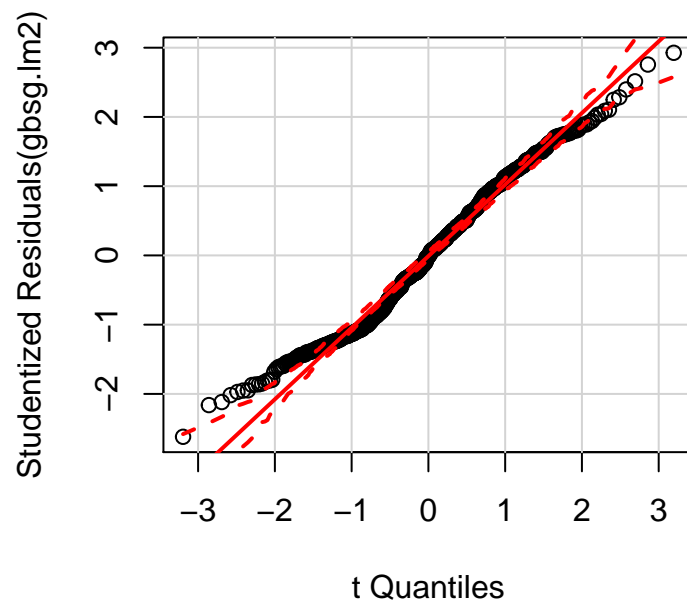
```
with(augment(gbsg.lm2),  
      shapiro.test(.std.resid))
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  .std.resid  
## W = 0.98328, p-value = 4.633e-07
```

```
# quantile quantile plot  
qqPlot(gbsg.lm)
```



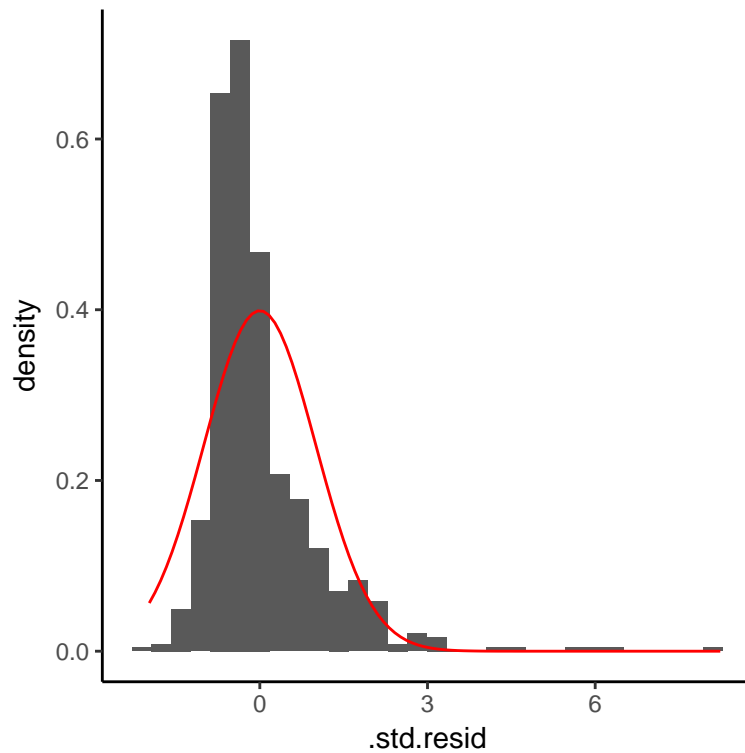
```
qqPlot(gbsg.lm2)
```



You could also plot a histogram of your residuals along with a normal distribution.

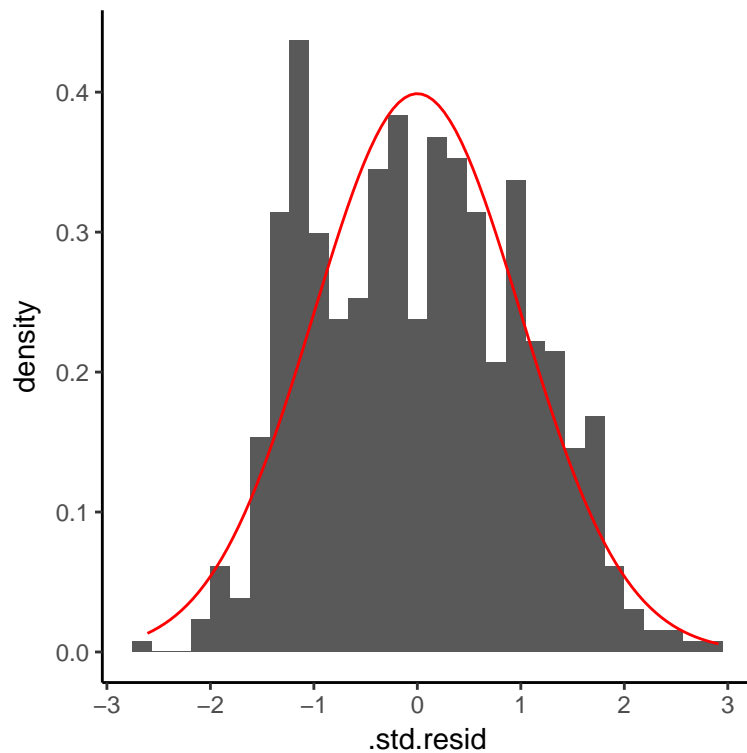
```
# take a look at the residuals
augment(gbsg.lm) %>%
  ggplot(aes(.std.resid)) + # plot studentized residuals on x-axis
  geom_histogram(aes(y = ..density..)) + # plot histogram as density, rather than frequency
  stat_function(fun = dnorm, color = 'red') # put a N(0,1) density over the top
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
augment(gbsg.lm2) %>%
  ggplot(aes(.std.resid)) + # plot studentized residuals on x-axis
  geom_histogram(aes(y = ..density..)) + # plot histogram as density, rather than frequency
  stat_function(fun = dnorm, color = 'red') # put a N(0,1) density over the top
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



No/little multicollinearity

Variance inflation factors give you a measure of how much the residuals are inflated when each predictor is added to the model (compared to the full model minus the predictor in question). This can give us an idea for what predictors might be correlated with each other. Anything greater than 2 should certainly give us pause.

```
# variance inflation factors
```

```
vif(gbsg.lm)
```

```
##      age      meno      size      grade      pgr      er      hormon
## 2.579321 2.513984 1.016654 1.050404 1.220625 1.327381 1.096945
```

```
vif(gbsg.lm2)
```

```
##      age      size      grade      meno      lpgr      ler      hormon
## 2.516079 1.013574 1.160316 2.527746 1.949148 1.948627 1.093311
```

No autocorrelation

Autocorrelation most often occurs when there is a correlation between observations at regular time intervals. The Durbin-Watson test will give us a measure of autocorrelation (the null hypothesis is that there is no autocorrelation).

```
# Durbin-Watson test for autocorrelation
```

```
durbinWatsonTest(gbsg.lm)
```

```
## lag Autocorrelation D-W Statistic p-value
## 1 -0.004115498 2.005419 0.984
## Alternative hypothesis: rho != 0
```

```
durbinWatsonTest(gbsg.lm2)
```

```
## lag Autocorrelation D-W Statistic p-value
## 1 -0.01422531 2.025106 0.802
## Alternative hypothesis: rho != 0
```

Homoscedasticity

Heteroscedasticity refers to a model where the variance about the predicted value of the outcome changes as the predicted value increases. We can check for this using a score test for non-constant error variance or graphically.

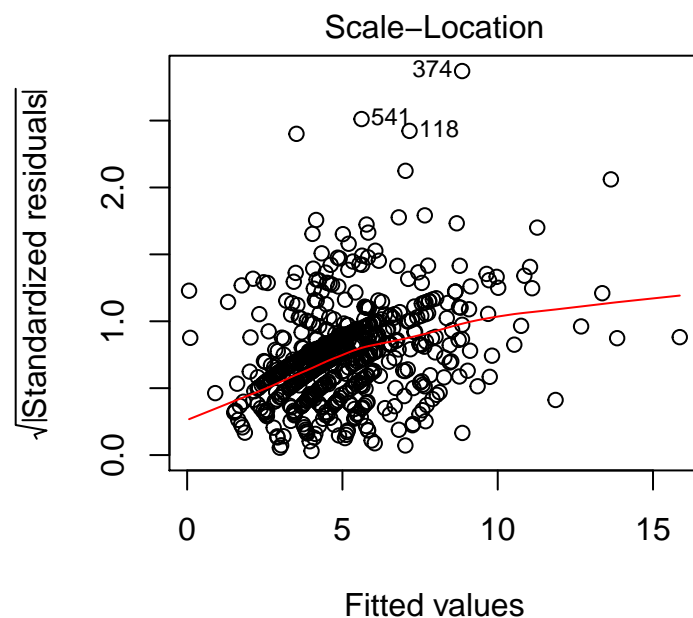
```
# statistical test for heteroscedasticity (null is that they are homoscedastic)
ncvTest(gbsg.lm)
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 165.5525 Df = 1 p = 6.928119e-38
```

```
ncvTest(gbsg.lm2)
```

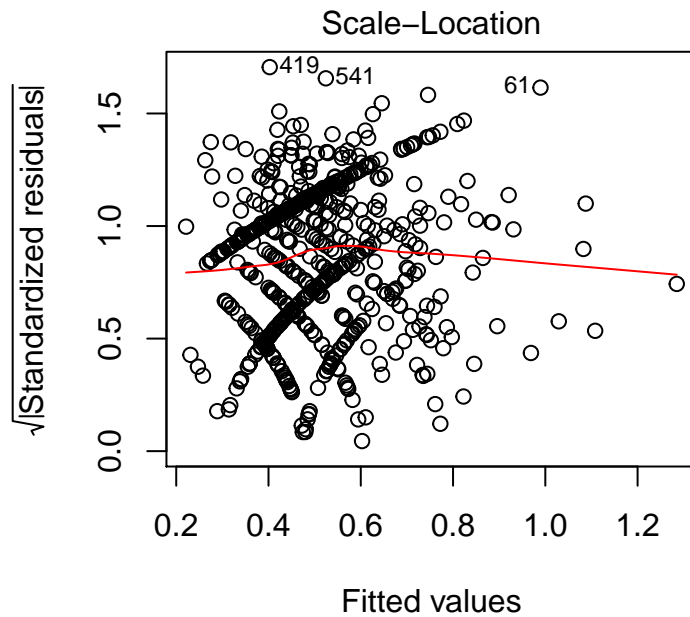
```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 7.317349 Df = 1 p = 0.006829208
```

```
# scale-location plot
plot(gbsg.lm, which = 3)
```



```
lm(nodes ~ age + meno + size + grade + pgr + er + hoi)
```

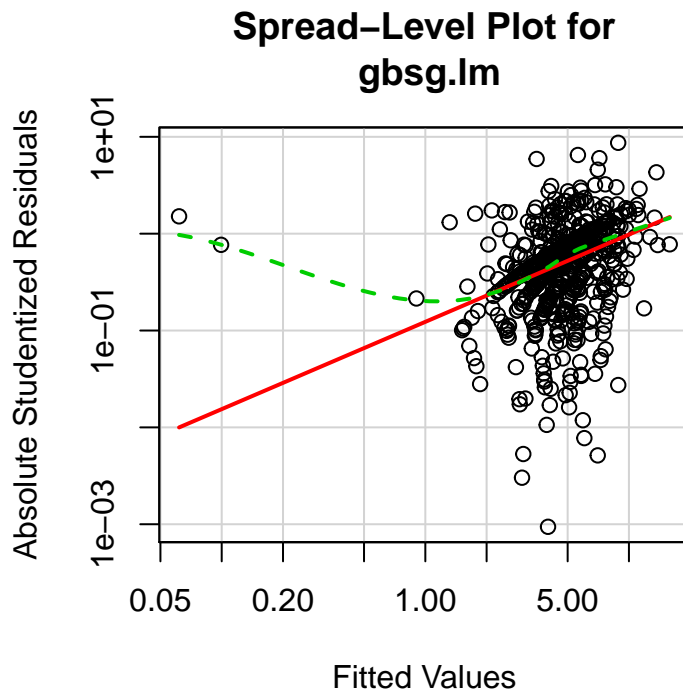
```
plot(gbsg.lm2, which = 3)
```

```
n(lnodes ~ age + size + grade + meno + lpgr + ler + hc
```

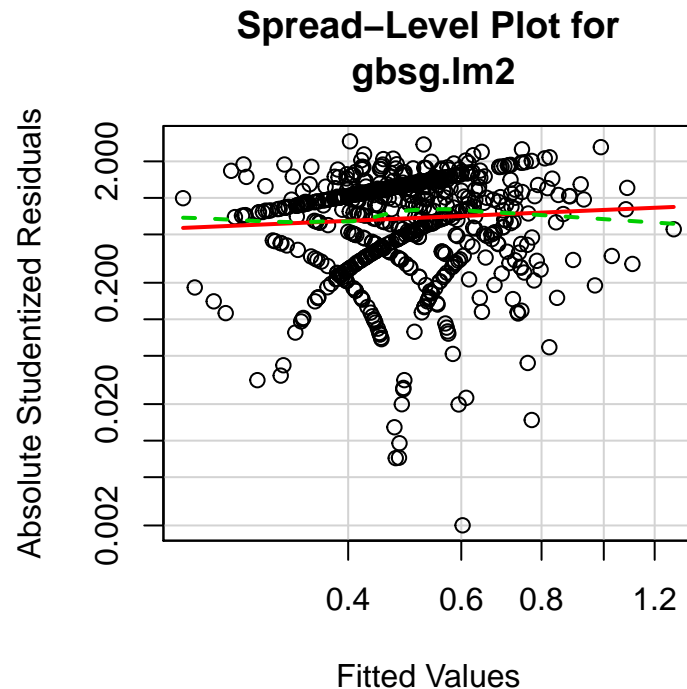
The `spreadLevelPlot` function gives us a slightly nicer picture of this. The solid red line is what we would hope to see under a homoscedastic model, and the green line is what is actually observed. These two plots (the one just above and the one just below) are representations of the exact same data. The plot below, however, is scaled differently to highlight departures from the expected.

```
# another scale-location plot
spreadLevelPlot(gbsg.lm)
```



```
##
## Suggested power transformation: 0.09857515
```

```
spreadLevelPlot(gbbsg.lm2)
```



```
##  
## Suggested power transformation: 0.774467
```

Outliers / Influential Points

The `outlierTest` function gives us a list of the most significant outliers in a model by row number. The maximum number of rows to return is controlled by the `n.max` function argument.

```
# returns the most significant outliers in the residuals  
outlierTest(gbbsg.lm)
```

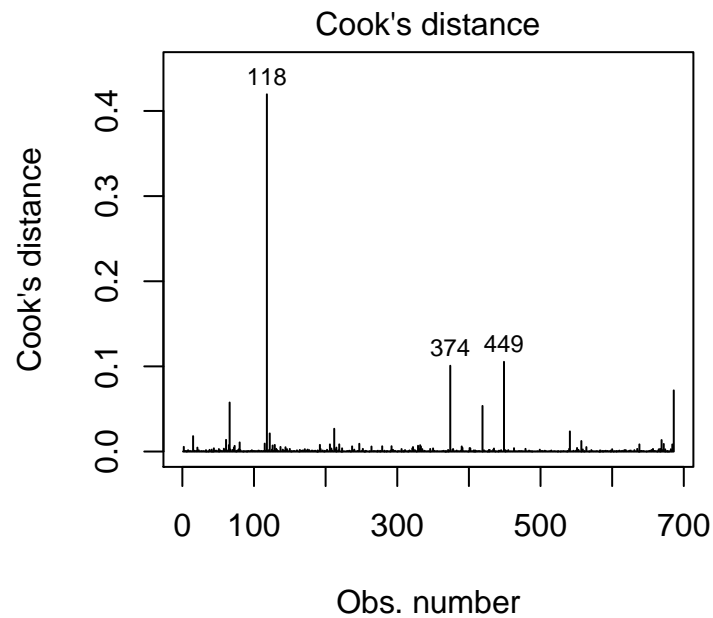
```
##      rstudent unadjusted p-value Bonferonni p  
## 374 8.677313      3.0022e-17  2.0595e-14  
## 541 6.496241      1.5957e-10  1.0946e-07  
## 118 6.021028      2.8416e-09  1.9493e-06  
## 419 5.907297      5.5069e-09  3.7777e-06  
## 66  4.579290      5.5528e-06  3.8092e-03  
## 449 4.299970      1.9594e-05  1.3442e-02
```

```
outlierTest(gbbsg.lm2)
```

```
##  
## No Studentized residuals with Bonferonni p < 0.05  
## Largest |rstudent|:  
##      rstudent unadjusted p-value Bonferonni p  
## 419 2.926259      0.0035456      NA
```

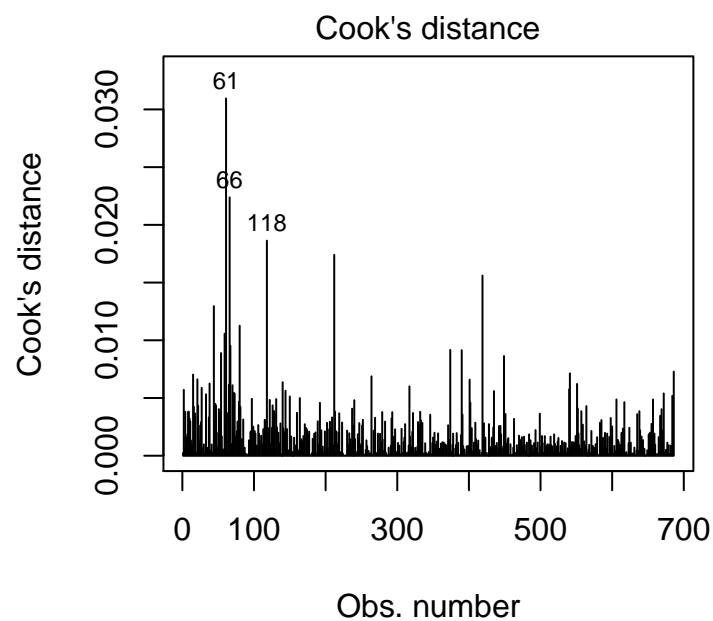
We can also graphically inspect the influence that each row of data (or each individual in the sample) has on the overall model. Cook's distance gives us a measure of how much each row of data influences the model. Ideally we would like them all to be close to the same, but we sometimes have some values that are overly influential (see this illustration of how one data point can influence a regression fit).

```
# Cook's Distance
plot(gbsg.lm, which = 4)
```



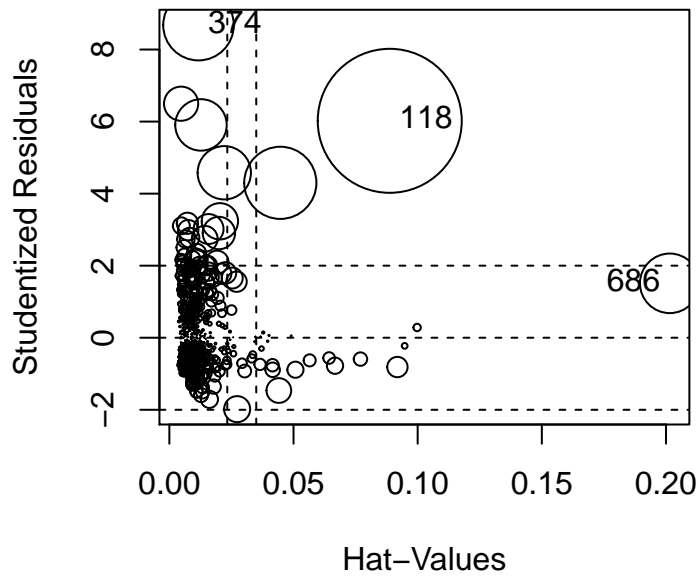
```
lm(nodes ~ age + meno + size + grade + pgr + er + ho
```

```
plot(gbsg.lm2, which = 4)
```



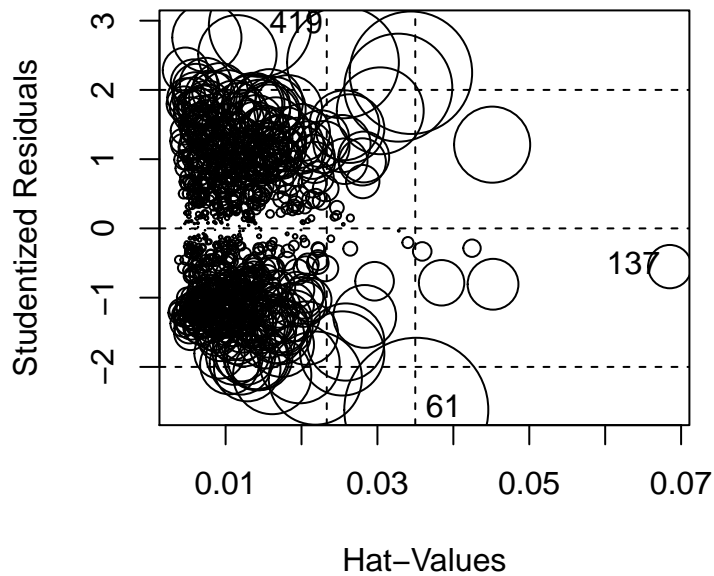
```
n(lnodes ~ age + size + grade + meno + lpgr + ler + hc
```

```
# Influence
influencePlot(gbsg.lm)
```



```
##      StudRes      Hat      CookD
## 118  6.021028 0.08875591 0.41956588
## 374  8.677313 0.01173853 0.10075428
## 686  1.510113 0.20151716 0.07180519
```

```
influencePlot(gbsg.lm2)
```



```
##      StudRes      Hat      CookD
## 61   -2.619349 0.03511812 0.030946805
## 137  -0.551649 0.06850720 0.002800518
## 419   2.926259 0.01452950 0.015607198
```

Influence of terms in the model

We also want to make our model as parsimonious as we can. Sometimes we will include a predictor variable because we believe that it is important, but usually we will include predictors only if they appear to be statistically important to our model. To get a quick look at the statistical significance of a predictor, we can

just look at the regression output.

```
tidy(gbsg.lm2)
```

```
##           term      estimate std.error statistic    p.value
## 1 (Intercept)  0.2188534351 0.158828460  1.3779233 1.686816e-01
## 2           age -0.0004743639 0.002313088 -0.2050782 8.375726e-01
## 3           size  0.0084830687 0.001039319  8.1621425 1.603046e-15
## 4           grade  0.0575594521 0.027277523  2.1101422 3.521185e-02
## 5 menopremenopausal -0.0262983051 0.047464585 -0.5540616 5.797194e-01
## 6           lpgr -0.0253118219 0.018632700 -1.3584624 1.747688e-01
## 7           ler   0.0026068067 0.019428389  0.1341751 8.933039e-01
## 8 hormonno tamoxifen -0.0304764028 0.032153444 -0.9478426 3.435473e-01
```

Added Variable Plots provide a graphical approach to identifying how much added predictive value a specific variable will give you. In order to generate an Added Variable Plot, we need to perform two additional regressions.

```
tmp <- data_frame(x1 = rnorm(100),
                  x2 = rnorm(100),
                  x3 = rnorm(100),
                  x4 = rnorm(100),
                  y = x1 + x2 + x3 + rnorm(100))

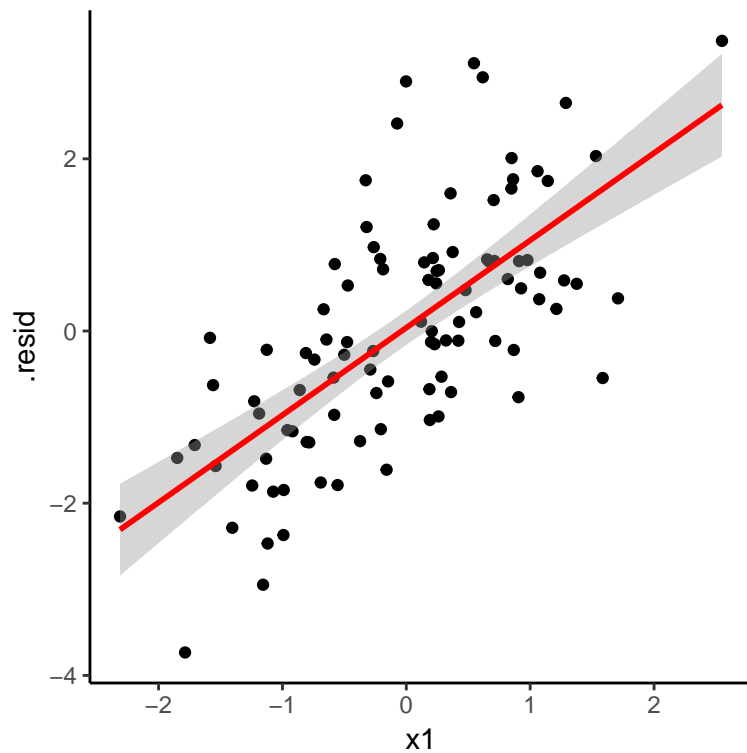
# full model
full_model <- lm(y ~ x1 + x2 + x3 + x4, data = tmp)
tidy(full_model)
```

```
##           term      estimate std.error statistic    p.value
## 1 (Intercept) -0.029346033 0.1031488 -0.28450192 7.766450e-01
## 2           x1  1.025008651 0.1115059  9.19241269 8.798286e-15
## 3           x2  0.966155517 0.1146042  8.43036674 3.701227e-13
## 4           x3  1.192778608 0.1097284 10.87028332 2.315127e-18
## 5           x4  0.009432959 0.1137804  0.08290493 9.341016e-01
```

```
# model without x1
no_x1 <- update(full_model, . ~ . - x1) %>%
  augment()

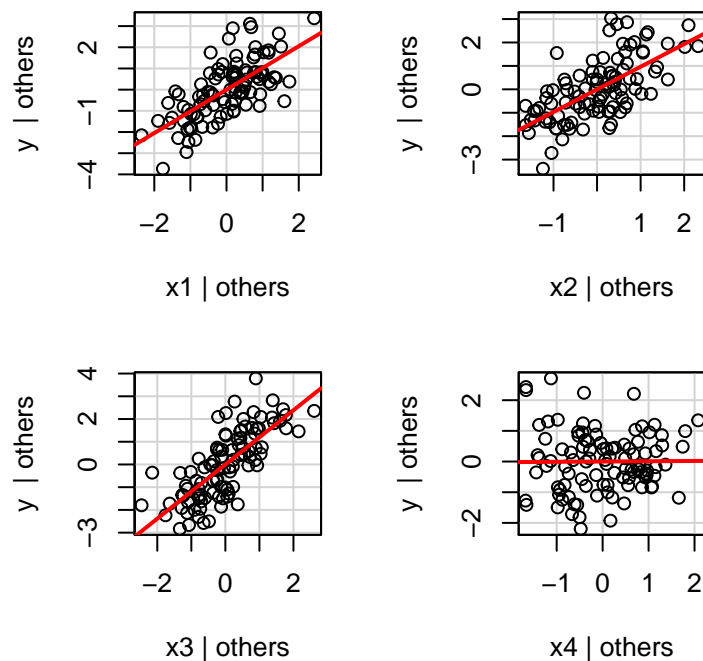
no_x1$x1 <- tmp$x1 # put back into the data_frame, but not into the model

# now regress x1 onto y conditioning on the other variables
ggplot(data = no_x1, aes(x = x1, y = .resid)) +
  geom_point() +
  geom_smooth(method = 'lm', color = 'red')
```



```
avPlots(full_model)
```

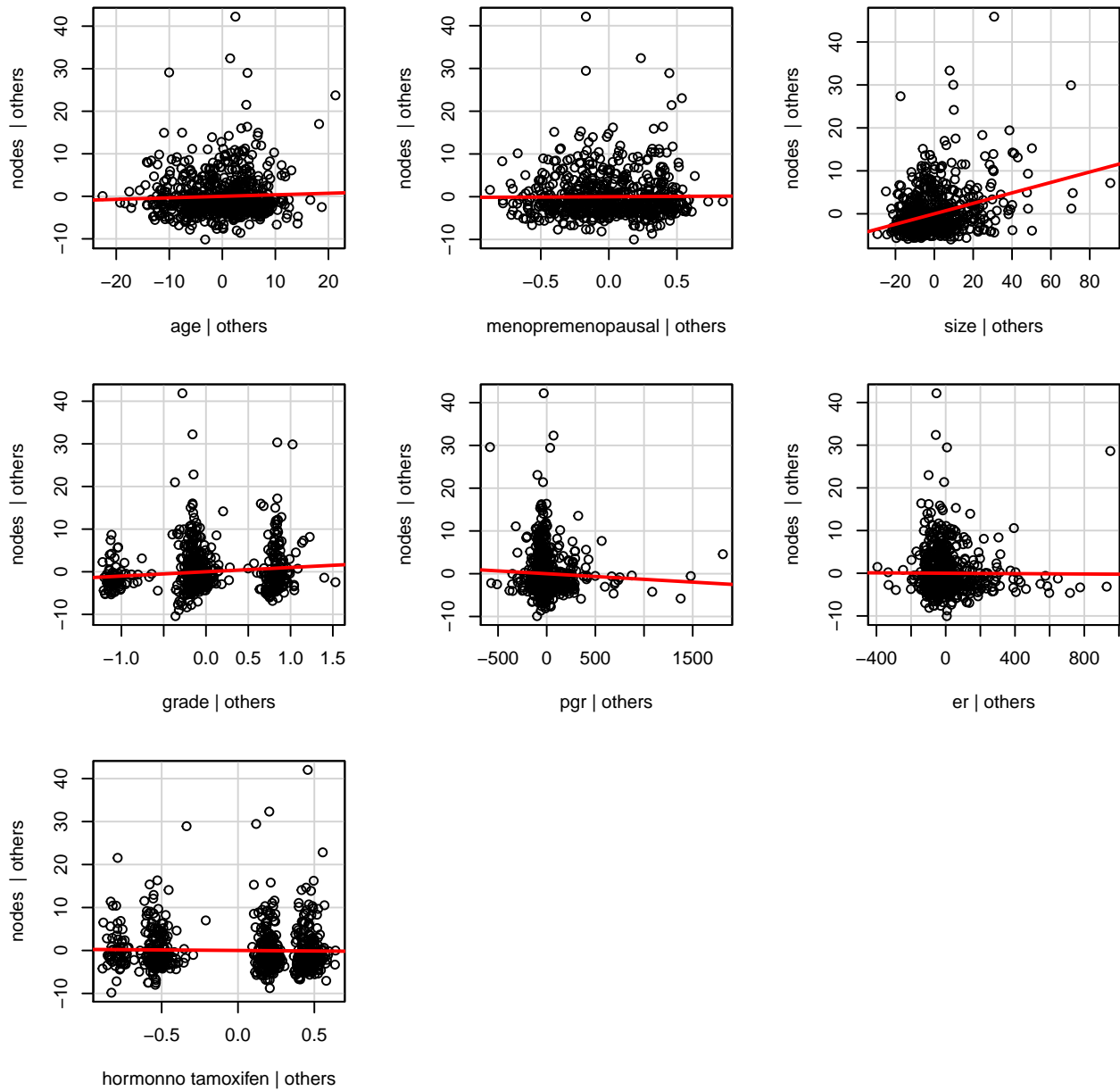
Added-Variable Plots



Now lets take a look at all of the predictors in our model:

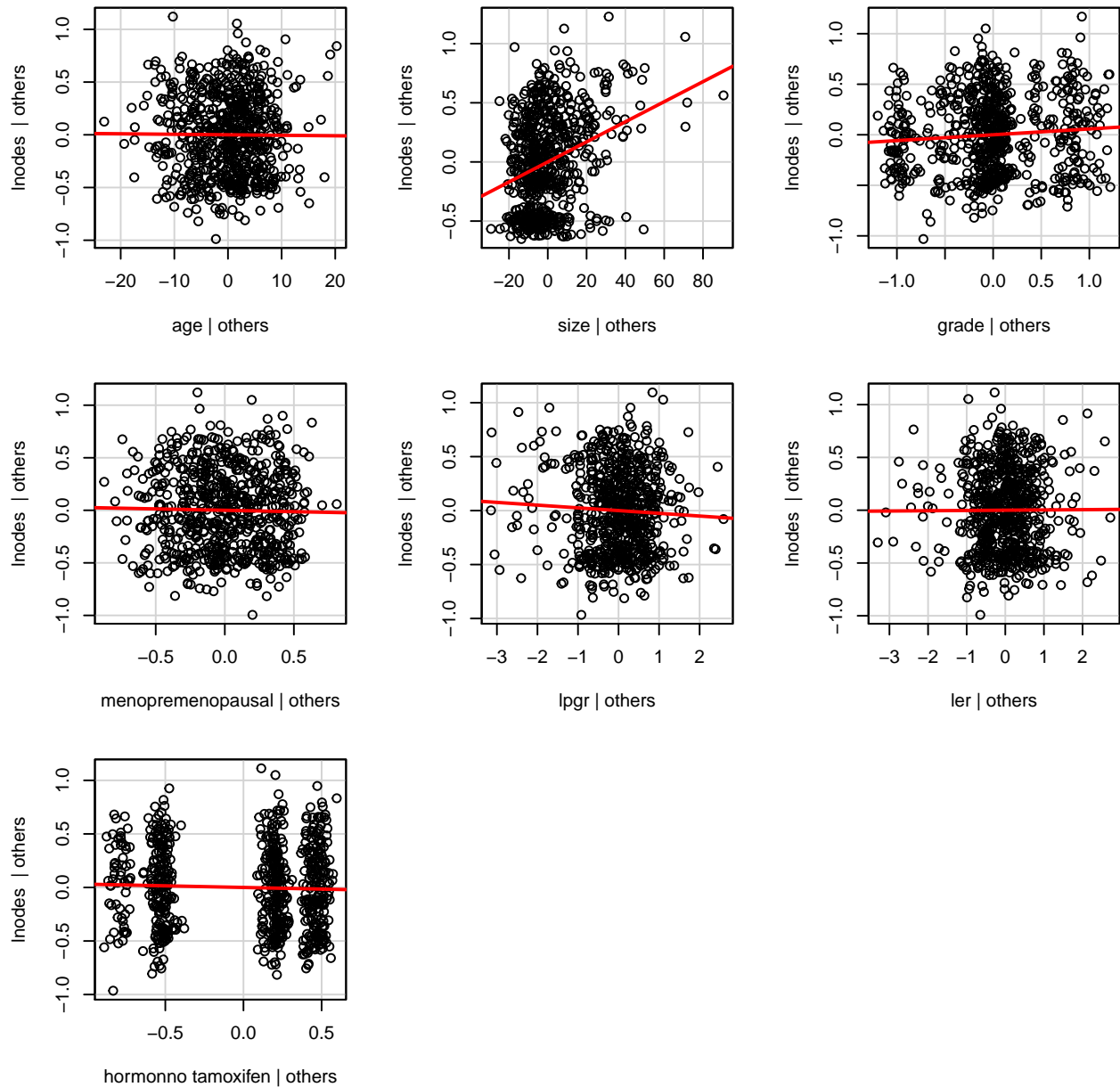
```
# added variable plots
avPlots(gbgs.lm)
```

Added-Variable Plots



```
avPlots(gbsg.lm2)
```

Added-Variable Plots



Homework

Choose the best model for 'nodes' in the GBSG dataset (or log nodes), checking all model assumptions. A good place to start is with the full model we began with in Discussion 3. Be sure to take a look at the rubric associated with this for specific things I will be looking for. Submit a link to a gist with your code. All figures and tables you wish to include should be generated by your code (i.e. you don't need to include any images – I will run your code and generate any figures you wish to include). All comments and rationale for your final model should be included in the code.