

Exercise 14 - Model diagnostics

Zoltan Kekecs

1 April 2019

Contents

1	Abstract	2
2	Data management and descriptive statistics	2
2.1	Loading packages	2
2.2	Custom functions	2
2.3	Load data about housing prices in King County, USA	3
2.4	Check the dataset	3
3	Model diagnostics	3
3.1	Build the model	3
4	Dealing with outliers	4
4.1	Identifying extreme cases	4
4.2	Identifying extreme cases with high leverage	4
5	Assumptions of linear regression	7
5.1	Normality	7
5.2	Linearity	14
5.3	Homoscedasticity	16
5.4	No multicollinearity	22

1 Abstract

This exercise will show you how to check whether the assumptions of linear regression hold true for your model, what is the consequence of the violation of the assumptions, and what to do in case of an assumption being violated.

The latest version of this document and the code the document refers to can be found in the GitHub repository of the class at: https://github.com/kekecsz/PSYP13_Data_analysis_class-2018

2 Data management and descriptive statistics

2.1 Loading packages

You will need the following packages for this exercise.

```
library(psych) # for describe
library(car)   # for residualPlots, vif, pairs.panels, ncvTest
library(lmtest) # bptest
library(sandwich) # for coefTest vcovHC estimator
library(boot)  # for bootstrapping
```

2.2 Custom functions

We will use these custom functions to get bootstrapped confidence intervals.

```
# function to obtain regression coefficients source:
# https://www.statmethods.net/advstats/bootstrapping.html
bs_to_boot <- function(model, data, indices) {
  d <- data[indices, ] # allows boot to select sample
  fit <- lm(formula(model), data = d)
  return(coef(fit))
}

# function to obtain adjusted R^2 source:
# https://www.statmethods.net/advstats/bootstrapping.html
# (partially modified)
adjR2_to_boot <- function(model, data, indices) {
  d <- data[indices, ] # allows boot to select sample
  fit <- lm(formula(model), data = d)
  return(summary(fit)$adj.r.squared)
}

# Computing the bootstrap BCa (bias-corrected and
# accelerated) bootstrap confidence intervals by Efron
# (1987) This is useful if there is bias or skew in the
# residuals.

confint.boot <- function(model, data = NULL, R = 1000) {
  if (is.null(data)) {
    data = eval(parse(text = as.character(model$call[3])))
  }
  boot.ci_output_table = as.data.frame(matrix(NA, nrow = length(coef(model)),
    ncol = 2))
  row.names(boot.ci_output_table) = names(coef(model))
```

```

names(boot.ci_output_table) = c("boot 2.5 %", "boot 97.5 %")
results.boot = results <- boot(data = data, statistic = bs_to_boot,
  R = 1000, model = model)

for (i in 1:length(coef(model))) {
  boot.ci_output_table[i, ] = unlist(unlist(boot.ci(results.boot,
    type = "bca", index = i))[c("bca4", "bca5")])
}

return(boot.ci_output_table)
}

```

2.3 Load data about housing prices in King County, USA

In this exercise we will use the same housing price dataset as in exercises 12 and 13.

The dataset is from Kaggle containing data about housing prices and variables that may be used to predict housing prices. The data is about apartments sold in King County, USA (Seattle and surrounding area).

We only use a portion of the full dataset now containing information about $N = 200$ accommodations.

You can load the data with the following code

```

# data from
# github/kekecsz/PSYP13_Data_analysis_class-2018/master/data_house_small_sub.csv.
data_house = read.csv("https://bit.ly/2DpwK0r")

```

2.4 Check the dataset

You should always check the dataset for coding errors or data that does not make sense, by eyeballing the data through the data view tool, checking descriptive statistics and through data visualization.

3 Model diagnostics

Whenever we arrive at a model that we use for statistical inference, we need to make sure that the assumptions of the linear regression hold true for the model we are working with.

Thus, you have to do model diagnostics for all important models in your analyses. Usually we would do model diagnostics on the final model we arrive at in the end of model selection.

In some cases, if you do result-based post hoc model selection, it might be important to run model diagnostics on the interim steps as well, the models based on which you reached your final model composition.

If you make any adjustments to the data, or the model based on the model diagnostics, remember that you have to re-run model diagnostics on the new model/new dataset.

3.1 Build the model

Here we first build a model to predict the price of the apartment by using only `sqft_living` and `grade` as predictors. We are going to run model diagnostics on this model.

```

mod_house2 <- lm(price ~ sqft_living + grade, data = data_house)

```

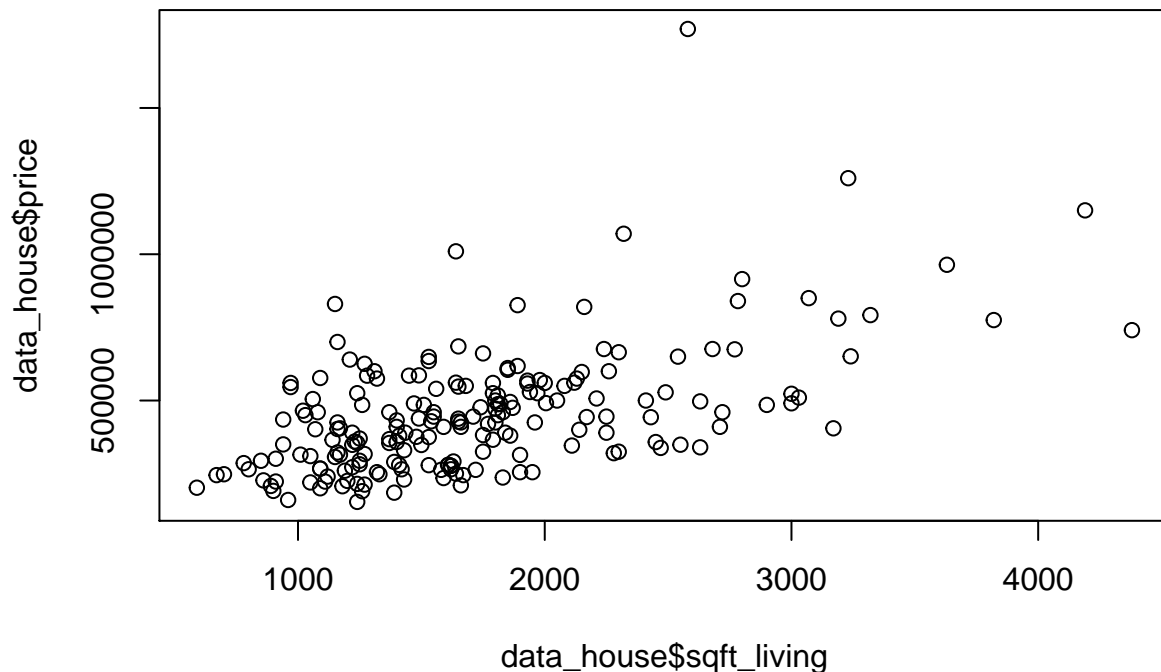
4 Dealing with outliers

4.1 Identifying extreme cases

Cases with extreme values can be identified by visualizing the data and the relationship of the outcome and the predictors one by one.

For example here we visualize price and squarefootage data on a scatterplot.

```
plot(data_house$price ~ data_house$sqft_living)
```



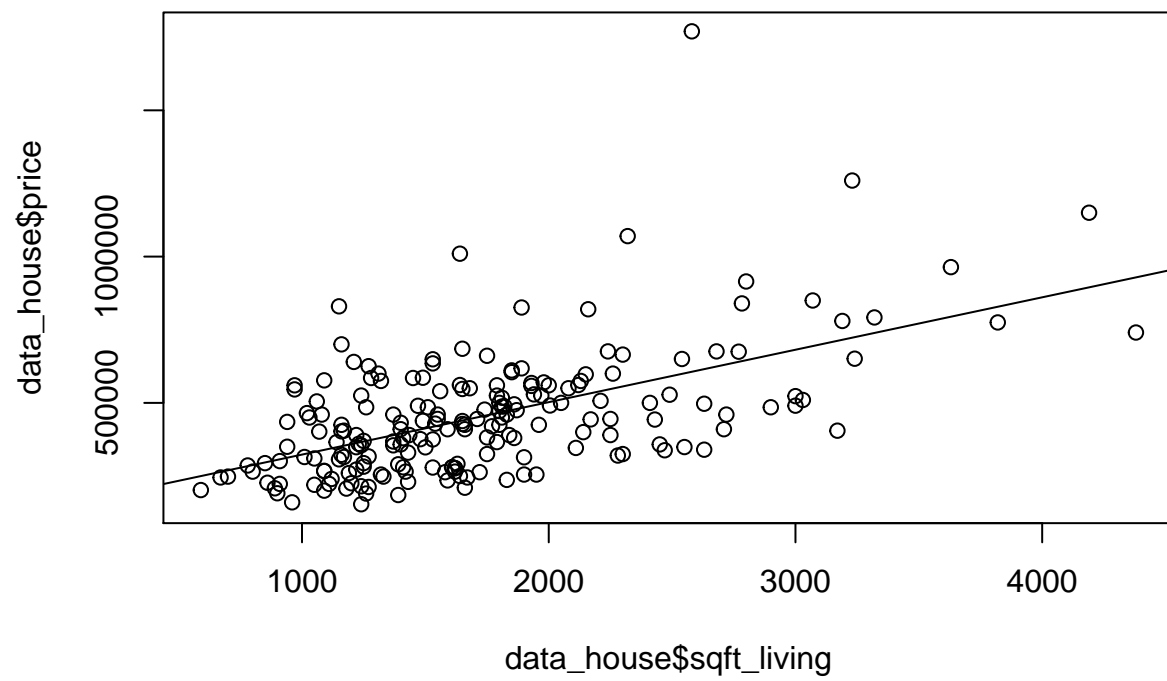
Notice that most apartment sales had a final sales price below 1 million USD. However, there were a few cases where the price was higher than that. These might be considered extreme cases, especially the one with the sales price of 1.7 million USD.

However, it might not be a huge problem if we have a few of these cases in the dataset if we have a lot of data to counteract their effects, especially if they don't have high leverage.

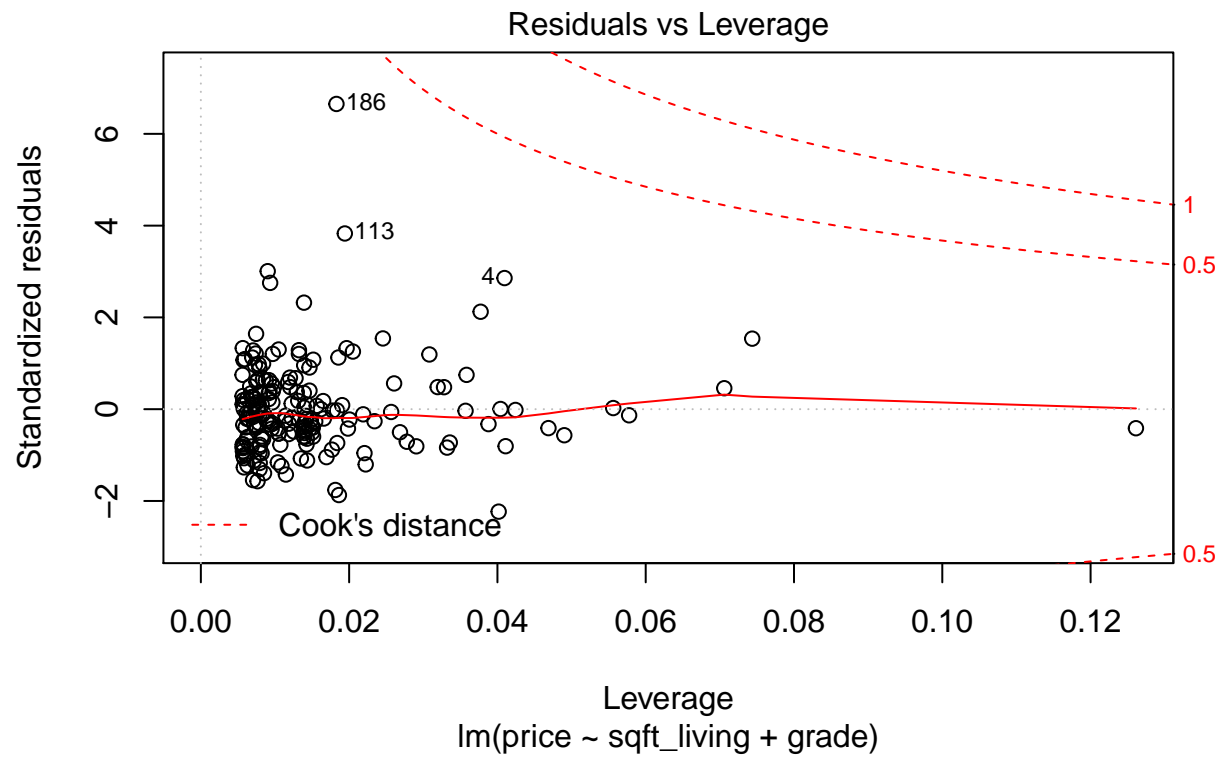
4.2 Identifying extreme cases with high leverage

More problematic are cases that not only have very different values from the other data, but they have a high influence on the regression line as well. Highly influential cases can be identified by looking at the scatterplot with the regression line, by plotting the residual - leverage plot, and by using Cook's distance.

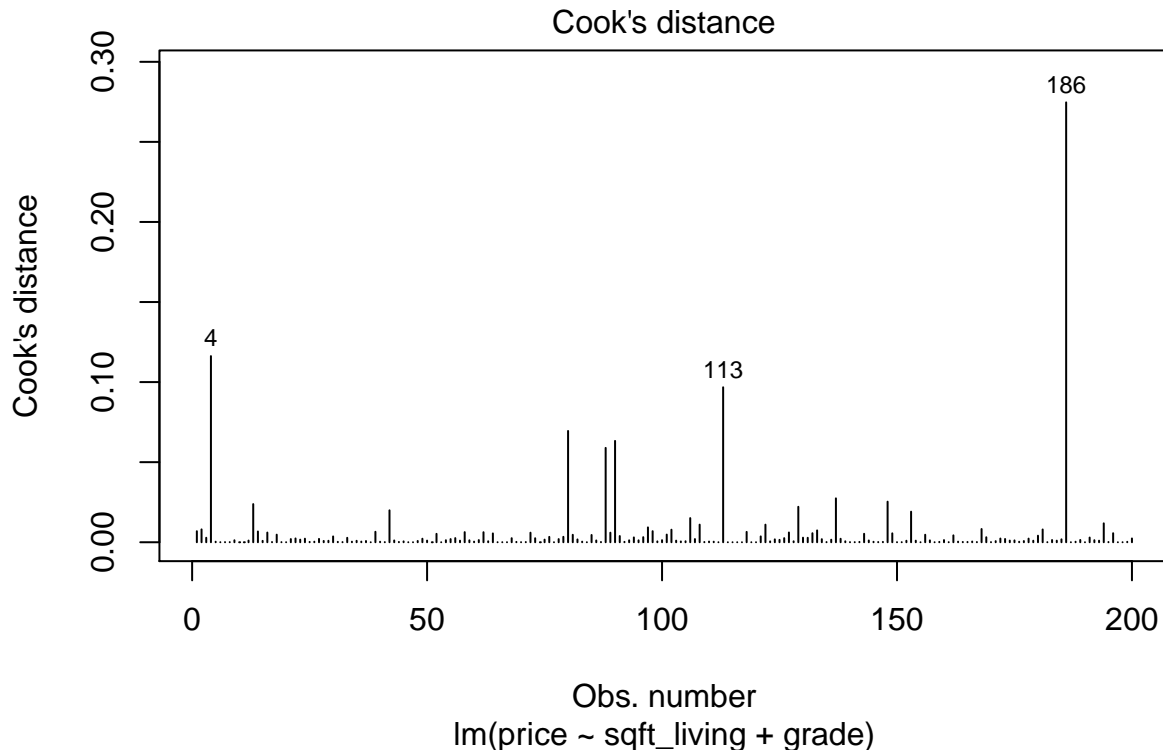
```
plot(data_house$price ~ data_house$sqft_living)
abline(lm(price ~ sqft_living, data = data_house))
```



```
plot(mod_house2, which = 5)
```



```
plot(mod_house2, which = 4)
```



Cases that are close to the middle of the regression line on the scatterplot have less leverage while cases at the ends have more. Cases that have high residual error and high leverage are influential cases. Cook's distance is a number that quantifies this, taking into account the combination of "extremeness" and leverage.

It is not well established what should be considered a problematic case, but there are some rules of thumb. Some say that cases with Cook's distance > 1 are influential, while some use a threshold of Cook's distance $> 4/N$. (In our case this is $4/200 = 0.02$).

We have no cases that would have a Cook's distance higher than 1, but we have several where it is higher than 0.02. So we might have a problem with outliers according to one of the criteria.

Lets test whether the assumptions of multiple regression hold true, and determine if we need to do anything with these cases.

5 Assumptions of linear regression

- **Normality:** The residuals of the model must be normally distributed
- **Linearity:** Ther relationship between the outcome variable and the predictor(s) must be linear
- **Homoscedasticity:** The variance of the residuals are similar at all values of the predictor(s)
- **No Multicollinearity:** None of the predictors can be linearly determined by the other predictor(s)

5.1 Normality

The residuals of the model must be normally distributed.

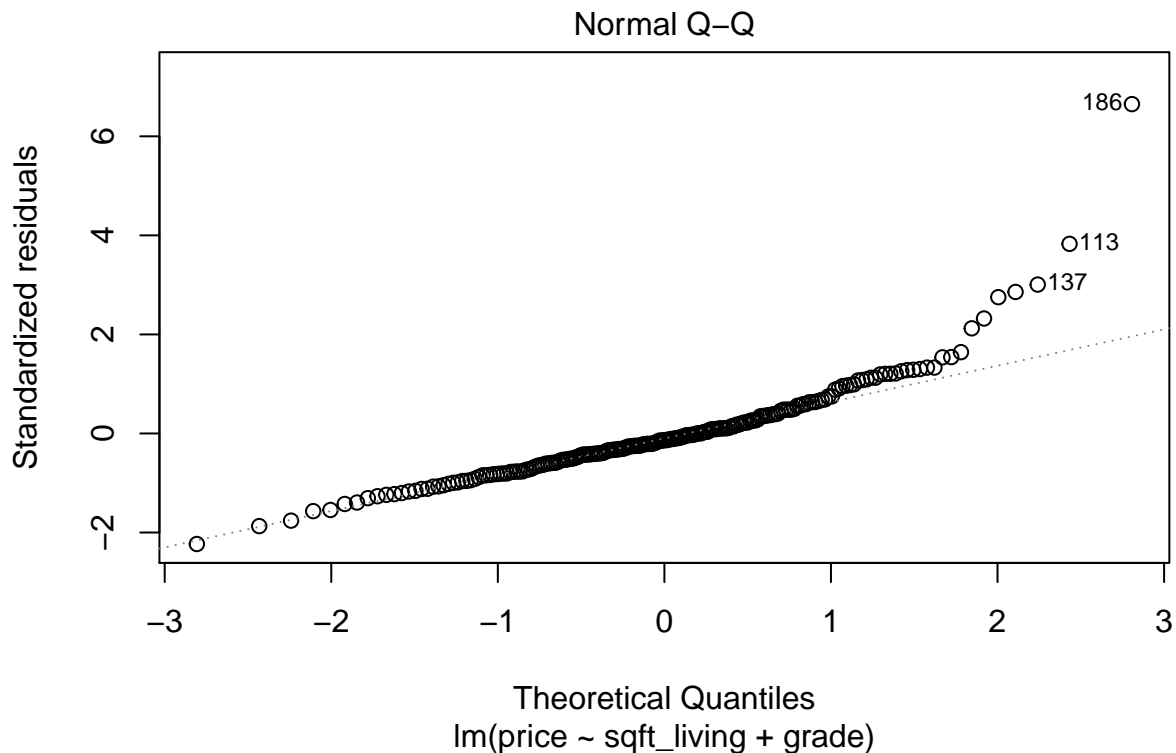
Note that here we talk about the distribution of the prediction error (residuals) of the model, and not the distribution of the individual predictors or the outcome.

We can check this assumption by plotting the QQ plot and checking if the cases are aligned with the theoretical diagonal. If the cases divert from the diagonal (the dashed line in the plot), the assumption of normality might be violated.

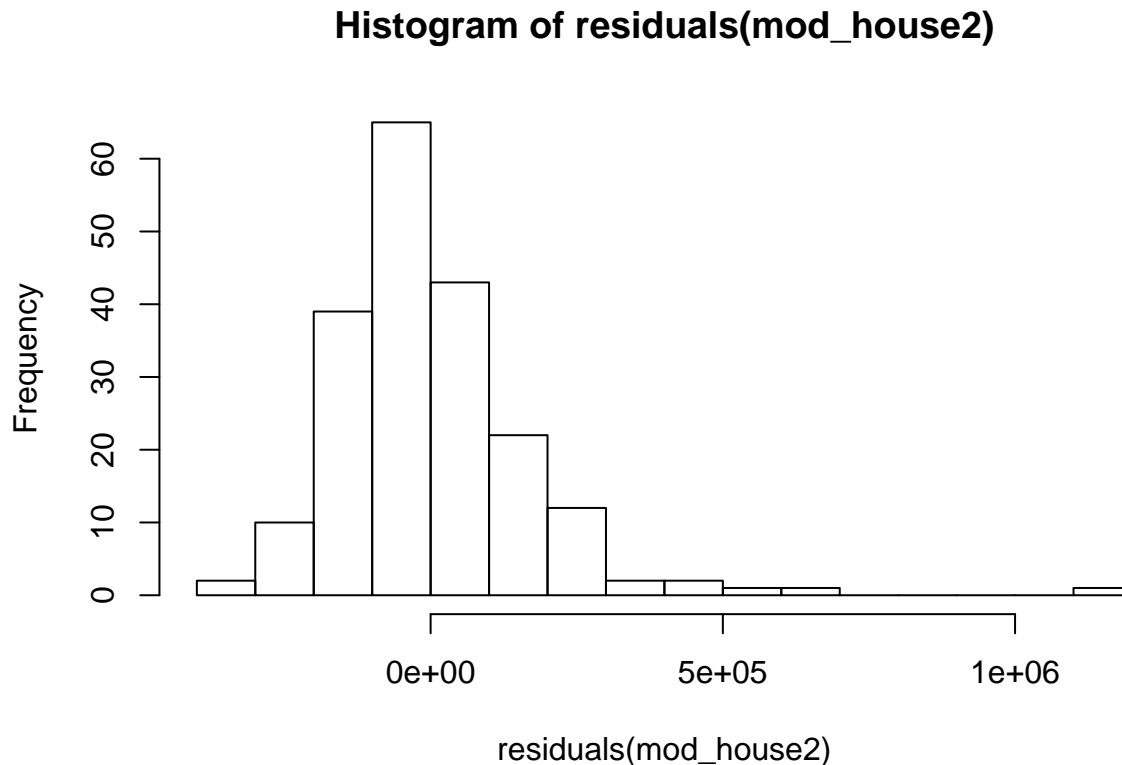
We should also look at the histogram of the residuals. We should see a roughly bell shaped (normal) distribution if the assumption of normality is met.

We can also ask for the skew and kurtosis statistics using the `describe()` function. Skew and kurtosis > 1 can indicate the violation of the assumption of normality.

```
# QQ plot
plot(mod_house2, which = 2)
```



```
# histogram
hist(residuals(mod_house2), breaks = 20)
```

```
# skew and kurtosis
describe(residuals(mod_house2))
```

```
##      vars   n mean      sd   median trimmed      mad      min      max
## X1      1 200    0 169214.6 -23118.96 -14686.4 125887.3 -371917.1 1120748
##      range skew kurtosis      se
## X1 1492665 2.05      9.73 11965.28
```

Notice that the plots and the statistics indicate a slight deviation from the assumption of normality mainly driven by a few irregular cases.

5.1.1 What happens if the assumption of normality is violated?

The estimates and confidence intervals might be less accurate if the assumption of normality is violated. The effect of this depends on the sample size as well. For large sample sizes ($N > 500$) the effect was very minor, while for smaller samples (N round 100) the effect is greater. For example in the study of Lumley, Diehr, Emerson and Chen (2002), the effect of extreme violation of normality (skewness = 8.8, kurtosis = 131) was that the 95% confidence intervals contained the true population mean in only 93.6% of cases instead 95% when N was 500, and 91.3% of the cases when N was 65.

The bottom line is that if the assumption of normality is violated, the confidence intervals and the p values are less reliable, but they are still informative.

Reference:

Lumley, T., Diehr, P., Emerson, S., & Chen, L. (2002). The importance of the normality assumption in large public health data sets. *Annual review of public health*, 23(1), 151-169.

5.1.2 What to do if the assumption of normality is violated?

1. you can interpret the results more cautiously, such as using 99% CI instead of 95% CI, or $p < 0.01$ as a threshold for significance
2. if the deviation from normality is driven by a few cases, you might want to try to exclude the outliers. If you are doing formal hypothesis testing, exclusion of variables should not be based on the p-value. Rules for exclusion can be pre-registered, or a sensitivity analysis can be presented, meaning that you conduct the same analysis twice on data including and excluding the problematic cases, and compare the results of the two analysis to see the influence of exclusion of outliers.
3. you can try transforming the predictors and/or the outcome variable to achieve more normally distributed residuals. But if you do this, be mindful when you interpret the unstandardized coefficients that they will be in the transformed units. The same goes for the error terms, so RSS of a model on transformed variables cannot be compared with that of one with untransformed variables.
4. you can use bootstrapping to get a robust estimate of the confidence intervals.

In our case, since non-normality is the result of a few extreme cases, we can build a model with these cases excluded to see if this fixes the problem. We exclude cases 186 and 113 here because they showed up as extreme cases according to the Cook's distance, and they also contributed to the deviation from normality according to the QQ plot.

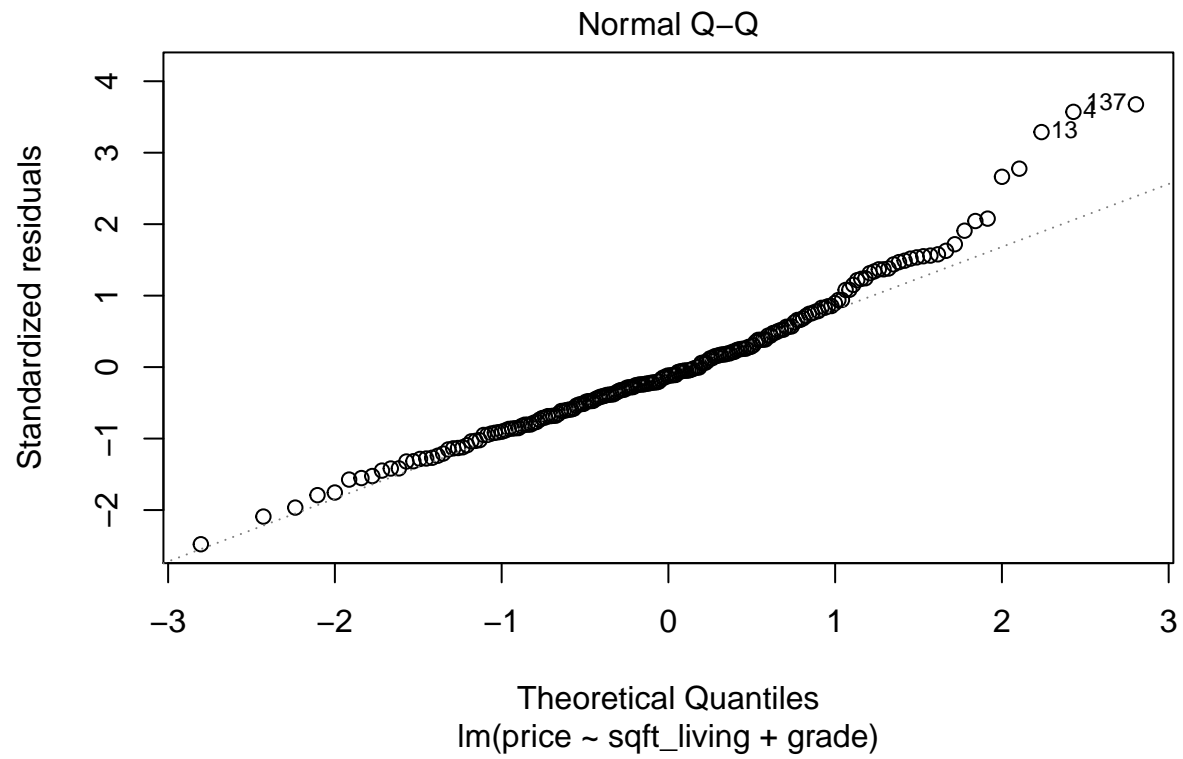
Here we refit the model without these cases two cases and recheck the assumption of normality. Notice that the residuals look much more normal in the model without the outliers.

```
data_house_nooutliers = data_house[-c(186, 113), ]
mod_house3 = lm(price ~ sqft_living + grade, data = data_house_nooutliers)

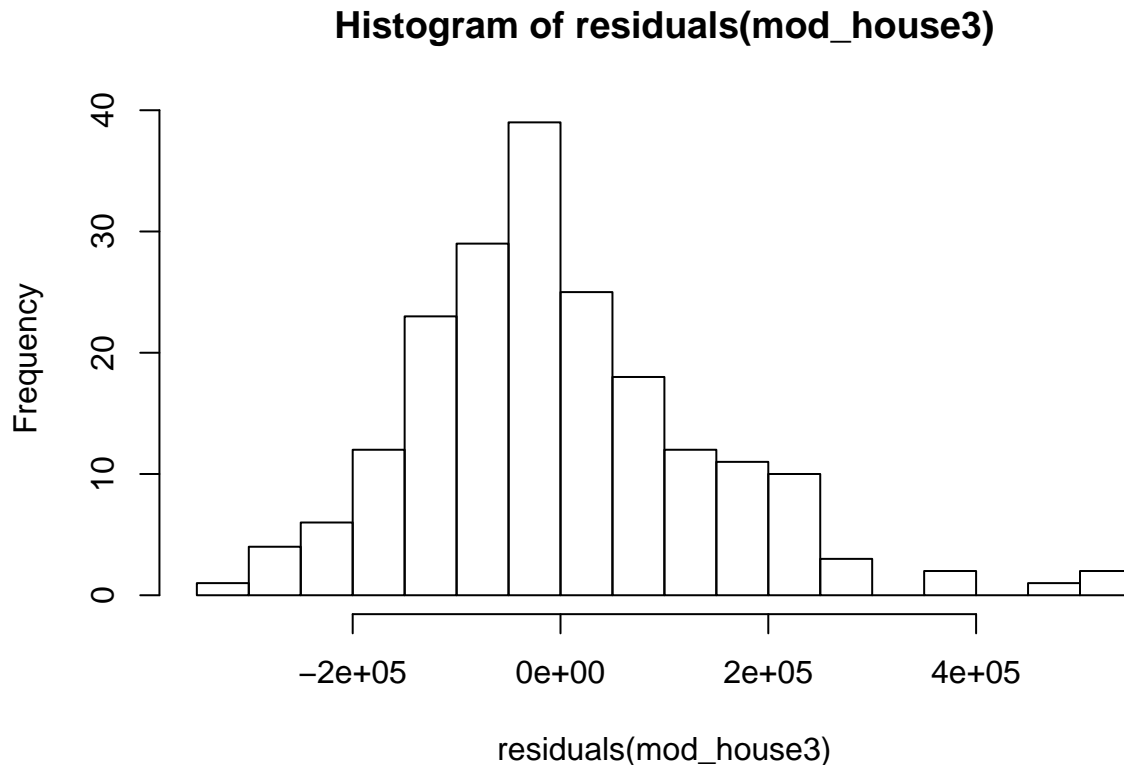
# recheck the assumption of normality of residuals
describe(residuals(mod_house3))
```

##	vars	n	mean	sd	median	trimmed	mad	min	max	
##	X1	1	198	0	142511.8	-16998.56	-8569.97	121394.9	-347654.3	524234.8
##			range	skew	kurtosis	se				
##	X1	871889.1	0.79	1.35	10127.87					

```
plot(mod_house3, which = 2)
```



```
hist(residuals(mod_house3), breaks = 20)
```



When we compare the two models in a sensitivity analysis, the exclusion of the outlier did not change the statistical inference, the previously significant predictors are still significant, and the overall model F test is significant in both cases. The adjusted R^2 improved considerably because our regression line now fits the remaining data much better. However, model fit should be assessed on new data or a test set as well to get a more accurate picture of true prediction efficiency, because similar outliers might be present in new data as well, and our model will not be very good at predicting these cases.

```
# comparing the models on data with and without the outliers
summary(mod_house2)
```

```
##
## Call:
## lm(formula = price ~ sqft_living + grade, data = data_house)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -371917 -100605  -23119   66886 1120748
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -174389.86   95255.17  -1.831  0.068646 .
## sqft_living    119.17     24.76    4.813 2.96e-06 ***
## grade         57352.79   16052.79    3.573 0.000444 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 170100 on 197 degrees of freedom
```

```
## Multiple R-squared:  0.358, Adjusted R-squared:  0.3515
## F-statistic: 54.94 on 2 and 197 DF,  p-value: < 2.2e-16
```

```
summary(mod_house3)
```

```
##
## Call:
## lm(formula = price ~ sqft_living + grade, data = data_house_nooutliers)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -347654  -95907  -16999   73466  524235
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -164126.22   81165.63  -2.022   0.0445 *
## sqft_living    108.95     20.92    5.208 4.82e-07 ***
## grade        57141.26   13653.24    4.185 4.31e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 143200 on 195 degrees of freedom
## Multiple R-squared:  0.413, Adjusted R-squared:  0.407
## F-statistic: 68.59 on 2 and 195 DF,  p-value: < 2.2e-16
```

We are going to work with the new model without the outliers in the upcoming tests of the assumptions.

5.1.3 Bootstrapping

In case you needed a robust estimate of the confidence intervals, you can use bootstrapping to get them. Bootstrapping means that we randomly sample from our own observations and rebuild the model on this random sample. We repeat this many times, and draw conclusions about the confidence limits based on this multitude of results.

You will need to run the custom functions in the top of the script for these to work.

Compare the regular and bootstrapped confidence intervals of the model coefficients.

```
# regular confidence intervals for the model coefficients
confint(mod_house3)
```

```
##              2.5 %      97.5 %
## (Intercept) -324201.40325 -4051.0325
## sqft_living    67.69491   150.2028
## grade        30214.27831 84068.2348
```

```
# bootstrapped confidence intervals for the model
# coefficients
confint.boot(mod_house3)
```

```
##              boot 2.5 % boot 97.5 %
## (Intercept) -334177.00307 -11545.269
## sqft_living    69.92251    158.985
## grade        32075.58477 82958.086
```

```
# regular adjusted R squared
summary(mod_house3)$adj.r.squared
```

```
## [1] 0.406965
```

```

# bootstrapping with 1000 replications
results.boot <- boot(data = data_house, statistic = adjR2_to_boot,
  R = 1000, model = mod_house3)

# get 95% confidence intervals for the adjusted R2
boot.ci(results.boot, type = "bca", index = 1)

## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 1000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = results.boot, type = "bca", index = 1)
##
## Intervals :
## Level      BCa
## 95%      ( 0.2228,  0.4784 )
## Calculations and Intervals on Original Scale

```

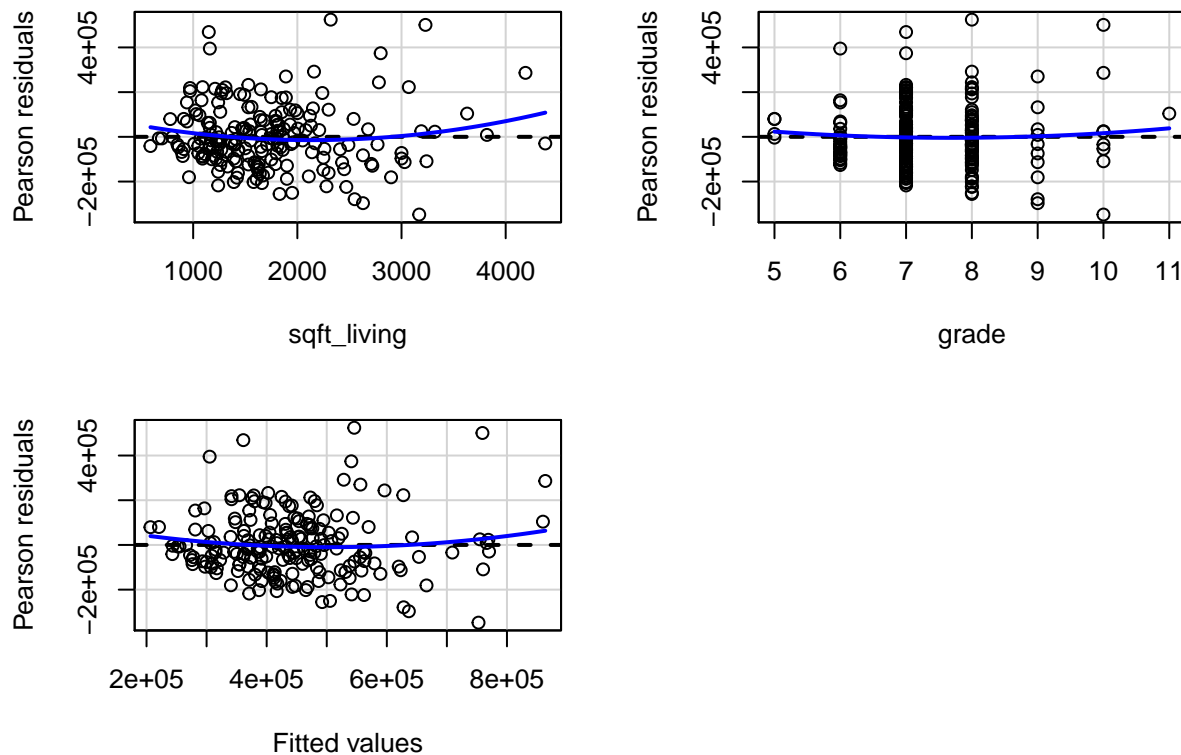
5.2 Linearity

The relationship between the outcome variable and the predictor(s) must be linear.

You can use the `residualPlots()` function for each predictor from the `car` package to explore linearity. It will return the scatterplot with a spline roughly indicating the relationship of the predictor and the outcome, and you will also see the residual-predicted value plot. In all of these plots you should see roughly flat lines if the assumption of linearity holds true.

Also, this function returns the results of a non-linearity test for each outcome-predictor pair. If the test is significant ($p < 0.05$), there might be a violation of the linearity assumption.

```
residualPlots(mod_house3)
```



```
##           Test stat Pr(>|Test stat|)
## sqft_living    1.6218      0.1065
## grade          0.6485      0.5174
## Tukey test     1.2576      0.2085
```

In our case, even though there is minor curvature visible on the plots, the tests are all non significant, so the linearity assumption seems to hold true for our model.

5.2.1 What happens if the assumption of linearity is violated?

If there is nonlinear relationships among the outcome and the predictors, the predictions of the model may be off, resulting in lower prediction accuracy. Also, the model coefficients will be also unreliable when used for prediction, and even though the stnadardized coefficients and the t-test p-values for some predictors would indicate no significant added predictive value for a predictor, in reality the predictor might still be related to the outcome, just in a non-linear way.

5.2.2 What to do if the assumption of linearity if violated?

If the assumption of linearity is violated, you can try to account for the non-linearity in the relationship by making your model more flexible.

1. you can include a higher order term of the predictor that seems to have a non-linear relationship with the outcome. (See exercise 12 on how to do this.) Usually, including the second and third order term should be enough to address curved relationship. You should avoid adding too much flexibility to reduce overfitting.
2. If higher order terms do not capture the relationship well, you might have to experiment with non-linear regression. You can learn more about nonlinear regression for example from this book: James, G.,

Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning. New York: springer. It is freely available at: <http://www-bcf.usc.edu/~gareth/ISL/>

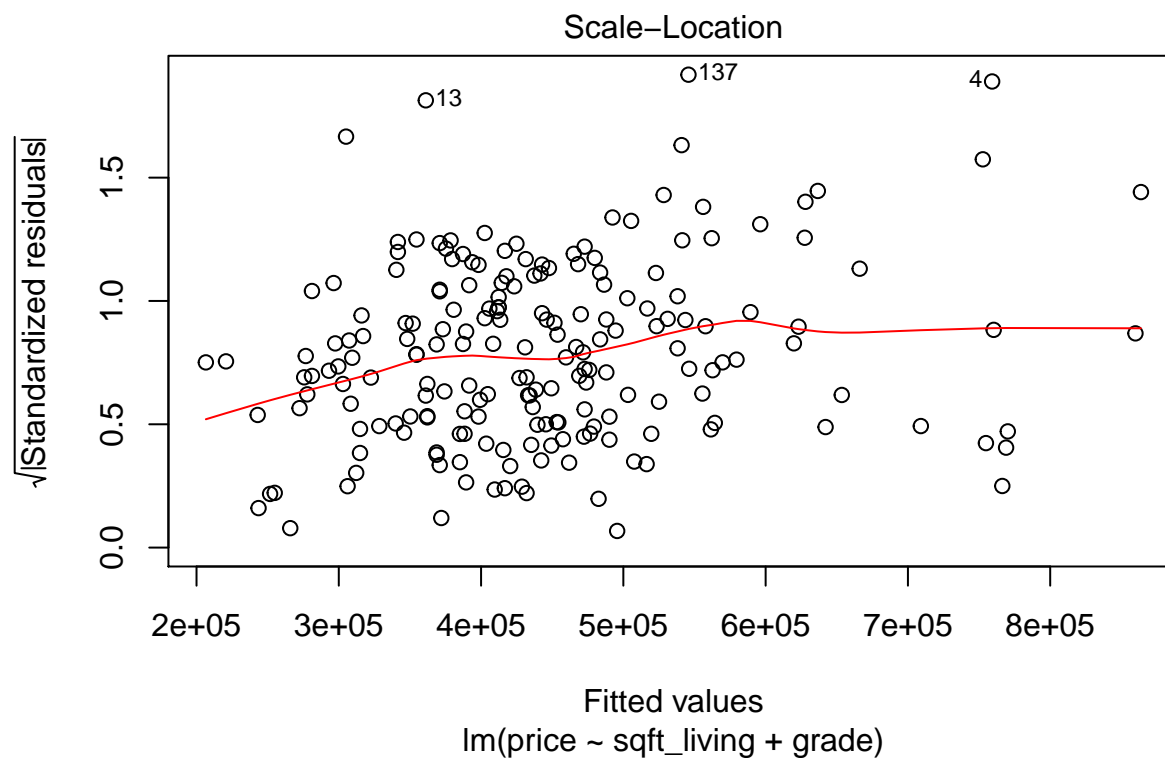
5.3 Homoscedasticity

In regression we assume that the standard deviations of the error terms (residuals) are constant and do not depend on the value of the predictors. So for example the variance of the residuals should be the same for 600 sqft apartments and 3000 sqft apartments.

We can check this assumption by looking at the plot of the standardized residuals and the predicted values, where we should see roughly equal variation at all values of the predicted values. There are also good statistical tests to check this assumption. The Breush-Pagan test is initiated with the `bptest()` function from the `lmtest` package, and the NCV test is performed by the `ncvTest()` function from base R (no package needed). A p-value < 0.05 in these tests would indicate a violation of the assumption of homoscedasticity, so it would mean that there is significant heteroscedasticity.

These test indicate that there is significant heteroscedasticity in `mod_house3`.

```
plot(mod_house3, which = 3)
```



```
ncvTest(mod_house3) # NCV test
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 17.01078, Df = 1, p = 3.7168e-05
```

```
bptest(mod_house3) # Breush-Pagan test
```

```
##
```



```
## studentized Breusch-Pagan test
##
## data:  mod_house3
## BP = 10.028, df = 2, p-value = 0.006645
```

5.3.1 What happens if the assumption of homoscedasticity is violated?

If there is heteroscedasticity, the prediction can be off, meaning that the model overall would be less efficient in predicting the outcome. This being said, we can still use the model for prediction, it will be just less good predicting new data.

More importantly, the model coefficients and their confidence intervals are no longer accurate.

So if we are interested in predicting the outcome, we could still do that with some success even if there is heteroscedasticity. However, interpreting the individual coefficients and our confidence in their added predictive value is no longer possible.

5.3.2 What to do if the assumption of homoscedasticity is violated?

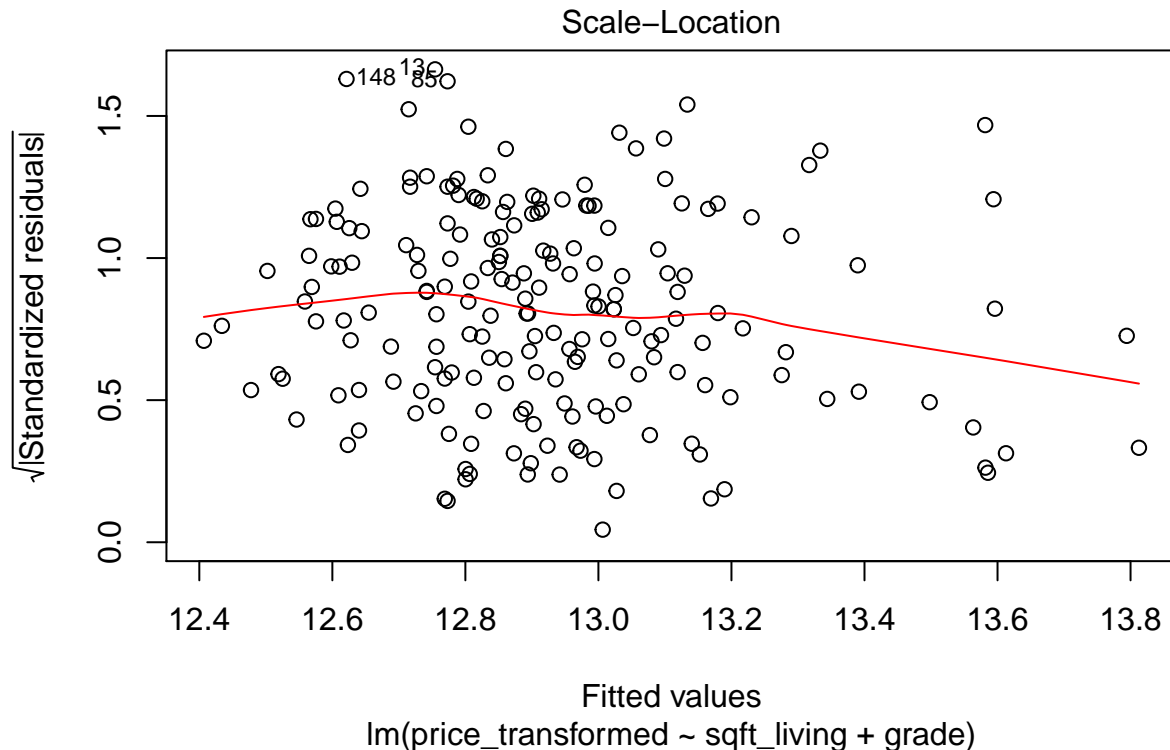
We have a couple of options to deal with heteroscedasticity.

1. **Transformation.** If we are mainly interested in more accurate predictions, we can do some transformation on the outcome and/or the predictors to make them more normally distributed, and thus, homogenize the variability in some parts of the dataset. For example in the example below we apply `log()` transformation to the outcome, price, and refit the model that way. The tests now show that the assumption is no longer violated. However in this case, we have to keep in mind when we interpret the model coefficient and the predicted (fitted) values that we are no longer predicting price, we the predictions are about $\log(\text{price})$. So the predicted values will have to be transformed back with the exponential function to be meaningful again, and the model coefficients now also have to be interpreted by keeping in mind that they refer to changes in $\log(\text{price})$.

```
data_house_nooutliers[, "price_transformed"] = log(data_house_nooutliers[,
  "price"])

mod_house4 = lm(price_transformed ~ sqft_living + grade, data = data_house_nooutliers)

plot(mod_house4, which = 3)
```



```
ncvTest(mod_house4) # NCV test
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.5274885, Df = 1, p = 0.46766
```

```
bptest(mod_house4) # Breush-Pagan test
```

```
##
## studentized Breusch-Pagan test
##
## data: mod_house4
## BP = 0.97249, df = 2, p-value = 0.6149
```

You can get the predictions for price, transformed back from `log()` to its regular scale with the `exp()` function

```
exp(predict(mod_house4))
```

2. **Using robust estimators.** If it is important to keep the outcome on its original scale to make the interpretation of the model coefficients more intuitive, we can use robust estimators to establish the Heteroscedasticity-consistent (HC) standard errors and use them to compute confidence intervals, and corrected p-values of the model coefficients. These new statistics are robust to the violation of homoscedasticity, so we call them robust estimates. In the example below we use the Huber-White Sandwich Estimator is used to get these robust values. For small samples (N around 50) you might have to use some other method to correct the standard error, for example Bell-McCaffrey estimates. For more details see the paper by Imbens and Kolesar (2016), and their R function at: <https://github.com/kolesarm/Robust-Small-Sample-Standard-Errors>

Reference: Imbens, G. W., & Kolesar, M. (2016). Robust standard errors in small samples: Some practical

advice. Review of Economics and Statistics, 98(4), 701-712.

```
# compute robust SE and p-values
mod_house3_sandwich_test = coeftest(mod_house3, vcov = vcovHC,
  type = "HC")
mod_house3_sandwich_test

##
## t test of coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -164126.218  81784.565 -2.0068  0.04615 *
## sqft_living   108.949    21.224  5.1332 6.864e-07 ***
## grade        57141.257  13001.055  4.3951 1.817e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

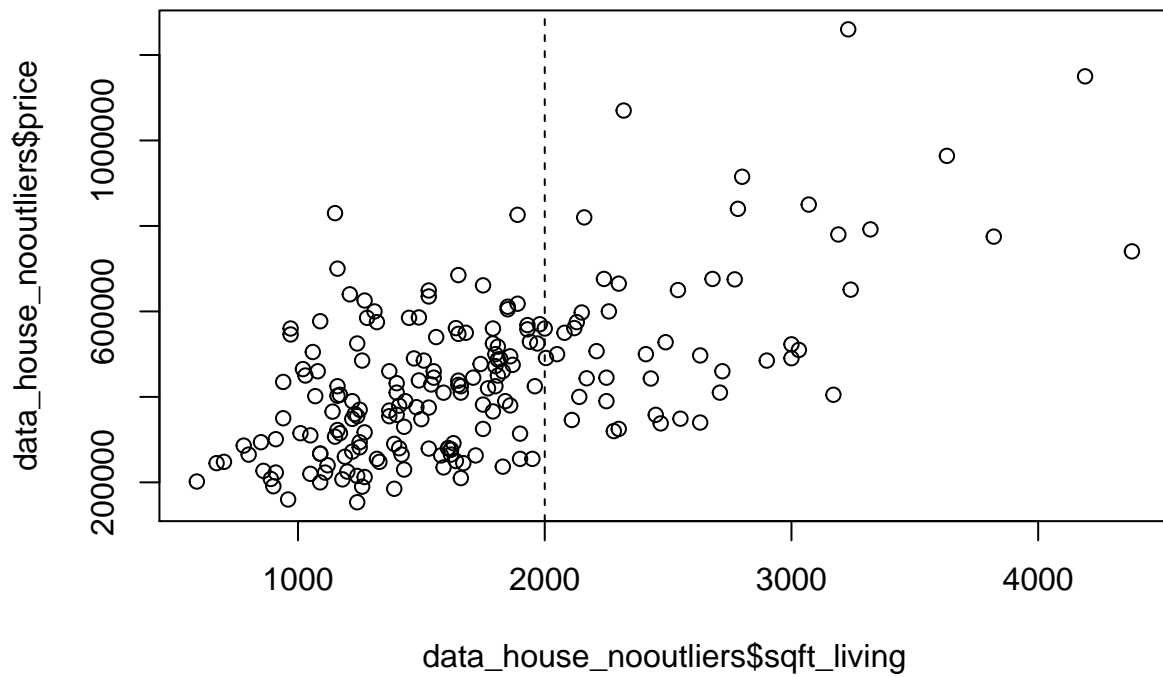
mod_house3_sandwich_se = unclass(mod_house3_sandwich_test)[,
  2]

# compute robust confidence intervals
CI95_lb_robust = coef(mod_house3) - 1.96 * mod_house3_sandwich_se
CI95_ub_robust = coef(mod_house3) + 1.96 * mod_house3_sandwich_se
cbind(mod_house3_sandwich_test, CI95_lb_robust, CI95_ub_robust)

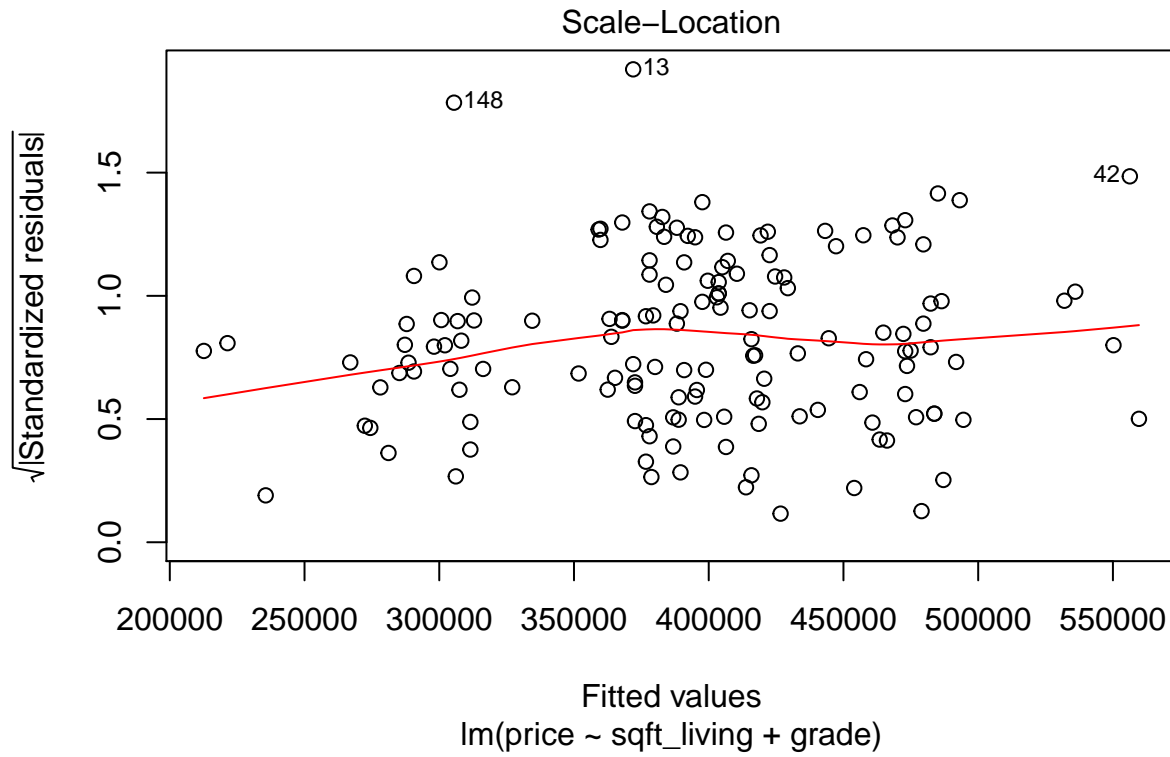
##              Estimate Std. Error  t value      Pr(>|t|) CI95_lb_robust
## (Intercept) -164126.2179 81784.56484 -2.006812 4.615088e-02 -324423.96495
## sqft_living   108.9489   21.22419  5.133242 6.864072e-07    67.34947
## grade        57141.2565 13001.05460  4.395125 1.816739e-05   31659.18951
##              CI95_ub_robust
## (Intercept)   -3828.4708
## sqft_living    150.5483
## grade         82623.3235
```

3. **Separate models for different parts of the dataset.** Another approach you might take is to visually inspect the plots of the relationships, to see if you can identify clusters with different variance. For example, in our case you might notice that the variances start to be larger at around 2000sqft apartment size. Here, we create two separate datasets for the part of the dataset containing 2000sqft or smaller apartments (`data_house_small`), and apartments larger than that (`data_house_large`), and fit separate models on these. So now if you would like to predict the price of apartments 2000sqft or smaller, you would use `mod_house3_small`, and for larger apartments we would use `mod_house3_large`. The homogeneity of variance tests indicate no heteroscedasticity in these two models.

```
plot(data_house_nooutliers$price ~ data_house_nooutliers$sqft_living)
abline(v = 2000, lty = "dashed")
```



```
data_house_small = data_house_nooutliers[data_house_nooutliers[,  
  "sqft_living"] <= 2000, ]  
data_house_large = data_house_nooutliers[data_house_nooutliers[,  
  "sqft_living"] > 2000, ]  
  
mod_house3_small = lm(price ~ sqft_living + grade, data = data_house_small)  
mod_house3_large = lm(price ~ sqft_living + grade, data = data_house_large)  
  
plot(mod_house3_small, which = 3)
```



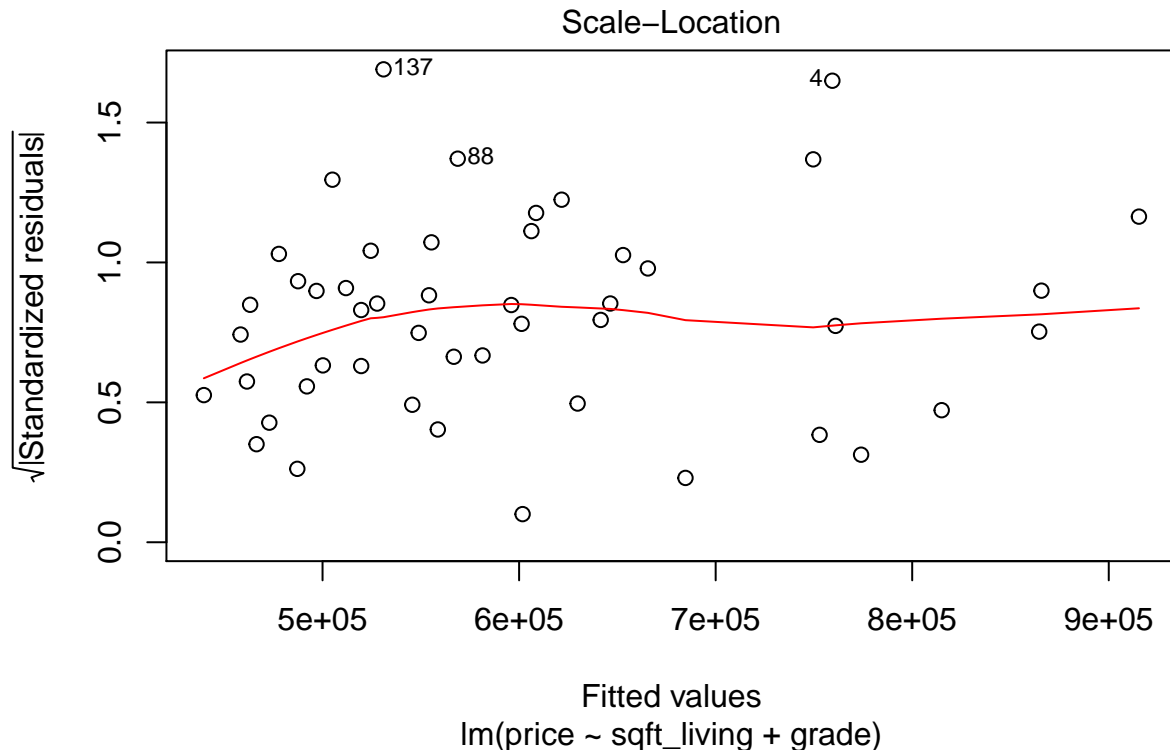
```
ncvTest(mod_house3_small) # NCV test
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.5286578, Df = 1, p = 0.46717
```

```
bptest(mod_house3_small) # Breush-Pagan test
```

```
##
## studentized Breusch-Pagan test
##
## data: mod_house3_small
## BP = 0.7299, df = 2, p-value = 0.6942
```

```
plot(mod_house3_large, which = 3)
```



```
ncvTest(mod_house3_large) # NCV test
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.8163077, Df = 1, p = 0.36626
```

```
bptest(mod_house3_large) # Breush-Pagan test
```

```
##
## studentized Breusch-Pagan test
##
## data: mod_house3_large
## BP = 0.86662, df = 2, p-value = 0.6484
```

5.4 No multicollinearity

We assume that none of the predictors can be linearly determined by the other predictor(s). This usually mean that the predictors should not be too highly correlated with each other, but there are some more subtle cases where even though the predictors are not super highly correlated pair-by-pair, but a combination of multiple predictors can still accurately determin another predictor.

So to evaluate whether this assumption holds true, we compute the variance inflation factor for each predictor using the `vif()` function in the `car` package.

There is no real consensus about what `vif` values indicates problematic multicollinearity, some advocate for `vif` values of 10 or larger are problematic (for example Montgomery and Peck ,1992). However, a more conservative approach is to explore and treat multicollinearity if VIF is above 3 (an approach suggested by Zuur, Ieno, and Elphick, 2010).

Reference: Montgomery, D.C. & Peck, E.A. (1992) Introduction to Linear Regression Analysis. Wiley, New York. Zuur, A. F., Ieno, E. N., & Elphick, C. S. (2010). A protocol for data exploration to avoid common statistical problems. Methods in ecology and evolution, 1(1), 3-14.

```
vif(mod_house3)
```

```
## sqft_living      grade
##      1.849221      1.849221
```

In the example above multicollinearity is not problematic.

5.4.1 What happens in the case of multicollinearity?

In regression we would often like to know the main effect of each predictor, and whether this effect (or in other words, whether the added predictive value of a certain predictor) is statistically significant or not. The model coefficient indicates the size and direction of the effect of increasing the value of the predictor by one point, when all other predictors are held constant. However, if correlation is high among the predictors, there are fewer cases in the dataset where one predictor has high values, while the others are low. This makes it difficult for the model to estimate the independent effect of the predictors that are highly correlated.

The bottom line is, that in the case of multicollinearity, the model coefficients and the t-test of their individual predictive value will become less reliable. Furthermore, the model coefficients will be very unstable, and will show big changes and even sign changes (a positive effect would turn into a negative effect) if we change the model a bit, such as by adding or removing a new predictor.

Fortunately, multicollinearity does not affect prediction accuracy. So if prediction accuracy is our only concern, we might not have to deal with multicollinearity. But if we are interested in the coefficients as well and drawing inference about the individual predictors and their effect and predictive value, we have to address it.

5.4.2 What to do in the case of multicollinearity?

There are two types of multicollinearity:

Structural multicollinearity, where multicollinearity is a result of entering a predictor into the model that is derived from one or more of the other predictors. For example higher order terms (e.g. grade^2), or interaction terms (e.g. $\text{long} \times \text{lat}$).

Data multicollinearity, where multicollinearity is present in the data itself rather than being an artifact of our model.

To demonstrating the two types of multicollinearity and how to address them, we will build some new models.

5.4.2.1 Example of handling structural multicollinearity

First, let's build a model where in addition to the `sqft_living` and `grade`, we also enter geolocation information into the model as predictors.

In the first model we only add the main effect of these predictors (no interaction term). Notice that longitude has a negative coefficient indicating that further east we go the lower the price of the apartment which makes sense if you look at the map of King County with Seattle being in the west of the county. Latitude has a positive coefficient, indicating that the further north we go the higher the price. The added predictive value of latitude is significant in this model. Also notice that there is no multicollinearity in the model according to `vif`.

```
mod_house_geolocation = lm(price ~ sqft_living + grade + long +
  lat, data = data_house_nooutliers)
summary(mod_house_geolocation)
```

```
##
## Call:
```

```
## lm(formula = price ~ sqft_living + grade + long + lat, data = data_house_nooutliers)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -239590  -79244  -14858   55982  483474
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -3.478e+07  6.672e+06  -5.214 4.74e-07 ***
## sqft_living  1.200e+02  1.783e+01   6.730 1.89e-10 ***
## grade        5.267e+04  1.156e+04   4.555 9.29e-06 ***
## long       -7.518e+04  5.262e+04  -1.429   0.155
## lat         5.350e+05  6.257e+04   8.551 3.72e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 120900 on 193 degrees of freedom
## Multiple R-squared:  0.5858, Adjusted R-squared:  0.5772
## F-statistic: 68.25 on 4 and 193 DF,  p-value: < 2.2e-16
```

```
vif(mod_house_geolocation)
```

```
## sqft_living      grade      long      lat
##    1.884575    1.860764    1.041486    1.024952
```

Now, we add an interaction term as well (using the * sign instead of +), where we include the interaction of longitude and latitude in the model. Notice that coefficients and the signs underwent a violent shift. Longitude now has a positive coefficient, latitude has a negative coefficient, and latitude's added predictive value no longer seems significant. However, the vif values of long, lat and long:lat are extremely high indicating severe multicollinearity. The big change in the results is because of the coefficient's instability due to multicollinearity. Even though the coefficients underwent violent change, the model R^2 only changed a little bit.

```
mod_house_geolocation_inter = lm(price ~ sqft_living + grade +
  long * lat, data = data_house_nooutliers)
summary(mod_house_geolocation_inter)
```

```
##
## Call:
## lm(formula = price ~ sqft_living + grade + long * lat, data = data_house_nooutliers)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -223574  -75687  -16930   58851  479323
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.974e+09  2.656e+09   1.873  0.0626 .
## sqft_living  1.206e+02  1.772e+01   6.809 1.23e-10 ***
## grade        5.436e+04  1.152e+04   4.717 4.59e-06 ***
## long        4.091e+07  2.173e+07   1.882  0.0613 .
## lat       -1.048e+08  5.585e+07  -1.876  0.0621 .
## long:lat    -8.618e+05  4.570e+05  -1.886  0.0608 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 120100 on 192 degrees of freedom
```



```
## Multiple R-squared:  0.5934, Adjusted R-squared:  0.5828
## F-statistic: 56.03 on 5 and 192 DF,  p-value: < 2.2e-16
```

```
vif(mod_house_geolocation_inter)
```

```
##      sqft_living      grade      long      lat      long:lat
## 1.885276e+00 1.872074e+00 1.799568e+05 8.275735e+05 1.125795e+06
```

Structural multicollinearity in this case is due to the fact that the interaction term long:lat is derived from both long and lat predictors, so they are now highly correlated.

This can be fixed by standardizing the variables. A good method to use here is “centering”, that is, subtracting the mean of the variable from the values of the variable. By doing this, we can still preserve the original scale of the variables, and thus, the coefficients will mean the same thing as before centering. (We could use Z transformation as well, but that would change the interpretation of the coefficients, because the predictors would be transformed to a scale between -1 to +1.)

So let's center longitude and latitude, and re-build the model. Notice that multicollinearity is gone, and the coefficients are now back to their original signs and roughly the same effect size as before entering the interaction. Latitude has a significant added predictive value again. But also notice that the R^2 was not at all effected by the removal of multicollinearity! So the prediction efficiency is the same, but the coefficients are more reliable now and their interpretation is more straightforward.

```
data_house_nooutliers$long_centered = data_house_nooutliers$long -
  mean(data_house_nooutliers$long)
data_house_nooutliers$lat_centered = data_house_nooutliers$lat -
  mean(data_house_nooutliers$lat)

mod_house_geolocation_inter_centered = lm(price ~ sqft_living +
  grade + long_centered * lat_centered, data = data_house_nooutliers)

vif(mod_house_geolocation_inter_centered)
```

```
##              sqft_living              grade
##              1.885276              1.872074
##              long_centered      lat_centered
##              1.058954              1.049776
## long_centered:lat_centered
##              1.050759
```

```
summary(mod_house_geolocation_inter_centered)
```

```
##
## Call:
## lm(formula = price ~ sqft_living + grade + long_centered * lat_centered,
##     data = data_house_nooutliers)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -223574  -75687  -16930   58851  479323
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -166830.36    68692.37  -2.429   0.0161 *
## sqft_living       120.63      17.72    6.809 1.23e-10 ***
## grade           54356.80   11522.67   4.717 4.59e-06 ***
## long_centered   -87942.98   52713.77  -1.668   0.0969 .
## lat_centered    516751.63   62905.43   8.215 3.08e-14 ***
```

```
## long_centered:lat_centered -861819.77  456978.57  -1.886   0.0608 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 120100 on 192 degrees of freedom
## Multiple R-squared:  0.5934, Adjusted R-squared:  0.5828
## F-statistic: 56.03 on 5 and 192 DF,  p-value: < 2.2e-16
```

5.4.2.2 Example of handling data multicollinearity

Now let's build a model where in addition to the `sqft_living` and `grade`, we also use `sqft_above` as a predictor (the squarefootage of the area of the apartment above ground level).

The vif is now above 3.

```
mod_house5 = lm(price ~ sqft_living + grade + sqft_above, data = data_house)

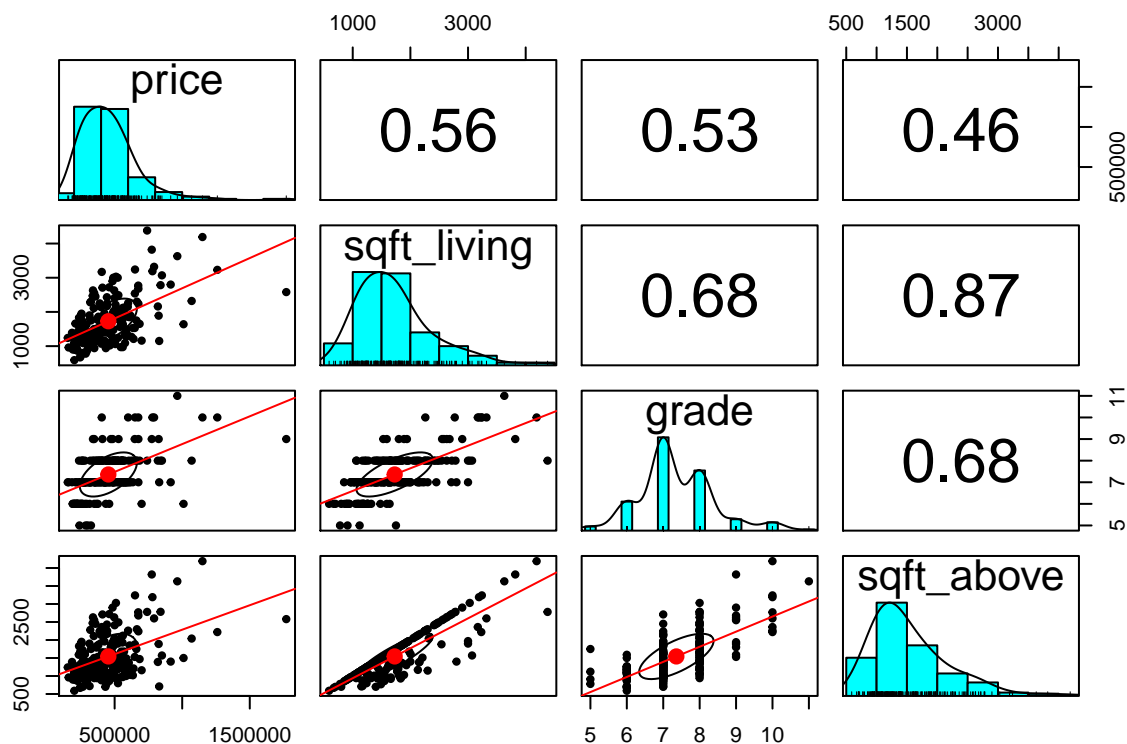
vif(mod_house5)
```

```
## sqft_living      grade sqft_above
##    4.505215    1.972816    4.532763
```

We can explore the reason for this by looking at the correlation matrix of the predictors. The `pairs.panels()` function returns a series of very useful diagnostic plots where we can not only see the correlation of the variables, but we can also see the distribution of the variables and the scatterplot of their pairwise relationships.

The correlation matrix clearly indicates that the correlation of `sqft_living` and `sqft_above` is very high.

```
pairs.panels(data_house[, c("price", "sqft_living", "grade",
                             "sqft_above")], col = "red", lm = T)
```



Compare the coefficients in the model summary of the two models. What can you observe? Do the coefficients in the new model make sense?

```
summary(mod_house5)
```

```
##
## Call:
## lm(formula = price ~ sqft_living + grade + sqft_above, data = data_house)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -417412  -90760  -25309   71154 1144030
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -215411.12   95961.49  -2.245  0.0259 *
## sqft_living    185.94     38.20    4.868 2.32e-06 ***
## grade         66522.66   16386.54    4.060 7.10e-05 ***
## sqft_above     -91.88     40.32   -2.279  0.0238 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 168300 on 196 degrees of freedom
## Multiple R-squared:  0.3746, Adjusted R-squared:  0.365
## F-statistic: 39.13 on 3 and 196 DF, p-value: < 2.2e-16
```

```
summary(mod_house3)
```

```
##
## Call:
## lm(formula = price ~ sqft_living + grade, data = data_house_nooutliers)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -347654  -95907  -16999   73466  524235
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -164126.22   81165.63  -2.022   0.0445 *
## sqft_living    108.95     20.92    5.208 4.82e-07 ***
## grade        57141.26   13653.24    4.185 4.31e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 143200 on 195 degrees of freedom
## Multiple R-squared:  0.413, Adjusted R-squared:  0.407
## F-statistic: 68.59 on 2 and 195 DF, p-value: < 2.2e-16
```

Because of the multicollinearity, we cannot trust the coefficients or the t-tests of the predictors.

You have a number of options to deal with multicollinearity:

1. removing highly correlated predictors
2. linearly combining predictors, such as using the mean of the two values for each observation (e.g. for case 3, sqft_living is 1060, sqft_above is 960, so the mean for this case would be 1000). However, it is not clear how we would interpret this value in this particular case of predicting housing prices.
3. you could use some other statistical method entirely (e.g. partial least squares regression or principal components analysis)

Here I suggest using method 1., realizing that there is hardly any difference in the two variables sqft_living and sqft_above, and when they are different it is due to the apartment having a basement. So we should just settle with using one of these variables, the one that makes more sense for the types of apartments the price of which we want to be able to predict. In this case I personally would keep sqft_living in the model, because I have a theory that the size of the living area is what ultimately influences the price, whether the living area is above or below ground, and the information about the apartment having a basement can be added by using the has_basement variable in a later model. However, those more familiar with the literature of predicting housing price might use previous research results to make this decision. Or the decision might be made based on practical grounds, for example the fact that size of the living area is more accessible information than size of the area above ground, so we would be able to use the model practically in more cases for prediction, if we chose sqft_living as a predictor.

Notice that I did not use model comparison to make a decision about which variable to exclude! This decision, like all model selection decisions, should be made based on theoretical grounds, or based on previous research results or practical issues.

You can find some additional resources about multicollinearity on these links:

<https://statisticalhorizons.com/multicollinearity>

<http://blog.minitab.com/blog/understanding-statistics/handling-multicollinearity-in-regression-analysis>

<http://statisticsbyjim.com/regression/multicollinearity-in-regression-analysis/>