

STOR 455 Homework #3

40 points - Due 2/21 at 5:00pm

Directions: You will be assigned to a group of three to four students for this assignment. Parts 1, 2, & 5 should be turned in individually to Gradescope by each student in your group. Parts 3 & 4 should be submitted as a group to Gradescope. There are separate places to submit the individual and group portions of the assignment.

Situation: Can we predict the selling price of a house in Ames, Iowa based on recorded features of the house? That is your task for this assignment. Each group will have a dataset with information on forty potential predictors and the selling price (in \$1,000's) for a sample of homes. The data set for your group is in AmesTrain???.csv (where ?? corresponds to your group number) and can be found in the AmesTrain zipped file under class 14 in Sakai. A separate file identifies the variables in the Ames Housing data and explains some of the coding.

```
library(readr)
library(car)

## 载入需要的程辑包: carData

train <- read_csv("AmesTrain15.csv")

## Rows: 600 Columns: 42

## -- Column specification -----
## Delimiter: ","
## chr (15): LotConfig, HouseStyle, ExteriorQ, ExteriorC, Foundation, BasementH...
## dbl (27): Order, Price, LotFrontage, LotArea, Quality, Condition, YearBuilt,...

##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.

head(train)

## # A tibble: 6 x 42
##   Order Price LotFrontage LotArea LotConfig HouseStyle Quality Condition
##   <dbl> <dbl>         <dbl>   <dbl> <chr>      <chr>          <dbl>   <dbl>
## 1 >
```

```
## 1  101 178.      0   2117 Inside   2Story      6      5
## 2  533 209      65   8127 Inside   2Story      7      5
## 3 2610 159      77   9786 Inside   1.5Fin      6      7
## 4 1186 174      80  10400 Inside   1Story      7      6
## 5 1429 225       0  11454 Corner   SLvl       8      5
## 6 1246 116      60  11556 Inside   1Story      5      6
## # ... with 34 more variables: YearBuilt <dbl>, YearRemodel <dbl>,
## #   ExteriorQ <chr>, ExteriorC <chr>, Foundation <chr>, BasementHt <c
## #   BasementC <chr>, BasementFin <chr>, BasementFinSF <dbl>,
## #   BasementUnFinSF <dbl>, BasementSF <dbl>, Heating <chr>, HeatingQC
## #   CentralAir <chr>, FirstSF <dbl>, SecondSF <dbl>, GroundSF <dbl>,
## #   BasementFBath <dbl>, BasementHBath <dbl>, FullBath <dbl>, HalfBat
## #   h <dbl>,
## #   Bedroom <dbl>, KitchenQ <chr>, TotalRooms <dbl>, Fireplaces <dbl>,
## #   ...
```

Part 1. Build an initial basic model

Your basic model can use any of the quantitative variables in the dataset but should NOT use the categorical variables, transformations, or interactions (we'll discuss these in class soon) – those will come in a later assignment.

```
library(dplyr)

##
## 载入程辑包: 'dplyr'

## The following object is masked from 'package:car':
##
##   recode

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

myDataNum = select_if(train, is.numeric)
myData = subset(myDataNum, select = -c(Order))
myData = subset(myData, select = -c(GroundSF))
head(myDataNum)

## # A tibble: 6 x 27
##   Order Price LotFrontage LotArea Quality Condition YearBuilt YearRem
##   <dbl> <dbl>      <dbl>   <dbl>   <dbl>      <dbl>      <dbl>      <db
##   1>
```

```
## 1  101 178.      0  2117      6      5      2000      2000
## 2  533 209      65  8127      7      5      2003      2003
## 3 2610 159      77  9786      6      7      1962      198
1
## 4 1186 174      80 10400      7      6      1970      197
0
## 5 1429 225      0 11454      8      5      1995      199
5
## 6 1246 116      60 11556      5      6      1952      195
2
## # ... with 19 more variables: BasementFinSF <dbl>, BasementUnFinSF <d
bl>,
## #   BasementSF <dbl>, FirstSF <dbl>, SecondSF <dbl>, GroundSF <dbl>,
## #   BasementFBath <dbl>, BasementHBath <dbl>, FullBath <dbl>, HalfBat
h <dbl>,
## #   Bedroom <dbl>, TotalRooms <dbl>, Fireplaces <dbl>, GarageCars <db
l>,
## #   GarageSF <dbl>, WoodDeckSF <dbl>, OpenPorchSF <dbl>, EnclosedPorc
hSF <dbl>,
## #   ScreenPorchSF <dbl>
```

```
head(myData)
```

```
## # A tibble: 6 x 25
##   Price LotFrontage LotArea Quality Condition YearBuilt YearRemodel
##   <dbl>      <dbl>   <dbl>   <dbl>      <dbl>      <dbl>      <dbl>
## 1 178.      0     2117      6        5      2000      2000
## 2 209      65     8127      7        5      2003      2003
## 3 159      77     9786      6        7      1962      1981
## 4 174      80    10400      7        6      1970      1970
## 5 225      0    11454      8        5      1995      1995
## 6 116      60    11556      5        6      1952      1952
## # ... with 18 more variables: BasementFinSF <dbl>, BasementUnFinSF <d
bl>,
## #   BasementSF <dbl>, FirstSF <dbl>, SecondSF <dbl>, BasementFBath <d
bl>,
## #   BasementHBath <dbl>, FullBath <dbl>, HalfBath <dbl>, Bedroom <db
l>,
## #   TotalRooms <dbl>, Fireplaces <dbl>, GarageCars <dbl>, GarageSF <d
bl>,
## #   WoodDeckSF <dbl>, OpenPorchSF <dbl>, EnclosedPorchSF <dbl>,
## #   ScreenPorchSF <dbl>
```

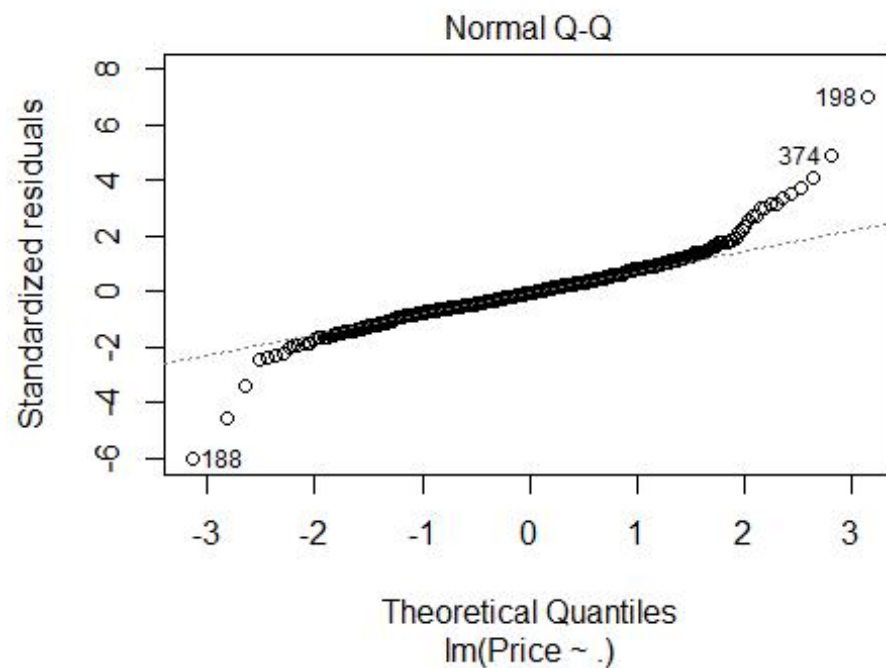
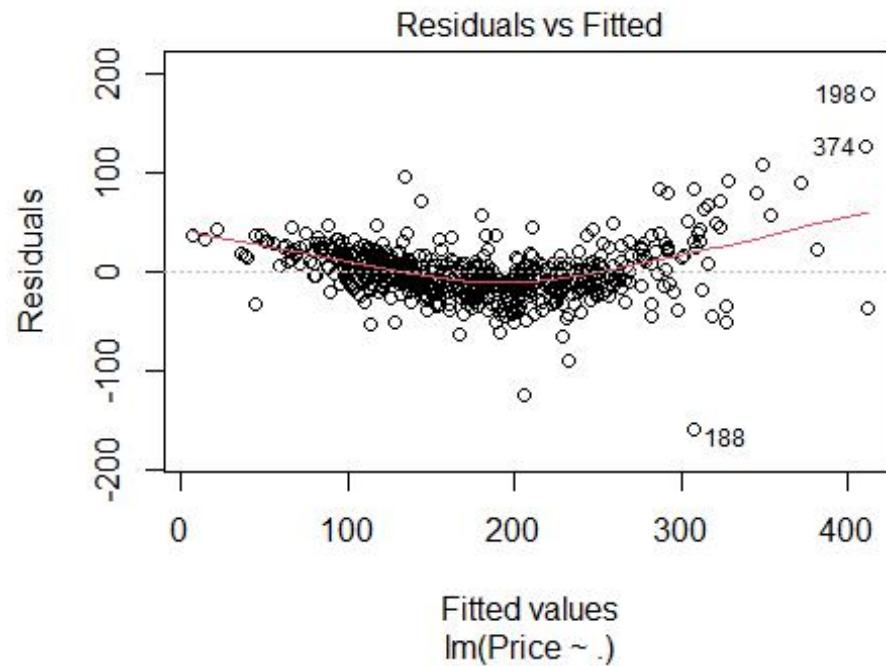
Use your data to select a set of predictors to include in your model. Keep track of the process you use and decisions you make to arrive at an initial set of predictors. Your report should include a summary of this process. You don't need to show all the output for every model you consider, but you should give **a clear description of the path you took and the criteria that you used to compare competing models**. Also, use at least two model selection methods to find a model (e.g. don't just check

all subsets, although it will work well here, this method will fail in future assignments).

In addition to the commentary on model selection, include the following information for this initial choice of a model: **the `summary()` output for your model, comments on which (if any) of the predictors in the model are not significant at a 5% level, and comments on what the VIF values tell you about the individual predictors in your model.**

Do not consider the Order variable (that is just an observation number) as one of your predictors. Avoid predictors that are exactly related. For example, if $\text{GroundSF} = \text{FirstSF} + \text{SecondSF}$ you will likely get trouble if you try to put all three in the same model.

```
modPrice = lm(Price~., data=myData)
plot(modPrice, 1:2)
```



Check linear conditions: - Linearity is questionable, since the residuals vs fitted plot is better described by a curve rather than a horizontal line. - Constant variance appears to be okay. - Normality of residuals is questionable, since the residuals have obvious skewness from the theoretical quantiles based on the qq plot.

(1) backward elimination

```
Full = lm(Price~., data=myData)
MSE = (summary(Full)$sigma)^2
backward_mod = step(Full, scale=MSE, trace=FALSE)
summary(backward_mod)

##
## Call:
## lm(formula = Price ~ LotFrontage + LotArea + Quality + Condition +
##     YearBuilt + YearRemodel + BasementFinSF + BasementSF + FirstSF +
##     SecondSF + Bedroom + Fireplaces + GarageSF + EnclosedPorchSF +
##     ScreenPorchSF, data = myData)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -157.706  -15.568   -1.546   12.088  181.205
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.510e+03  1.471e+02  -10.261  < 2e-16 ***
## LotFrontage    1.049e-01  3.683e-02   2.849  0.004543 **
## LotArea        6.678e-04  1.030e-04   6.484  1.91e-10 ***
## Quality        1.306e+01  1.347e+00   9.692  < 2e-16 ***
## Condition      4.335e+00  1.179e+00   3.676  0.000259 ***
## YearBuilt      4.009e-01  6.411e-02   6.253  7.76e-10 ***
## YearRemodel    3.262e-01  7.752e-02   4.207  2.99e-05 ***
## BasementFinSF  2.083e-02  3.131e-03   6.652  6.65e-11 ***
## BasementSF     2.503e-02  6.148e-03   4.071  5.32e-05 ***
## FirstSF        6.038e-02  7.541e-03   8.006  6.43e-15 ***
## SecondSF       5.766e-02  4.371e-03  13.192  < 2e-16 ***
## Bedroom       -3.714e+00  1.929e+00  -1.925  0.054682 .
## Fireplaces     4.226e+00  2.165e+00   1.952  0.051419 .
## GarageSF       3.706e-02  6.850e-03   5.411  9.15e-08 ***
## EnclosedPorchSF 4.373e-02  2.236e-02   1.956  0.050946 .
## ScreenPorchSF  5.201e-02  1.775e-02   2.929  0.003528 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 27.22 on 584 degrees of freedom
## Multiple R-squared:  0.8649, Adjusted R-squared:  0.8614
## F-statistic: 249.3 on 15 and 584 DF,  p-value: < 2.2e-16
```

(2) forward selection

```
none = lm(Price~1, data=myData)
forward_mod = step(none, scope=list(upper=Full), scale=MSE, direction="
forward", trace=FALSE)
summary(forward_mod)

##
## Call:
```

```
## lm(formula = Price ~ Quality + FirstSF + SecondSF + YearBuilt +
##     BasementFinSF + GarageSF + LotArea + YearRemodel + BasementSF +
##     Condition + ScreenPorchSF + LotFrontage + Bedroom + EnclosedPorch
SF +
##     Fireplaces, data = myData)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -157.706  -15.568   -1.546   12.088  181.205
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.510e+03  1.471e+02 -10.261  < 2e-16 ***
## Quality       1.306e+01  1.347e+00   9.692  < 2e-16 ***
## FirstSF       6.038e-02  7.541e-03   8.006 6.43e-15 ***
## SecondSF      5.766e-02  4.371e-03  13.192  < 2e-16 ***
## YearBuilt     4.009e-01  6.411e-02   6.253 7.76e-10 ***
## BasementFinSF  2.083e-02  3.131e-03   6.652 6.65e-11 ***
## GarageSF      3.706e-02  6.850e-03   5.411 9.15e-08 ***
## LotArea       6.678e-04  1.030e-04   6.484 1.91e-10 ***
## YearRemodel   3.262e-01  7.752e-02   4.207 2.99e-05 ***
## BasementSF    2.503e-02  6.148e-03   4.071 5.32e-05 ***
## Condition     4.335e+00  1.179e+00   3.676 0.000259 ***
## ScreenPorchSF  5.201e-02  1.775e-02   2.929 0.003528 **
## LotFrontage   1.049e-01  3.683e-02   2.849 0.004543 **
## Bedroom      -3.714e+00  1.929e+00  -1.925 0.054682 .
## EnclosedPorchSF 4.373e-02  2.236e-02   1.956 0.050946 .
## Fireplaces     4.226e+00  2.165e+00   1.952 0.051419 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 27.22 on 584 degrees of freedom
## Multiple R-squared:  0.8649, Adjusted R-squared:  0.8614
## F-statistic: 249.3 on 15 and 584 DF, p-value: < 2.2e-16
```

(3) stepwise regression

```
stepwise_mod = step(none, scope=list(upper=Full), scale=MSE, trace=FALSE)
summary(stepwise_mod)

##
## Call:
## lm(formula = Price ~ Quality + FirstSF + SecondSF + YearBuilt +
##     BasementFinSF + GarageSF + LotArea + YearRemodel + BasementSF +
##     Condition + ScreenPorchSF + LotFrontage + Bedroom + EnclosedPorch
SF +
##     Fireplaces, data = myData)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -157.706 -15.568 -1.546 12.088 181.205
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.510e+03 1.471e+02 -10.261 < 2e-16 ***
## Quality      1.306e+01 1.347e+00  9.692 < 2e-16 ***
## FirstSF      6.038e-02 7.541e-03  8.006 6.43e-15 ***
## SecondSF     5.766e-02 4.371e-03 13.192 < 2e-16 ***
## YearBuilt    4.009e-01 6.411e-02  6.253 7.76e-10 ***
## BasementFinSF 2.083e-02 3.131e-03  6.652 6.65e-11 ***
## GarageSF     3.706e-02 6.850e-03  5.411 9.15e-08 ***
## LotArea      6.678e-04 1.030e-04  6.484 1.91e-10 ***
## YearRemodel  3.262e-01 7.752e-02  4.207 2.99e-05 ***
## BasementSF   2.503e-02 6.148e-03  4.071 5.32e-05 ***
## Condition    4.335e+00 1.179e+00  3.676 0.000259 ***
## ScreenPorchSF 5.201e-02 1.775e-02  2.929 0.003528 **
## LotFrontage  1.049e-01 3.683e-02  2.849 0.004543 **
## Bedroom     -3.714e+00 1.929e+00 -1.925 0.054682 .
## EnclosedPorchSF 4.373e-02 2.236e-02  1.956 0.050946 .
## Fireplaces   4.226e+00 2.165e+00  1.952 0.051419 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 27.22 on 584 degrees of freedom
## Multiple R-squared:  0.8649, Adjusted R-squared:  0.8614
## F-statistic: 249.3 on 15 and 584 DF, p-value: < 2.2e-16
```

Summary: Backward elimination: `lm(formula = Price ~ LotFrontage + LotArea + Quality + Condition + YearBuilt + YearRemodel + BasementFinSF + BasementSF + FirstSF + SecondSF + Bedroom + Fireplaces + GarageSF + EnclosedPorchSF + ScreenPorchSF, data = myData)` Forward selection: `lm(formula = Price ~ Quality + FirstSF + SecondSF + YearBuilt + BasementFinSF + GarageSF + LotArea + YearRemodel + BasementSF + Condition + ScreenPorchSF + LotFrontage + Bedroom + EnclosedPorchSF + Fireplaces, data = myData)` Stepwise regression: `lm(formula = Price ~ Quality + FirstSF + SecondSF + YearBuilt + BasementFinSF + GarageSF + LotArea + YearRemodel + BasementSF + Condition + ScreenPorchSF + LotFrontage + Bedroom + EnclosedPorchSF + Fireplaces, data = myData)`

All three models give the same set of predictors; these predictors are my initial selection of predictors.

```
summary(backward_mod)
```

```
##
## Call:
## lm(formula = Price ~ LotFrontage + LotArea + Quality + Condition +
##     YearBuilt + YearRemodel + BasementFinSF + BasementSF + FirstSF +
##     SecondSF + Bedroom + Fireplaces + GarageSF + EnclosedPorchSF +
##     ScreenPorchSF, data = myData)
##
```



```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -157.706  -15.568   -1.546   12.088  181.205
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.510e+03  1.471e+02 -10.261  < 2e-16 ***
## LotFrontage    1.049e-01  3.683e-02   2.849  0.004543 **
## LotArea        6.678e-04  1.030e-04   6.484  1.91e-10 ***
## Quality        1.306e+01  1.347e+00   9.692  < 2e-16 ***
## Condition      4.335e+00  1.179e+00   3.676  0.000259 ***
## YearBuilt       4.009e-01  6.411e-02   6.253  7.76e-10 ***
## YearRemodel    3.262e-01  7.752e-02   4.207  2.99e-05 ***
## BasementFinSF   2.083e-02  3.131e-03   6.652  6.65e-11 ***
## BasementSF      2.503e-02  6.148e-03   4.071  5.32e-05 ***
## FirstSF         6.038e-02  7.541e-03   8.006  6.43e-15 ***
## SecondSF        5.766e-02  4.371e-03  13.192  < 2e-16 ***
## Bedroom       -3.714e+00  1.929e+00  -1.925  0.054682 .
## Fireplaces      4.226e+00  2.165e+00   1.952  0.051419 .
## GarageSF        3.706e-02  6.850e-03   5.411  9.15e-08 ***
## EnclosedPorchSF 4.373e-02  2.236e-02   1.956  0.050946 .
## ScreenPorchSF   5.201e-02  1.775e-02   2.929  0.003528 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 27.22 on 584 degrees of freedom
## Multiple R-squared:  0.8649, Adjusted R-squared:  0.8614
## F-statistic: 249.3 on 15 and 584 DF,  p-value: < 2.2e-16
```

Comments on which of the predictors in the model are not significant at a 5% level:
The following predictors have p values greater than 0.05: Bedroom, Fireplaces, EnclosedPorchSF.

```
vif(backward_mod)
```

```
##      LotFrontage      LotArea      Quality      Condition      Y
##      1.103602      1.134202      2.805858      1.421339      3.
232290
##      YearRemodel  BasementFinSF  BasementSF      FirstSF
##      2.174244      1.332166      4.520978      5.288575      2.
857212
##      Bedroom      Fireplaces      GarageSF  EnclosedPorchSF  Scree
enPorchSF
##      1.804941      1.471388      1.858104      1.311854      1.
099532
```

Comments on what the VIF values tell you about the individual predictors in your model: A high VIF (VIF > 5) indicates that the associated independent variable is

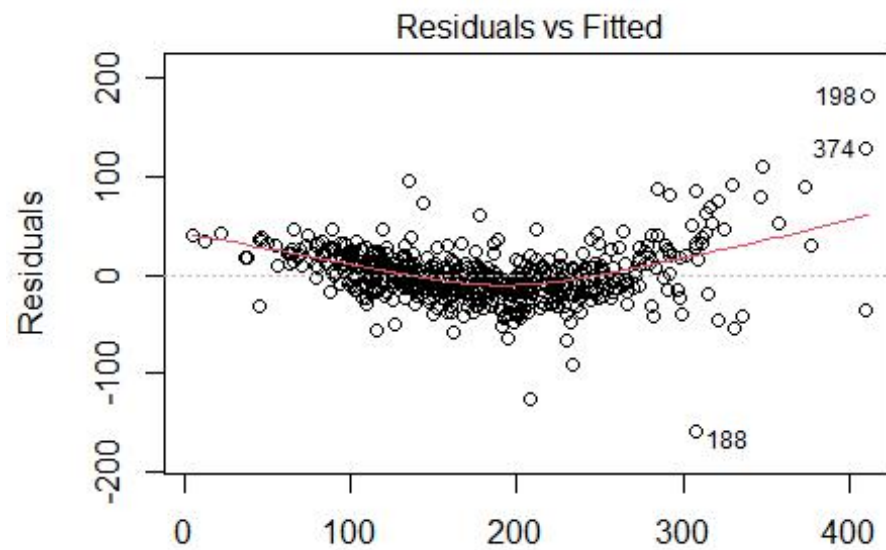
highly collinear with the other variables in the model. Here, the only variable that is highly collinear with the other variables in the model is FirstSF.

Part 2. Residual analysis for your basic model

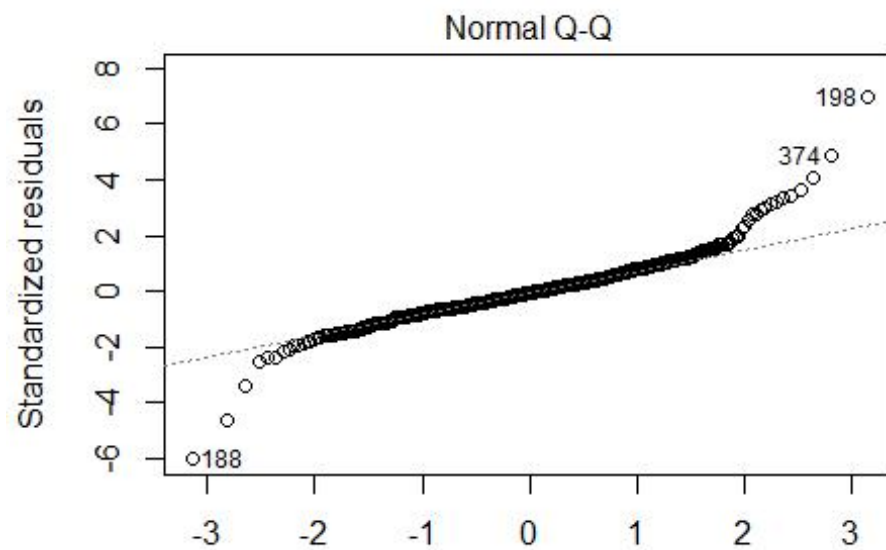
Do a residual analysis for the model you chose in Part 1. Include any plots relevant to checking model conditions - with interpretations. Also check whether any of the data cases are unusual with respect to studentized residuals. Since there are a lot of data points don't worry about the "mild" cases for studentized residuals, but indicate what specific criteria you are using to identify "unusual" points.

Adjust your model (either the predictors included or data values that are used to fit it, but not yet using transformations) on the basis of your residual analysis – but don't worry too much about trying to get all conditions "perfect". For example, don't automatically just delete any points that might give large residuals! If you do refit something, be sure to document what changed and include the new summary() output.

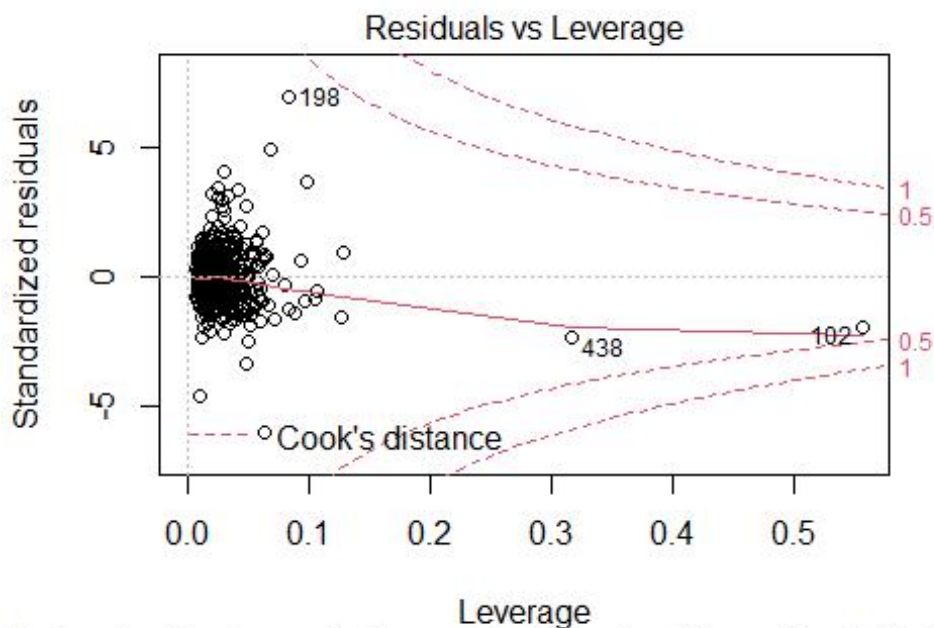
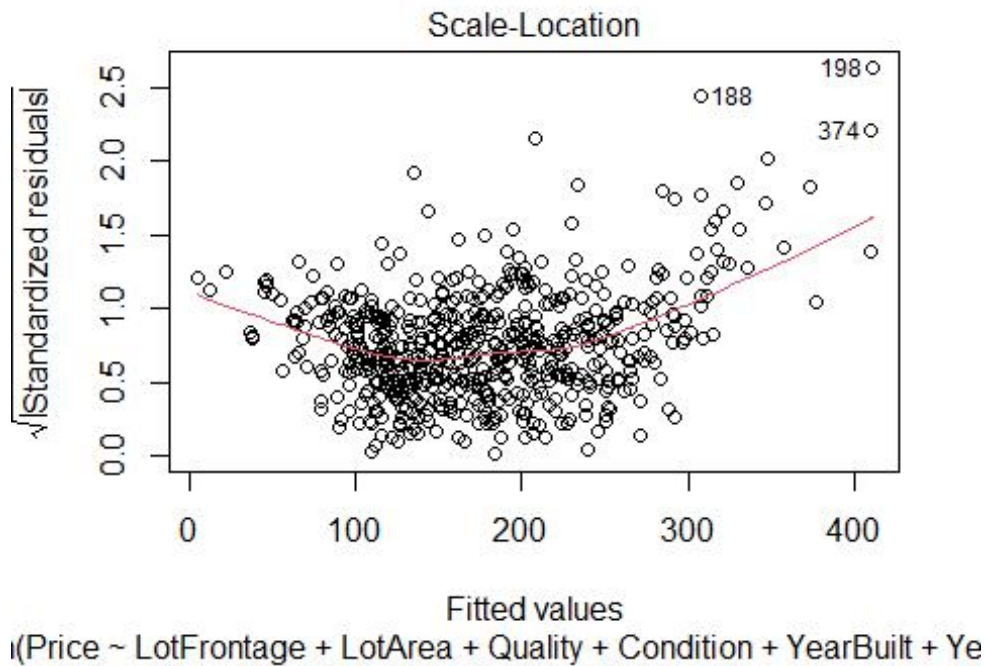
```
plot(backward_mod)
```



Fitted values
 l(Price ~ LotFrontage + LotArea + Quality + Condition + YearBuilt + Ye



Theoretical Quantiles
 l(Price ~ LotFrontage + LotArea + Quality + Condition + YearBuilt + Ye



(Price ~ LotFrontage + LotArea + Quality + Condition + YearBuilt + Ye

Linear conditions: - Linearity is questionable, since the residuals vs fitted plot is better described by a curve rather than a horizontal line. - Constant variance appears to be okay. - Normality of residuals is questionable, since the residuals have obvious skewness from the theoretical quantiles based on the qq plot.

Check whether any of the data cases are unusual with respect to studentized residuals:

```
head(sort(abs(rstudent(backward_mod)), decreasing=TRUE), 15)

##      198      188      374      204      62      572      70      268
## 7.257964 6.174013 5.007439 4.703305 4.150080 3.720407 3.443322 3.4200
68
##      581      202      537      535      386      228      126
## 3.377791 3.260223 3.182584 3.047404 2.968657 2.782562 2.757820
```

The first 12 houses listed above are outliers - they have studentized residuals larger than 3 or smaller than -3.

For adjustment, we decide not to do anything because every cook's distance is smaller than 0.5, so no data points have a drastic effect on the whole model.

Part 3: Find a "fancier model":

In addition to the quantitative predictors from Part 1, you may now consider models with:

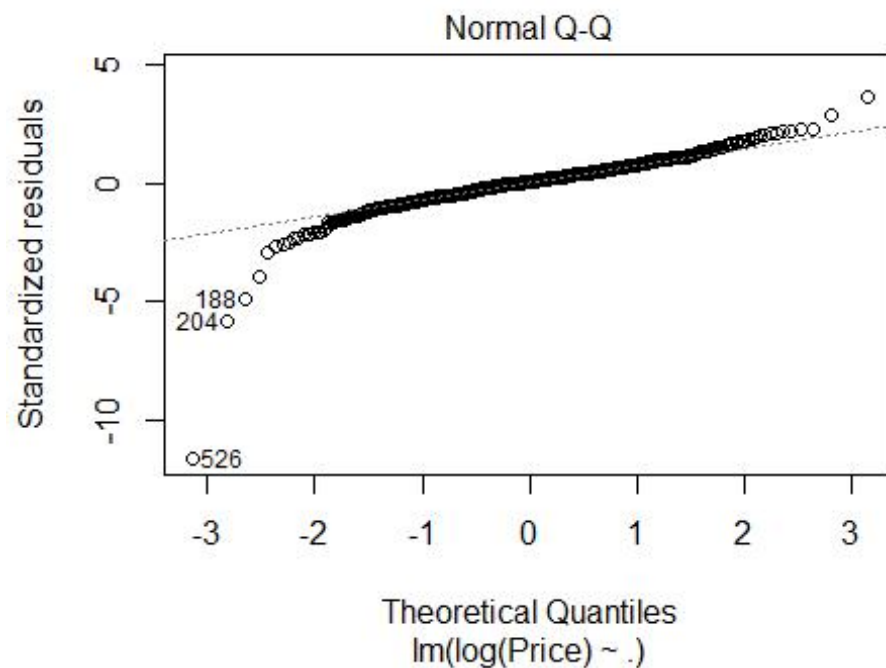
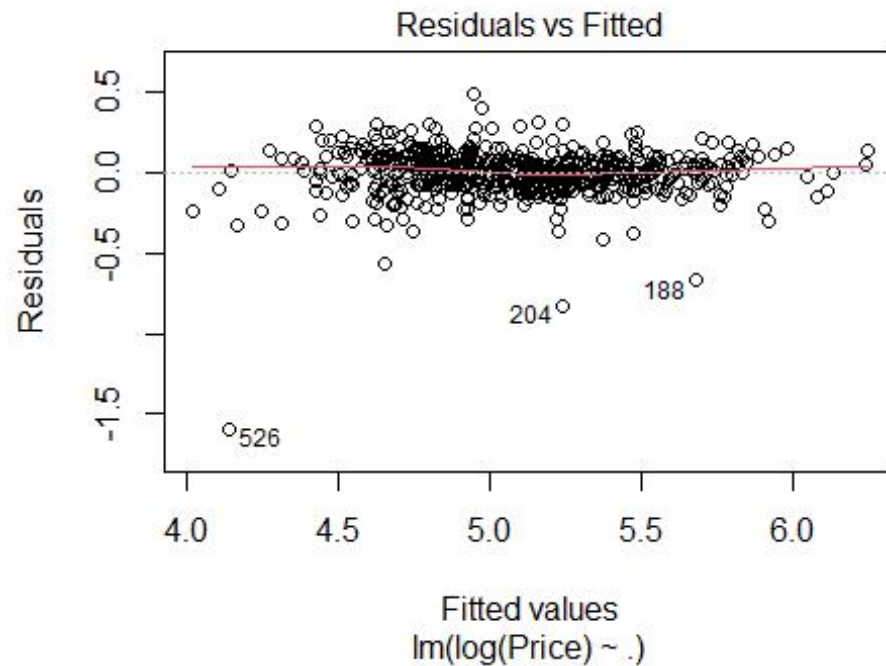
- Transformations of predictors. You can include functions of quantitative predictors. Probably best to use the I() notation so you don't need to create new columns when you run the predictions for the test data. For example: `lm(Price~LotArea+I(LotArea^2)+sqrt(LotArea)+log(LotArea),...`
- Transformations of the response. You might address curvature or skewness in residual plots by transforming the response prices with a function like `log(Price)`, `sqrt(Price)`, `Price^2`, etc.. These should generally not need the I() notation to make these adjustments.
- Combinations of variables. This might include for example creating a new variable which would count the total bathrooms in the house in a single predictor.

Do not haphazardly use transformation on predictors, but examine the relationships between the predictors and response to determine when a transformation would be warranted. Again use multiple model selection methods to determine a best model, but now with transformed variables are possible predictors in the model.

Discuss the process that you used to transform the predictors and/or response so that you could use this process in the future on a new data set.

Transformation: `log(Price)`

```
trans_mod = lm(log(Price)~., data = myData)
plot(trans_mod, 1:2)
```



Check linear conditions: - Linearity is much improved, since the residuals vs fitted plot almost has a perfect horizontal trend. - Constant variance appears to be okay. - Normality of residuals is much improved, since the residuals have less skewness from the theoretical quantiles based on the qq plot.

(1) backward elimination

```
Full_tr = lm(log(Price)~., data=myData)
MSE_tr = (summary(Full_tr)$sigma)^2
backward_mod_tr = step(Full_tr, scale=MSE_tr, trace=FALSE)
summary(backward_mod_tr)

##
## Call:
## lm(formula = log(Price) ~ LotFrontage + LotArea + Quality + Condition
+
##   YearBuilt + YearRemodel + BasementUnFinSF + BasementSF +
##   FirstSF + SecondSF + BasementFBath + Fireplaces + GarageCars +
##   EnclosedPorchSF + ScreenPorchSF, data = myData)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.57767 -0.06609  0.00754  0.07580  0.49306
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -6.469e+00  7.888e-01  -8.201 1.53e-15 ***
## LotFrontage     3.964e-04  1.940e-04   2.044 0.041439 *
## LotArea        2.321e-06  5.459e-07   4.251 2.48e-05 ***
## Quality        7.404e-02  7.086e-03  10.449 < 2e-16 ***
## Condition      6.106e-02  6.197e-03   9.853 < 2e-16 ***
## YearBuilt      3.754e-03  3.411e-04  11.005 < 2e-16 ***
## YearRemodel    1.314e-03  4.076e-04   3.223 0.001341 **
## BasementUnFinSF -7.132e-05  2.073e-05  -3.440 0.000623 ***
## BasementSF      1.477e-04  3.652e-05   4.045 5.92e-05 ***
## FirstSF        3.542e-04  3.753e-05   9.436 < 2e-16 ***
## SecondSF       2.781e-04  1.812e-05  15.346 < 2e-16 ***
## BasementFBath   2.295e-02  1.598e-02   1.437 0.151325
## Fireplaces     2.507e-02  1.146e-02   2.188 0.029053 *
## GarageCars     4.897e-02  1.076e-02   4.549 6.56e-06 ***
## EnclosedPorchSF 3.091e-04  1.176e-04   2.629 0.008792 **
## ScreenPorchSF  2.339e-04  9.396e-05   2.489 0.013073 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1438 on 584 degrees of freedom
## Multiple R-squared:  0.8746, Adjusted R-squared:  0.8713
## F-statistic: 271.4 on 15 and 584 DF,  p-value: < 2.2e-16
```

(2) forward selection

```
none_tr = lm(log(Price)~1, data=myData)
forward_mod_tr = step(none_tr, scope=list(upper=Full_tr), scale=MSE_tr,
  direction="forward", trace=FALSE)
summary(forward_mod_tr)
```

```
##
## Call:
## lm(formula = log(Price) ~ Quality + FirstSF + SecondSF + YearBuilt +
##   Condition + BasementFinSF + GarageCars + LotArea + EnclosedPorchS
F +
##   ScreenPorchSF + YearRemodel + BasementSF + Fireplaces + LotFronta
ge +
##   BasementUnFinSF, data = myData)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.58089 -0.06519  0.00795  0.07608  0.49476
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -6.575e+00  7.847e-01  -8.379 3.99e-16 ***
## Quality       7.353e-02  7.090e-03  10.370 < 2e-16 ***
## FirstSF      3.512e-04  3.753e-05   9.358 < 2e-16 ***
## SecondSF     2.776e-04  1.814e-05  15.302 < 2e-16 ***
## YearBuilt    3.747e-03  3.415e-04  10.973 < 2e-16 ***
## Condition    6.061e-02  6.214e-03   9.755 < 2e-16 ***
## BasementFinSF 3.361e-05  3.496e-05   0.961 0.336746
## GarageCars   4.781e-02  1.076e-02   4.443 1.06e-05 ***
## LotArea      2.348e-06  5.480e-07   4.285 2.13e-05 ***
## EnclosedPorchSF 3.170e-04  1.177e-04   2.693 0.007280 **
## ScreenPorchSF 2.327e-04  9.413e-05   2.472 0.013706 *
## YearRemodel  1.378e-03  4.053e-04   3.399 0.000722 ***
## BasementSF   1.399e-04  4.470e-05   3.131 0.001830 **
## Fireplaces   2.465e-02  1.147e-02   2.149 0.032076 *
## LotFrontage  4.142e-04  1.938e-04   2.138 0.032953 *
## BasementUnFinSF -6.264e-05  3.336e-05  -1.878 0.060907 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1439 on 584 degrees of freedom
## Multiple R-squared:  0.8743, Adjusted R-squared:  0.8711
## F-statistic: 270.8 on 15 and 584 DF, p-value: < 2.2e-16
```

(3) stepwise regression

```
stepwise_mod_tr = step(none_tr, scope=list(upper=Full_tr), scale=MSE_tr,
  trace=FALSE)
summary(stepwise_mod_tr)
```

```
##
## Call:
## lm(formula = log(Price) ~ Quality + FirstSF + SecondSF + YearBuilt +
##   Condition + GarageCars + LotArea + EnclosedPorchSF + ScreenPorchS
F +
##   YearRemodel + BasementSF + Fireplaces + LotFrontage + BasementUnF
inSF +
```



```
##      BasementFBath, data = myData)
##
## Residuals:
##      Min        1Q    Median        3Q        Max
## -1.57767 -0.06609  0.00754  0.07580  0.49306
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -6.469e+00  7.888e-01  -8.201 1.53e-15 ***
## Quality       7.404e-02  7.086e-03  10.449 < 2e-16 ***
## FirstSF      3.542e-04  3.753e-05   9.436 < 2e-16 ***
## SecondSF     2.781e-04  1.812e-05  15.346 < 2e-16 ***
## YearBuilt    3.754e-03  3.411e-04  11.005 < 2e-16 ***
## Condition    6.106e-02  6.197e-03   9.853 < 2e-16 ***
## GarageCars   4.897e-02  1.076e-02   4.549 6.56e-06 ***
## LotArea      2.321e-06  5.459e-07   4.251 2.48e-05 ***
## EnclosedPorchSF 3.091e-04  1.176e-04   2.629 0.008792 **
## ScreenPorchSF 2.339e-04  9.396e-05   2.489 0.013073 *
## YearRemodel  1.314e-03  4.076e-04   3.223 0.001341 **
## BasementSF   1.477e-04  3.652e-05   4.045 5.92e-05 ***
## Fireplaces   2.507e-02  1.146e-02   2.188 0.029053 *
## LotFrontage  3.964e-04  1.940e-04   2.044 0.041439 *
## BasementUnFinSF -7.132e-05  2.073e-05  -3.440 0.000623 ***
## BasementFBath 2.295e-02  1.598e-02   1.437 0.151325
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1438 on 584 degrees of freedom
## Multiple R-squared:  0.8746, Adjusted R-squared:  0.8713
## F-statistic: 271.4 on 15 and 584 DF,  p-value: < 2.2e-16
```

Model selection: Backward elimination: `lm(formula = log(Price) ~ LotFrontage + LotArea + Quality + Condition + YearBuilt + YearRemodel + BasementUnFinSF + BasementSF + FirstSF + SecondSF + BasementFBath + Fireplaces + GarageCars + EnclosedPorchSF + ScreenPorchSF, data = myData)` Forward selection: `lm(formula = log(Price) ~ Quality + FirstSF + SecondSF + YearBuilt + Condition + BasementFinSF + GarageCars + LotArea + EnclosedPorchSF + ScreenPorchSF + YearRemodel + BasementSF + Fireplaces + LotFrontage + BasementUnFinSF, data = myData)` Stepwise regression: `lm(formula = log(Price) ~ Quality + FirstSF + SecondSF + YearBuilt + Condition + GarageCars + LotArea + EnclosedPorchSF + ScreenPorchSF + YearRemodel + BasementSF + Fireplaces + LotFrontage + BasementUnFinSF + BasementFBath, data = myData)`

The only difference between the 3 models is that backward elimination and stepwise regression includes BasementFBath, while forward selection includes BasementFinSF. Since the model produced by backward elimination and stepwise regression has a slightly higher multiple R^2 than the other model, we are going to use this model as our fancier model.

Our fancier model is:

```
fancier_mod = backward_mod_tr
```

Part 4. Residual analysis for your fancier model

Repeat the residual analysis from Part 2 on your new model constructed in Part 3. A residual analysis was likely (hopefully) part of your process for determining your “fancier” model. That does not need to be repeated here as long as you clearly discuss your process.

We already did this analysis in part 3. Linear conditions are satisfied pretty well.

Part 5. Final model

Suppose that you are interested in a house Ames that has characteristics listed below and want to **find a 95% prediction interval for the price of this house**.

A 2 story 11 room home, built in 1983 and remodeled in 1999 on a 21540 sq. ft. lot with 400 feet of road frontage. Overall quality is good (7) and condition is average (5). The quality and condition of the exterior are both good (Gd) and it has a poured concrete foundation. There is an 757 sq. foot basement that has excellent height, but is completely unfinished and has no bath facilities. Heating comes from a gas air furnace that is in excellent condition and there is central air conditioning. The house has 2432 sq. ft. of living space above ground, 1485 on the first floor and 947 on the second, with 4 bedrooms, 2 full and one half baths, and 1 fireplace. The 2 car, built-in garage has 588 sq. ft. of space and is average (TA) for both quality and construction. The only porches or decks is a 384 sq. ft. open porch in the front.

#95% prediction interval for the price of an individual house

```
newHouse=data.frame(LotFrontage=400,  
                    LotArea=21540,  
                    Quality=7,  
                    Condition=5,  
                    YearBuilt=1983,  
                    YearRemodel=1999,  
                    BasementUnFinSF=757,  
                    BasementSF=757,  
                    FirstSF=1485,  
                    SecondSF=947,  
                    BasementFBath=0,  
                    Fireplaces=1,  
                    GarageCars=2,  
                    EnclosedPorchSF=0,  
                    ScreenPorchSF=0)
```

#prediction interval: just this one

```
predict.lm(fancier_mod, newHouse, interval="prediction")
```

```
##          fit          lwr          upr
## 1 5.602833 5.288915 5.916751
```

A 95% prediction interval for the log(price) of this house is [5.288915, 5.916751].

```
print(exp(5.288915))
```

```
## [1] 198.1283
```

```
print(exp(5.916751))
```

```
## [1] 371.2037
```

A 95% prediction interval for the Price of this house is [198.1283, 371.2037].