

Exercise 15 - Overfitting

Zoltan Kekecs

14 november 2019

Contents

1	Abstract	2
2	Data management and descriptive statistics	2
2.1	Loading packages	2
2.2	Load data about housing prices in King County, USA	2
2.3	Check the dataset	2
3	Overfitting	2
3.1	First rule of model selection:	2
3.2	Comparing model performance on the training set and the test set	3
3.3	Result-based models selection	7
3.4	Testing performance on the test set	8
3.5	BOTTOM LINE	8

1 Abstract

This exercise will demonstrate how overfitting works. The basic idea is that by providing too much flexibility for the model, we allow it to fit to the sample it is being trained on too well. However, this close fit to the sample will also fit some of the patterns that are only present in our sample (due to the sampling error), and are not present in the whole population.

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-2019

2 Data management and descriptive statistics

2.1 Loading packages

You will need to load the following packages for this exercise:

```
library(tidyverse) # for tidy format
```

2.2 Load data about housing prices in King County, USA

In this exercise we will predict the price of apartments and houses.

We use a dataset from Kaggle containing data about housing prices and variables that may be used to predict housing prices. This dataset contains house sale prices for King County, USA (Seattle and surrounding area) which includes Seattle. It includes homes sold between May 2014 and May 2015. More info about the dataset here: <https://www.kaggle.com/harlfoxem/housesalesprediction>

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_house = read_csv("https://bit.ly/2DpwK0r")
```

2.3 Check the dataset

You should always get familiar with the dataset you are using, and check for any inconsistencies that need to be corrected.

In the code below we convert the area metrics that are in square feet in the original dataset to square meters. We also specify that the variable `has_basement` is a factor.

```
data_house = data_house %>% mutate(sqm_living = sqft_living *  
  0.09290304, sqm_lot = sqft_lot * 0.09290304, sqm_above = sqft_above *  
  0.09290304, sqm_basement = sqft_basement * 0.09290304, sqm_living15 = sqft_living15 *  
  0.09290304, sqm_lot15 = sqft_lot15 * 0.09290304, has_basement = factor(has_basement))  
  
data_house %>% summary()
```

3 Overfitting

3.1 First rule of model selection:

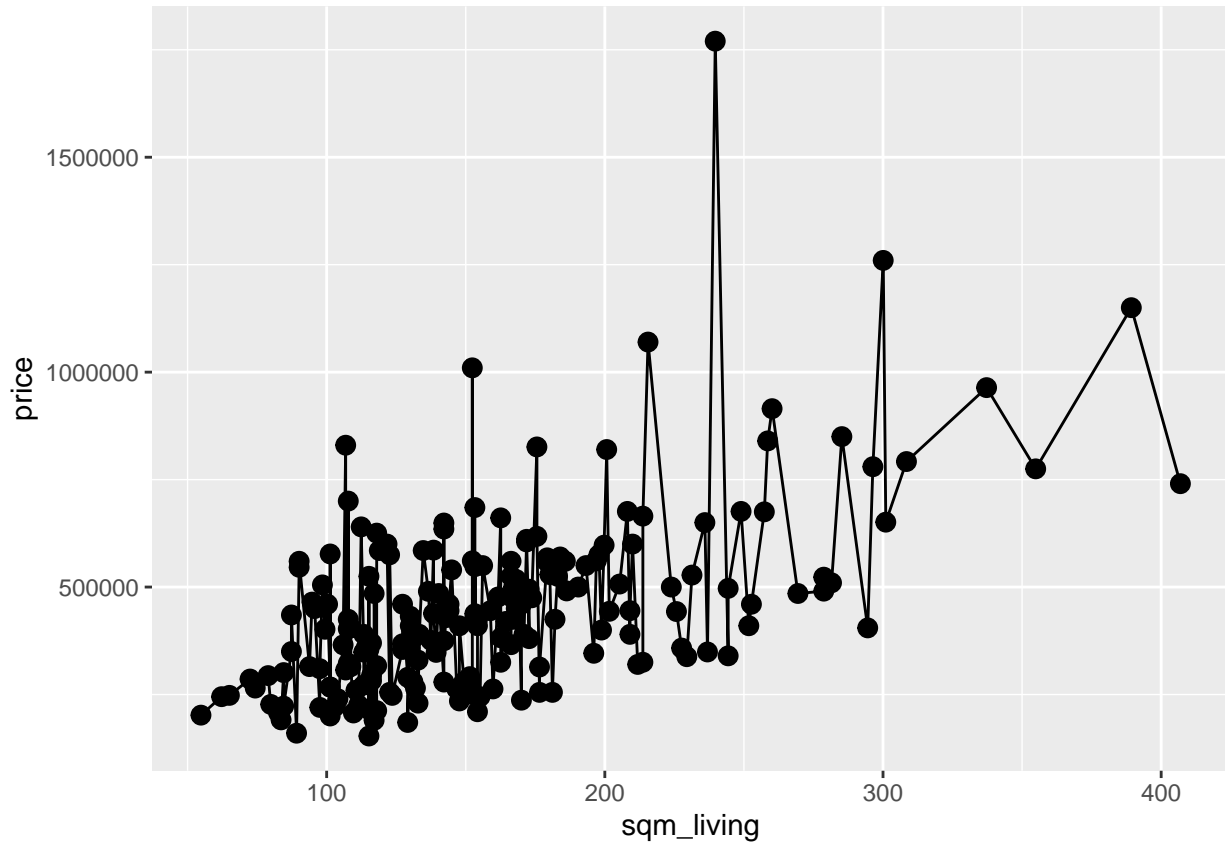
Always go with the model that is grounded in theory and prior research, result-driven model selection can lead to bad predictions on new datasets due to overfitting!

“Predicting” variability of the outcome **in your original data** is easy. If you fit a model that is flexible enough, you will get perfect fit on your initial data.

For example you can fit a line that would cover your data perfectly, reaching 100% model fit... to a dataset where you already knew the outcome.

See the example “prediction line” below.

```
data_house %>% ggplot() + aes(y = price, x = sqm_living) + geom_point(size = 3) +  
  geom_line()
```



However, when you try to apply the same model to new data, it will produce bad model fit. In most cases, worse, than a simple regression.

3.2 Comparing model performance on the training set and the test set

In this context, data on which the model was built is called the training set, and the new data where we test the true prediction efficiency of a model is called the test set. The test set can be truly newly collected data, or it can be a set aside portion of our old data which was not used to fit the model.

Linear regression is very inflexible, so it is less prone to overfitting. This is one of its advantages compared to more flexible prediction approaches.

In the next part of the exercise we will demonstrate that the more predictors you have, the higher your R^2 will be, even if the predictors have nothing to do with your outcome variable.

First, we will generate some random variables for demonstration purposes. These will be used as predictors in some of our models in this exercise. It is important to realize that these variables are randomly generated, and have no true relationship to the sales price of the apartments. Using these random numbers we can demonstrate well how people can be misled by good prediction performance of models containing many predictors.

```
rand_vars = as.data.frame(matrix(rnorm(mean = 0, sd = 1, n = 50 *
  nrow(data_house)), ncol = 50))
data_house_withrandomvars = cbind(data_house, rand_vars)
```

We create a new data object from the first half of the data ($N = 100$). We will use this to fit our models on. This is our training set. We set aside the other half of the dataset so that we will be able to test prediction performance on it later. This is called the test set.

```
training_set = data_house_withrandomvars[1:100, ] # training set, using half of the data
test_set = data_house_withrandomvars[101:200, ] # test set, the other half of the dataset
```

Now we will perform a hierarchical regression where first we fit our usual model predicting price with `sqm_living` and `grade` on the training set. Next, we fit a model containing `sqm_living` and `grade` and the 50 randomly generated variables that we just created.

(the names of the random variables are `V1`, `V2`, `V3`, ...)

```
mod_house_train <- lm(price ~ sqm_living + grade, data = training_set)
mod_house_rand_train <- lm(price ~ sqm_living + grade + V1 +
  V2 + V3 + V4 + V5 + V6 + V7 + V8 + V9 + V10 + V11 + V12 +
  V13 + V14 + V15 + V16 + V17 + V18 + V19 + V20 + V21 + V22 +
  V23 + V24 + V25 + V26 + V27 + V28 + V29 + V30 + V31 + V32 +
  V33 + V34 + V35 + V36 + V37 + V38 + V39 + V40 + V41 + V42 +
  V43 + V44 + V45 + V46 + V47 + V48 + V49 + V50, data = training_set)
```

Now we can compare the model performance. First, if we look at the normal R^2 indexes of the models or the RSS, we will find that the model using the random variables (`mod_house_rand_train`) was much better at predicting the training data. The error was smaller in this model, and the overall variance explained is bigger. You can even notice that some of the random predictors were identified as having significant added prediction value in this model, even though they are not supposed to be related to price at all, since we just created them randomly. This is because some of these variables are aligned with the outcome to some extent by random chance.

```
summary(mod_house_train)
```

```
##
## Call:
## lm(formula = price ~ sqm_living + grade, data = training_set)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -420718  -84545  -15651   81155  471385
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -318985.3   123358.6  -2.586  0.011201 *
## sqm_living    1245.5     353.4    3.524  0.000650 ***
## grade       77791.0    21506.0    3.617  0.000475 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 148300 on 97 degrees of freedom
## Multiple R-squared:  0.4765, Adjusted R-squared:  0.4657
## F-statistic: 44.15 on 2 and 97 DF,  p-value: 2.327e-14
```

```
summary(mod_house_rand_train)
```

```
##
## Call:
## lm(formula = price ~ sqm_living + grade + V1 + V2 + V3 + V4 +
##      V5 + V6 + V7 + V8 + V9 + V10 + V11 + V12 + V13 + V14 + V15 +
##      V16 + V17 + V18 + V19 + V20 + V21 + V22 + V23 + V24 + V25 +
##      V26 + V27 + V28 + V29 + V30 + V31 + V32 + V33 + V34 + V35 +
##      V36 + V37 + V38 + V39 + V40 + V41 + V42 + V43 + V44 + V45 +
##      V46 + V47 + V48 + V49 + V50, data = training_set)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -259872  -58982  -10714   64579  197501
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -288487.2   145955.9  -1.977  0.053979 .
## sqm_living    1885.2     418.5    4.505 4.39e-05 ***
## grade        57736.0    26220.2    2.202 0.032611 *
## V1             3389.0    18344.9    0.185 0.854230
## V2             3462.1    18752.4    0.185 0.854321
## V3            26589.4    16540.5    1.608 0.114635
## V4            -2416.6    16876.3   -0.143 0.886749
## V5            23766.6    16911.7    1.405 0.166499
## V6            26254.5    18147.7    1.447 0.154614
## V7           -31103.7    16182.8   -1.922 0.060677 .
## V8            -7519.0    18657.0   -0.403 0.688768
## V9           -27183.4    19278.9   -1.410 0.165121
## V10           -2844.9    16081.1   -0.177 0.860339
## V11           -7237.1    15606.7   -0.464 0.644993
## V12          -46241.5    17559.6   -2.633 0.011406 *
## V13             6588.9    20087.0    0.328 0.744356
## V14           -8576.9    17843.2   -0.481 0.632972
## V15          -22675.6    18042.7   -1.257 0.215045
## V16          -21359.0    17252.8   -1.238 0.221864
## V17            33483.2    20622.6    1.624 0.111147
## V18            46476.5    19674.9    2.362 0.022359 *
## V19          -12154.3    19066.8   -0.637 0.526921
## V20            24378.7    17657.1    1.381 0.173911
## V21            13263.3    17400.9    0.762 0.449737
## V22           -8343.7    18310.1   -0.456 0.650712
## V23            44027.1    19077.8    2.308 0.025462 *
## V24          -27467.4    18057.1   -1.521 0.134924
## V25          -32102.8    19348.3   -1.659 0.103734
## V26           -1563.8    18577.3   -0.084 0.933273
## V27            65446.5    17434.9    3.754 0.000479 ***
## V28          -11612.0    15849.6   -0.733 0.467420
## V29          -28379.9    17684.6   -1.605 0.115242
## V30           -3504.5    16606.3   -0.211 0.833775
## V31          -14930.8    19094.8   -0.782 0.438178
## V32          -10822.5    17653.5   -0.613 0.542795
## V33             71105.4    17657.2    4.027 0.000205 ***
## V34          -19880.1    19548.2   -1.017 0.314369
## V35          -13351.2    16544.0   -0.807 0.423727
## V36          -18523.2    19513.7   -0.949 0.347354
```

```
## V37          4873.7      19477.8    0.250 0.803510
## V38         -23227.4      21987.3   -1.056 0.296186
## V39         -28143.9      21123.2   -1.332 0.189163
## V40         -50622.5      16803.8   -3.013 0.004163 **
## V41          2477.4      15688.1    0.158 0.875199
## V42          11568.7      15165.2    0.763 0.449368
## V43          41992.2      18390.2    2.283 0.026972 *
## V44         -13044.4      18633.2   -0.700 0.487338
## V45          27414.9      17423.1    1.573 0.122315
## V46         -30879.6      15816.7   -1.952 0.056869 .
## V47         -10480.0      17942.3   -0.584 0.561951
## V48          27065.0      18249.2    1.483 0.144730
## V49         -22737.1      18347.7   -1.239 0.221412
## V50         -23236.6      16420.1   -1.415 0.163622
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 122900 on 47 degrees of freedom
## Multiple R-squared:  0.8257, Adjusted R-squared:  0.6329
## F-statistic: 4.283 on 52 and 47 DF,  p-value: 6.906e-07

pred_train <- predict(mod_house_train)
pred_train_rand <- predict(mod_house_rand_train)
RSS_train = sum((training_set[, "price"] - pred_train)^2)
RSS_train_rand = sum((training_set[, "price"] - pred_train_rand)^2)
RSS_train

## [1] 2.133707e+12
RSS_train_rand

## [1] 710290417396

That is why we need to use model fit indexes that are more sensitive to the number of variables we included as predictors, to account for the likelihood that some variables will show a correlation by chance. Such as adjusted R2, or the AIC. The anova() test is also sensitive to the number of predictors in the models, so it is not easy to fool by adding a bunch of random data as predictors. Better yet, the AIC indicates that smaller model is significantly better than the more complicated model.

summary(mod_house_train)$adj.r.squared

## [1] 0.4657206

summary(mod_house_rand_train)$adj.r.squared

## [1] 0.6329346

AIC(mod_house_train)

## [1] 2670.159

AIC(mod_house_rand_train)

## [1] 2660.165

anova(mod_house_train, mod_house_rand_train)

## Analysis of Variance Table
##
## Model 1: price ~ sqm_living + grade
```

```
## Model 2: price ~ sqm_living + grade + V1 + V2 + V3 + V4 + V5 + V6 + V7 +
##      V8 + V9 + V10 + V11 + V12 + V13 + V14 + V15 + V16 + V17 +
##      V18 + V19 + V20 + V21 + V22 + V23 + V24 + V25 + V26 + V27 +
##      V28 + V29 + V30 + V31 + V32 + V33 + V34 + V35 + V36 + V37 +
##      V38 + V39 + V40 + V41 + V42 + V43 + V44 + V45 + V46 + V47 +
##      V48 + V49 + V50
##   Res.Df      RSS Df Sum of Sq      F Pr(>F)
## 1      97 2.1337e+12
## 2      47 7.1029e+11 50 1.4234e+12 1.8838 0.01519 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

3.3 Result-based models selection

(Result-based models selection is only shown here with demonstration purposes, to show how it can mislead researchers. Whenever possible, stay away from using such approaches, and rely on theoretical considerations and previous data when building models.)

After seeing the performance of `mod_house_rand_train`, and not knowing that it contains random variables, one might be tempted to build a model with only the predictors that were identified as having a significant added predictive value, to improve the model fit indices (e.g. adjusted R^2 or AIC). And that would achieve exactly that: it would result in the virtual improvement of the indexes, but not the actual prediction efficiency, so the better indexes would be just an illusion resulting from the fact that we “concealed” from the statistical tests that we have tried to use a lot of predictors in a previous model.

Excluding variables that seem “useless” based on the results will blind the otherwise sensitive measures of model fit. This is what happens when using automated model selection procedures, such as backward regression.

In the example below we use backward regression. This method first fits a complete model with all of the specified predictors, and then determines which predictor has the smallest amount of unique added explanatory value to the model, and excludes it from the list of predictors, refitting the model without this predictor. This procedure is iterated until there is no more predictor that can be excluded without significantly reducing model fit, at which point the process stops.

```
mod_back_train = step(mod_house_rand_train, direction = "backward")
```

The final model with the reduced number of predictors will have much better model fit indexes than the original complex model, because the less useful variables were excluded, and only the most influential ones were retained, resulting in a small and powerful model. Or at least this is what the numbers would suggest us on the training set.

Lets compare the prediction performance of the final model returned by backward regression (`mod_back_train`) with the model only containing our good old predictors, `sqm_living` and `grade` (`mod_house_train`) on the training set.

```
anova(mod_house_train, mod_back_train)
```

```
## Analysis of Variance Table
##
## Model 1: price ~ sqm_living + grade
## Model 2: price ~ sqm_living + grade + V3 + V7 + V12 + V15 + V17 + V18 +
##      V20 + V23 + V25 + V27 + V29 + V33 + V39 + V40 + V43 + V45 +
##      V46 + V48 + V50
##   Res.Df      RSS Df Sum of Sq      F    Pr(>F)
## 1      97 2.1337e+12
## 2      78 9.4654e+11 19 1.1872e+12 5.1489 1.031e-07 ***
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
summary(mod_house_train)$adj.r.squared
```

```
## [1] 0.4657206
```

```
summary(mod_back_train)$adj.r.squared
```

```
## [1] 0.7052536
```

```
AIC(mod_house_train)
```

```
## [1] 2670.159
```

```
AIC(mod_back_train)
```

```
## [1] 2626.878
```

All of the above model comparison methods indicate that the backward regression model (`mod_back_train`) performs better. We know that this model can't be too much better than the smaller model, since it only contains a number of randomly generated variables in addition to the two predictors in the smaller model. So if we would only rely on these numbers, we would be fooled to think that the backward regression model is better.

3.4 Testing performance on the test set

A surefire way of determining actual model performance is to test it on new data, data that was not used in the “training” of the model. Here, we use the set aside test set to do this.

Note that we did not re-fit the models on the test set, we use the models fitted on the training set to make our predictions using the `predict()` function on the `test_set`!!!

```
# calculate predicted values
```

```
pred_test <- predict(mod_house_train, test_set)
```

```
pred_test_back <- predict(mod_back_train, test_set)
```

```
# now we calculate the sum of squared residuals
```

```
RSS_test = sum((test_set[, "price"] - pred_test)^2)
```

```
RSS_test_back = sum((test_set[, "price"] - pred_test_back)^2)
```

```
RSS_test
```

```
## [1] 3.639381e+12
```

```
RSS_test_back
```

```
## [1] 5.528998e+12
```

This test reveals that the backward regression model has more error than the model only using `sqm_living` and `grade`.

3.5 BOTTOM LINE

1. Model selection should be done pre-analysis, based on theory, previous results from the literature, or conventions on the field. Post-hoc result-driven predictor selection can lead to overfitting.
2. The only good test of a model's true prediction performance is to test the accuracy of its predictions on new data (or a set-aside test set using the above described training-set test-set approach)