

MA2822:
ADVANCED STATISTICS
ECOLE CENTRALE PARIS

Regression Analysis of Used Car Prices

Professor
Christine KERIBIN

Authors
Erik HENDEY BRÖTE
Patrick VON PLATEN

March 2, 2017

Contents

1	Introduction	2
2	Charectaristics of and Correlations in the Data	3
3	Study of price as a function of the prescene of ABS-technology	5
4	Study of Price as a Function of Mileage	6
4.1	Simple Regression	6
4.2	Inference from Simple Model	6
4.3	Determination of Mystery Variable <i>kop1</i>	7
5	Study of Price as a Function of all variables	8
5.1	Preliminary selection of different models	8
5.2	Validation of the models	9
6	Conclusion	10

1 Introduction

In this project we study the price of used cars based on a number of relevant characteristics. Using regression analysis tools implemented in the open source language *R* (which is widely used for statistical analysis), we aim to model the sales price as well as detect correlations between factors. The dataset includes 172 observations of cars with a 12 characteristics including the following: *price*, in thousands of Euros; *age*, in months; *km*, the total mileage in kilometers; *TIA*, the number of months until the next vehicle inspection. The data also includes dummy variables (which are equal to 1 or 0 depending on whether the factor is true or false) such as *ABS* and *SunRoof* representing the presence of ABS-technology or a sunroof respectively.

2 Charectaristics of and Correlations in the Data

Before procededing it is a good idea to get a general idea of the correlations between the variables. Figure 3 is a scatter plot showing the dependencies between the *price*, *age*, *km*, and *TIA*. It shows a relativety strong postive correlation between *price* and *age* and some negative correlation between *age* and *km* as well as postitive correlation between *price* and *km*. These relationships are in accord with our existing knowledge of the used car market. *TIA* seems to explain the other variables in this model poorly.

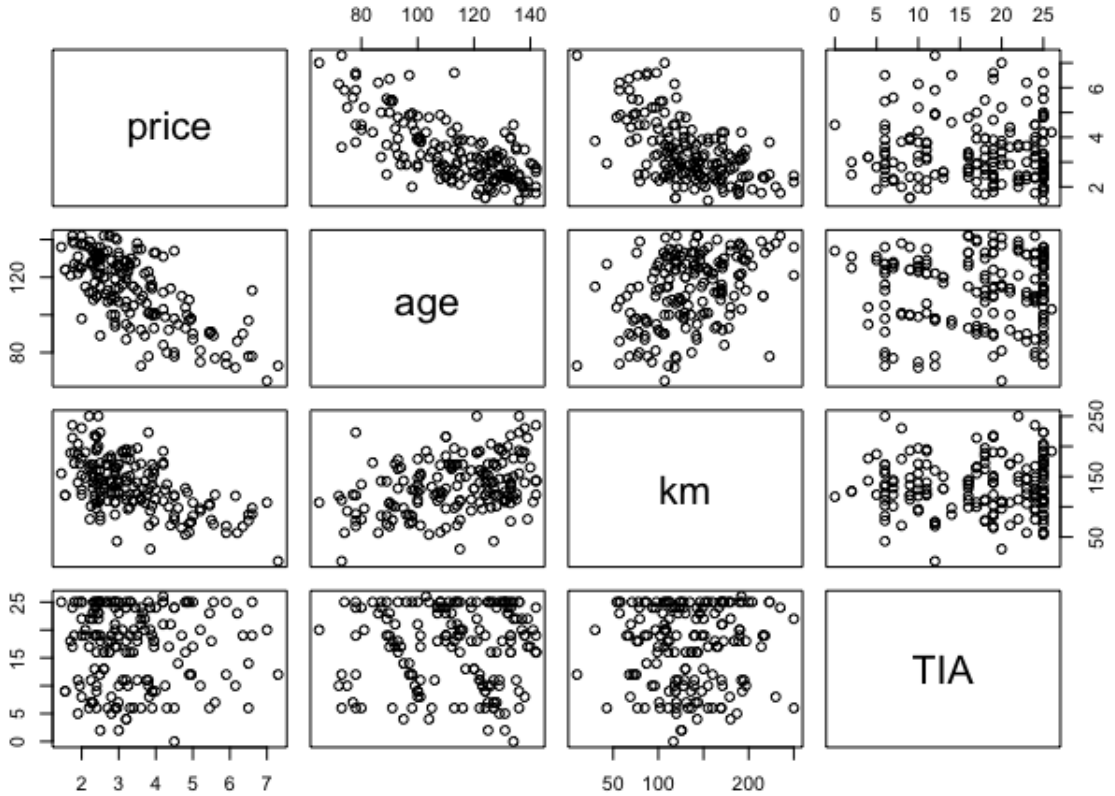


Figure 1: Scatter plot of the variables *price*, *age*, *km*, and *TIA*.

In addition to this we can use boxplots to analyze whether the presence of ABS technology or a sunroof, represented here by the dummy variables *ABS* and *SunRoof*, have an effect on price. Figure 2 displays this analysis for each variable. To the left the boxplots for vehicles with and without ABS show a few perhaps expected charectaristics of the dataset. The median price of a car with ABS is significantly higher than their counterparties (compare 3 200 EUR to 2 625 EUR). Also, the variance of the price of cars having an ABS-technology is much lower than for those that without it. Since 50% of all “ABS” cars have a price that ranges between 2.6 and 4 with pretty much no “ABS” car having a price higher than 6, where as 50% of

all “Non-ABS” cars have a price that ranges between 2.5 and 4.5 and a maximum of 7. It can be concluded that the price is somewhat correlated to the dummy variable *ABS*, but not very strongly.

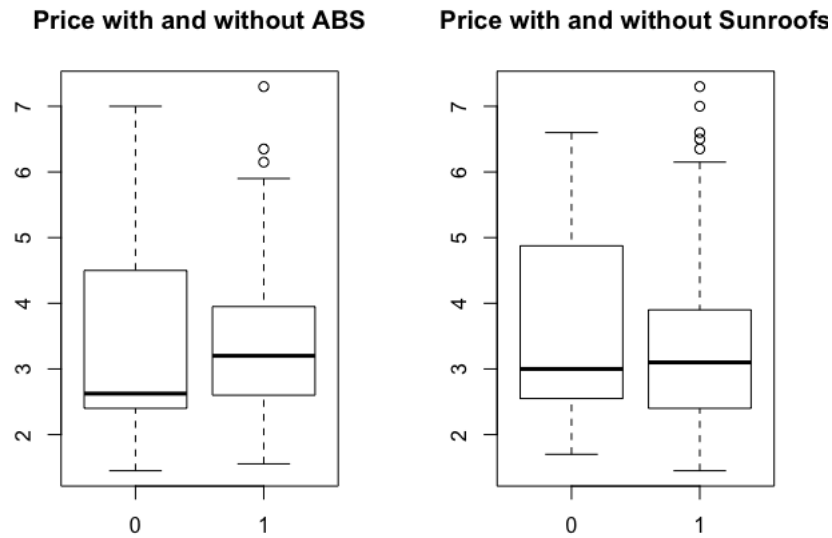


Figure 2: Pairs of boxplots of price of vehicles with and without ABS (left) and sunroofs (right).

A similar analysis can be done for the presence of a sunroof. The second set of boxplots in Figure 2 show values for vehicles with and without them. The medians are quite similar for “Non-OpenRoof” and “OpenRoof” cars. The “box range” around the median is smaller for cars with sunroofs, but apart from that they are no visible differences. The correlation between price and OpenRoof seems to be rather small; further analysis is needed to confirm this observed relationship.

3 Study of price as a function of the prescene of ABS-technology

As seen in Section 2 cars with ABS seem to generally have a higher price than cars without this technology. We will further analyze how well *price* can be explained by the *ABS*.

To better understand how the *R* software works we created two simple regression models with all variables provided with *ABS* as qualitative and quantitative variable. In implementation qualitative means that the variable is given a numerical value, 1 representing the presence of ABS or 0 otherwise. The quantitative variable is a logical data type meaning it can be either *TRUE* or *FALSE*. No difference can be seen in the *summary* function (which provides values and different significance metrics for each coefficient) between those two models. It may be pertinent to remember that the logical variable cannot be calculated upon, e.g. we cannot calculate its mean value. This should not make a difference in the final model as regression requires that the logical variable be implemented as a dummy variable (i.e. qualitatively) anyways.

The simplest of regressions can be a good way to analyze the explanatory power of *ABS*. We created a regression model of *price* as only a function of *ABS* corresponding to the following model.

$$price = \beta_0 + \beta_1 \times ABS + \epsilon \quad (1)$$

Here β_0 is the intercept and β_1 is the slope of the equation. ϵ is the residual term. Regression on our dataset resulted in the following estimations. $\hat{\beta}_0 = 3.45363$ and $\hat{\beta}_1 = -0.08321$. It is clear that the variable *ABS* does not account for any of the variance of the price, since its coefficient of determination, R^2 is very low (0.00097). That means that the model does not do a very good job predicting the price. This fact is also reflected by the p-value for *ABS*, $p = 0.686$. the p-value is the probability of obtaining a result as more extreme given the null hypothesis (which in this case is that the coefficient $\beta_1 = 0$.) For these reasons it is clear that *ABS* on its own explains the price very poorly. Arguably, the intercept alone would be a better model as the model would just create a nearly constant line.

4 Study of Price as a Function of Mileage

4.1 Simple Regression

A simple regression of *price* on mileage represented by the variable *km* shows quite different results. As in Section 3, we begin by performing a simple regression using the equation

$$price = \beta_0 + \beta_1 \times km + \epsilon \quad (2)$$

where β_0 is the intercept, β_1 is the slope and ϵ is the residual term. Applying this to our data results in β_0 and β_1 values of 5.559202 and -0.016051. Using the same analysis and reasoning as in 3, *km* is obviously a better parameter to describe the price of a car. The coefficient of determination R^2 is much higher (0.3317) and the p-value is much lower ($< 2e^{16}$), suggesting that *km* should be in model explaining the price.

4.2 Inference from Simple Model

We will now using the simple model from the previous section create confidence intervals of prices for cars with difference milages. This implemented in *R* using the *predict* function. For a cars with a milages of 50 000 km and 135 000 km a confidence intervals at the default significance level 95% are shown in Table 1.

Mileage	fit	lower	upper
50 000 km	4.756639	4.426392	5.086886
136 000 km	3.392284	3.238505	3.546063

Table 1: Confidence intervals at 95% confidence level for cars with milages of 50 000 km and 135 000 km

We see that the confidence intervall for a car that has 135 000 km mileage is much smaller than the other at the same sig. level. This implies that the variance of the price for cars having 135 000 km is lower that can be justified by the fact that the price of cars already having a high usage (in terms of km) are naturally lower and don't vary as much. Whereas for cars having lower usage (50 000 km) the price can vary much more since other factors play a more important role in determining its price (car type, car performance.) Also, it should be noted that 135 000 km is the mean of all of the data. Since the sample follows a t-distribution, it is normal that there is more data around 135 000 km (the distribution is denser). Therefore the 95% confidence intervall has a shorter lenght than the one at 50 000 km.

4.3 Determination of Mystery Variable *kop1*

The data set given includes the variable *kop1* given without explanation of what it characterizes. In fact, it reduces to the variable *km* as it is a centered and reduced version of it.

$$kop1 = \frac{km - \text{mean}(km)}{SD(km)} \quad (3)$$

In fact, simple regression models based on *kop1* and *km* result in the exact same model. It can be shown that the models are logically the same since R^2 and other important functions are equal. Algebraically this can be shown through the following reasoning. The regression model is given by the following equation:

$$price = \gamma_0 + \gamma_1 \times kop1 + e \quad (4)$$

Where γ_i are the coefficients and e is the residual. Using relationship 3 the regression model reduces to

$$price = \gamma_0 + \gamma_1 \times (km + \text{mean}(km)) + e = (\gamma_0 + \gamma_1 \times \text{mean}(km)) + \gamma_1 \times km + e \quad (5)$$

Given that $\text{mean}(km)$ is a constant, we note that this is equivalent with the regression model 2 as $\beta_0 = (\gamma_0 + \gamma_1 \times \text{mean}(km))$, $\beta_1 = \gamma_1$ and the residuals $\epsilon = e$.

It is obvious that centering and reducing shouldn't change anything. The regression looks at how changes in the X value affect the Y value

5 Study of Price as a Function of all variables

In this section, we will first calculate a linear regression model that takes all variables of the dataset into account in order to predict the price of the used car. In the following, we will use different methods to decide on which variables will be used for our “final model”. Therefore, different factors will be taken into consideration, such as

- the “*R-squared*” value showing how well the predictions fit the actual data and detected *high bias*,
- factors making sure the data is not *overfitting*, such as *BIC*, *AIC* and *adjusted “R-squared”*
- and factors using cross-validation methods in order to check whether the model is useful for data the model was not “trained on” (*BIC*, *AIC*)

5.1 Preliminary selection of different models

The first model we want to test later is simply the model that uses all variables in order to predict the price:

$$M_1 = lm(price \sim .) \quad (6)$$

The second model is obtained by looking at the p-values of every input variable using:

Listing 1: `summary()` in R

```
1 summary(Model1)
```

It can be seen that there are three variables that have a significant lower p-value than the other variables, being the *age*, *km* and *ageop2*. Therefore our second model is as follows:

$$M_2 = lm(price \sim age + km + ageop2) \quad (7)$$

Our third model is achieved by using `stepFunction`: “`stepAIC`”. This function calculates the *AIC* value of the model at every step and removes all the variables, so that the remaining model has a lower *AIC* and therefore is a “better fit”. Using this method, we get the following model:

$$M_3 = lm(price \sim age + km + ageop2 + kop2) \quad (8)$$

Our fourth and last model is achieved by using a model of preselected variables and executing the `stepAIC` function on this model whereas the “`stepAIC`” function will build on a preselected choice of variables. It is in some way a mixture of the methods used to obtain `model2` and `model3`. We preselect the variables *age*, *km*

since they show by far the highest p-values (see listing 1). Applying the “stepAIC” function, we get:

$$M_4 = lm(price \sim age + km + ageop2 + kop2) \quad (9)$$

This is the same model as model 3 and can therefore be discarded.

5.2 Validation of the models

First, we want to calculate the different “*R-squared*” values to select models having low biases when predicting the data. For all three models we get similiar results

	M_1	M_2	M_3
R^2	0.65833	0.64338	0.65581

The closer “*R-squared*” is to 1, the better it is. Obviously, the first model has the one that is closest to one since it uses all the variables for its model. We can see though that all “*R-squared*” are very close to each other and can therefore not yet determine the best model. We need more analysis.

To get in-depth analysis, we will use the “CV” function of the package “forecast”. This function calculates the *BIC*, *AIC*, *AICc* and *adjusted “R-squared”* for every model and returns a so-called “PRESS” value that takes into consideration all the factors mention above in the first part of the section 5 (especially the second and third factor). The comparable lowest “PRESS” value stands for the best model. The result of the “CV” function gave us:

```
> CV(Mfirst)
      CV      AIC      AICc      BIC      AdjR2
0.5970033 -88.0296131 -86.3796131 -53.4071739  0.6393558
> CV(Msecond)
      CV      AIC      AICc      BIC      AdjR2
0.5768241 -92.6608708 -92.2994250 -76.9233984  0.6370135
> CV(Mthird)
      CV      AIC      AICc      BIC      AdjR2
0.5615731 -96.7643674 -96.2552765 -77.8794006  0.6475706
```

Figure 3: *PRESS* analysis of the models.

As we can see, we get the lowest CV for the third model.

6 Conclusion

Having taking into consideration model evaluation factors, such as

- *R-squared*
- *R-squared adjusted*
- *BIC*
- *AIC*
- *AICc*
- and others

we chose **model3** as our model that will best predict the used car prices.

Last, but not least, we want to say that even if this project is rather “small”, it gave us a great example of how to structually approach a modelling problem and how to arrive at a result that is logically well founded.