

Statistical Data Analysis – solved problems

Goals: The text provides a pool of solved problems for labs in the course on Statistical Data Analysis. The exercises help to deepen knowledge gained in parallel Rmd files. At the same time, they serve as illustrative examples of future exam questions.

1 Linear and non-linear regression

Problem 1. (10 p) You built a linear model that predicts the median value of owner-occupied homes in \$1000's in a certain town (*medv*). The model works with the only independent variable (*lstat*) that captures the percentage of population with lower (economical) status in the given town. The model was built from a training set based on 506 towns and is this:

```
lm(formula = medv ~ lstat, data = Boston)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	34.55384	0.56263	61.41	<2e-16	***
<i>lstat</i>	-0.95005	0.03873	-24.53	<2e-16	***

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

Residual standard error: 6.216 on 504 degrees of freedom

Multiple R-squared: 0.5441, Adjusted R-squared: 0.5432

F-statistic: 601.6 on 1 and 504 DF, p-value: < 2.2e-16

- (a) (2 p) Verbally describe the relationship between *lstat* and *medv*. Decide whether *lstat* affects *medv* and quantify how. Is it a statistically significant relationship? Why?

The model says that with each additional percentage of people with lower status the median value of homes decreases on average by 950\$. This relationship is statistically significant, based on the F-statistic as well as the *lstat*'s t-value we can reject the null hypothesis that there is no relationship between *lstat* and *medv*.

- (b) (1 p) How do you understand the meaning of Intercept? Is the value of this coefficient a reliable figure to be interpreted literally? Explain.

The value of Intercept says that the average median value of homes in a town with 0 percentage of people with lower status is around 34,553\$. This value looks reasonable, however, its true reliability

depends on how far the model extrapolates (do we have any towns that at least approach no lower status representation in our training set?) and how far the model meets the linear regression assumptions (is the relationship between these variables truly linear?).

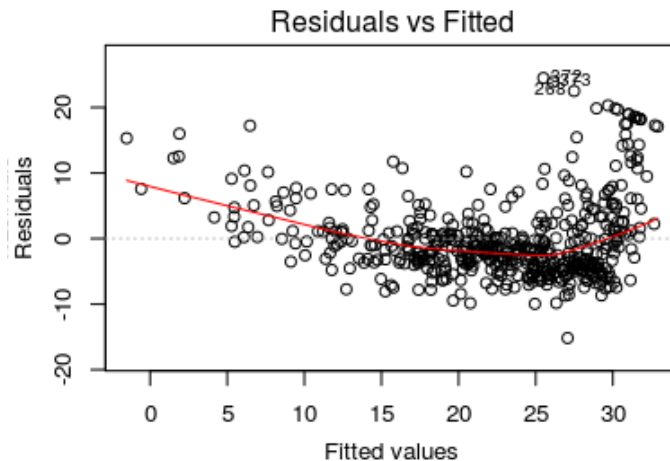
- (c) (1 p) How much do we improve our median value forecast compared to the simple average forecast that ignores the knowledge of lstat? In other words, how much the knowledge of lstat helps?

The value of R-squared shows that we will reduce the variance of the median home value estimates by about half.

- (d) (2 p) Calculate/estimate 95% confidence interval for β_{lstat} . What is this interval good for?

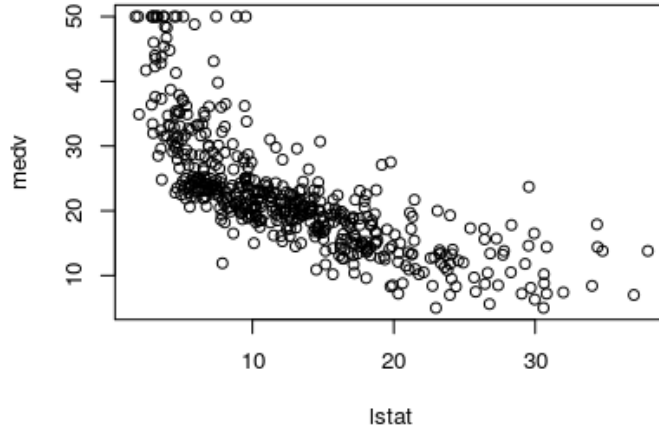
A rough estimate could be $[-0.95005 - 2 * 0.03873, -0.95005 + 2 * 0.03873] = [-1.02751, -0.87259]$. A more precise estimate puts $|t_{\alpha/2, m-2}| = |t_{0.025, 504}|$ instead of 2 into the formula above, however, the value 1.964682 is close to the rough estimate. This confidence interval has about 95% chance to contain the true value of β_{lstat} . This interval helps us to assume on the strength of relationship between lstat and medv, the interval does not contain 0, the relationship could be considered significant.

- (e) (2 p) Look at the model residual plot in the figure below (it plots differences between the actual and predicted values of the dependent variable). What conclusions can be drawn from the figure?

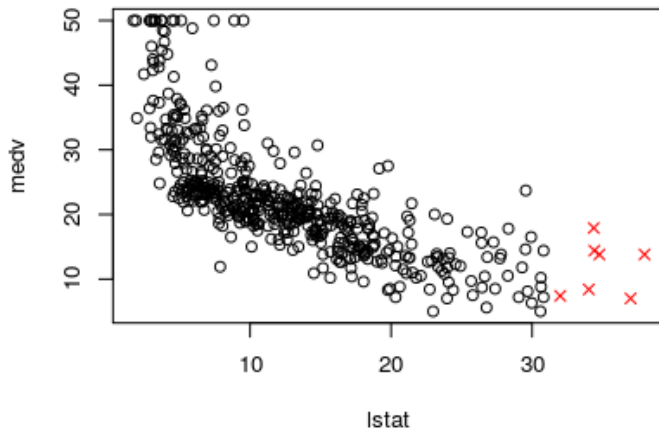


The plot shows that the assumption of linearity has not been met. The residuals should follow the normal distribution for all the values of lstat, they are heavily skewed in the plot. We should conclude that the relationship is non-linear and introduce new terms into the regression formula (lstat² is a good idea to start with).

- (f) (2 p) Explain the concept of influential observations. Denote a couple of influential points in the scatter plot below and explain how would you find them.



An influential observation is an observation whose deletion from the dataset would noticeably change the model parameter estimates. It could either be an outlier (a data point that differs significantly from other observations) or a high-leverage point (an observation made at extreme values of independent variables). The most influential observations can be found in the figure below.



Problem 2. (10 p) You are a mechanical locksmith and you are trying to find out how the shaft machining error is related to the machine tool parameter setting. You have compiled a multivariate linear model. The model expresses the relationship between the production error (the difference between the ideal shaft diameter and the actual shaft diameter, *ProdError*) and the setting of ten different continuous machine parameters (*P1-P10*). Below is the output you received:

```
summary(lm(ProdError ~ P1+P2+P3+P4+P5+P6+P7+P8+P9+P10), data=d)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.05270	0.09576	-0.550	0.5835
X1	0.01298	0.08924	0.145	0.8847

X2	0.01596	0.10939	0.146	0.8843
X3	-0.02865	0.09079	-0.316	0.7531
X4	0.04611	0.09548	0.483	0.6303
X5	0.14151	0.09343	1.515	0.1334
X6	-0.02375	0.10277	-0.231	0.8178
X7	0.25522	0.10516	2.427	0.0172 *
X8	0.06672	0.08972	0.744	0.4590
X9	0.09949	0.10171	0.978	0.3306
X10	-0.04003	0.09317	-0.430	0.6685

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

Residual standard error: 0.9039 on 89 degrees of freedom

Multiple R-squared: 0.1145, Adjusted R-squared: 0.01502

F-statistic: 1.151 on 10 and 89 DF, p-value: 0.3346

- (a) (2 p) Decide whether at least one of the machine parameters (independent variables) is useful for estimating a manufacturing error (ProdError). In other words, formally decide whether you can decline $H_0 : \beta_1 = \beta_2 = \dots = \beta_{10} = 0$. Justify correctly.

The reasoning should be based on the F-statistic and its corresponding p-value. The null hypothesis cannot be rejected, the model does not seem to be useful. The reasoning that stems from the statistics reached for the individual variables could be misleading due to multiple comparisons. For 10 variables, truly valid H_0 and $\alpha = 0.05$, there is only $0.95^{10} = 0.6$ probability that there will be no type I error in the individual coefficient tests, 40% of trials will find at least one falsely significant coefficient.

- (b) (2 p) Let us compare the full model constructed above with the intercept model and with the model that employs only the variable P7 identified as the most relevant. Let us compare them with F-test through an ANOVA run. Interpret the ANOVA table below.

```
lm.const<-lm(ProdError ~ 1,data=d) # the intercept model
lm.sel<-lm(ProdError ~ P7,data=d) # the P7 model
anova(lm.const,lm.sel,lm.all)
Analysis of Variance Table
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	99	82.114				
2	98	76.076	1	6.0384	7.3911	0.007879 **
3	89	72.711	9	3.3647	0.4576	0.899016

ANOVA easily compares nested models where the independent variables of a simpler model make a subset of the independent variables of a more complex model. We order the models from the most simple to the most complex and ANOVA compares all the pairs of neighboring models. The ANOVA table suggests that lm.sel outperforms lm.const while lm.all does not further improve lm.sel. This

conclusion is in contradiction with the conclusion in the previous answer. The contradiction arises from a methodological fault that we did. We used the same dataset to select $P7$ as the best variable and to test whether it performs well. This approach suffers from bias and could be misleading.

- (c) (2 p) The dataset under consideration contains 100 samples. How do the type I error and type II error in the individual coefficient tests change with increasing number of samples if we maintain a constant level of significance α ?

Type I error is a controlled parameter and its probability remains unchanged with the α value unchanged. However, the power of the test will increase, so the type II error will decrease. At the same time, the robustness of RSS, R^2 and consequently the F-test power will increase as well.

- (d) (4 p) Describe in detail the way in which you would validate your models over the samples that you currently have. You can create additional auxiliary models. Describe the validation method, define the error function, and specify with which baseline you will compare the calculated error.

Let us assume that we want to compare $lm.const$, $lm.sel$ and $lm.all$. Let us assume that our sample set is small and thus the hold-out method that splits the sample set on training and testing set is inappropriate (we need to use as many training samples as possible, the same holds for testing set). Then, a good option seems to be to run 10-fold cross-validation. We will always train our models on 9 folds and test them on the remaining one. We will gradually shift the testing fold. The dependent variable is continuous, we can use the root mean square error (RMSE) or mean absolute percentage error (MAPE). The error will always be calculated over the testing fold and averaged over the folds. If we repeat 10-fold cross-validation multiple times, we can statistically test whether performances of the individual models truly differ.

Watch out. Feature selection is a part of training process. It cannot be done only once before cross-validation, it must be repeated again and again for each split. Consequently, we will have 10 different $lm.sel$ models to test, the set of relevant variables included into the model may change over folds as well as their regression coefficients. These 10 models will serve to estimate the performance of the final $lm.sel$ model. Only the final model (to be reported and deployed) could be based on all the available samples, and will thus certainly employ the variable $P7$.

Problem 3. (10 p) There is a cubic spline with one knot ξ given by the formula: $f(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta_4(x - \xi)_+^3$.

- (a) (1 p) Define the basis function $(x - \xi)_+^3$.

The definition is: $(x - \xi)^3$ for $x > \xi$, otherwise 0.

- (b) (1 p) How many degrees of freedom does the given cubic spline have? Why?

A cubic spline with K knots has $K + 4$ parameters or degrees of freedom. Our spline has one knot and thus it has 5 independent parameters/degrees of freedom. The number of parameters can be seen from the formula above too, there are β_0, \dots, β_4 there.

- (c) (1 p) What are the properties of the cubic spline at the knot?

The spline is continuous at the knot and it has a continuous first and second derivative there. The properties follow from the general properties for d-degree splines.

- (d) (2 p) Write down a cubic spline with one knot as a piecewise polynomial. Note: you will only change the form of notation, name the parameters differently from the spline parameters above.

$$f_1(x) = a_1 + b_1x + c_1x^2 + d_1x^3 \text{ for } x < \xi$$

$$f_2(x) = a_2 + b_2x + c_2x^2 + d_2x^3 \text{ for } x \geq \xi$$

- (e) (3 p) Express the piecewise polynomial parameters using the cubic spline parameters $\beta_0, \beta_1, \dots, \beta_4$.

The procedure is straightforward: the spline must match f_1 before the knot and f_2 after the knot. For the first polynomial it is trivial, because the basis function is zero before the first knot: $a_1 = \beta_0, b_1 = \beta_1, \dots, d_1 = \beta_3$. For the second polynomial it holds: $a_2 + b_2x + c_2x^2 + d_2x^3 = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta_4(x - \xi)^3$. By developing the last term we get: $\beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta_4(x^3 - 3x^2\xi + 3x\xi^2 - \xi^3) = (\beta_0 - \beta_4\xi^3) + (\beta_1 + 3\beta_4\xi^2)x + (\beta_2 - 3\beta_4\xi)x^2 + (\beta_3 + \beta_4)x^3$, of which follows: $a_2 = \beta_0 - \beta_4\xi^3, b_2 = \beta_1 + 3\beta_4\xi^2, c_2 = \beta_2 - 3\beta_4\xi, d_2 = \beta_3 + \beta_4$.

- (f) (2 p) Proof that the piecewise cubic polynomial found in the previous two steps maintains the knot properties of a cubic spline.

Continuity $f_1(\xi) = f_2(\xi)$ can be proven by substituting for coefficients a, b, c, d : $f_1(\xi) = \beta_0 + \beta_1\xi + \beta_2\xi^2 + \beta_3\xi^3, f_2(\xi) = \beta_0 - \beta_4\xi^3 + (\beta_1 + 3\beta_4\xi^2)\xi + (\beta_2 - 3\beta_4\xi)\xi^2 + (\beta_3 + \beta_4)\xi^3 = \beta_0 + \beta_1\xi + \beta_2\xi^2 + \beta_3\xi^3 = f_1(\xi)$.

Continuity of the first derivative $f_1'(\xi) = f_2'(\xi)$ can be confirmed by substituting for the coefficients and deriving: $f_1'(\xi) = \beta_1 + 2\beta_2\xi + 3\beta_3\xi^2 = f_2'(\xi)$.

Continuity of the second derivative $f_1''(\xi) = 2\beta_2 + 6\beta_3\xi = f_2''(\xi)$.