

Simple Linear Regression with R

Martin Schweinberger

December 28, 2017

Regression models are used to evaluate if one, or more predictors (or independent variables) correlate significantly with a dependent or outcome variable. The underlying logic is almost exactly identical to ANOVA (Analysis of Variance) designs, but in linguistics regression models have traditionally been used in analyses based on corpus data (e.g. sociolinguistic studies) while ANOVAs were the method of choice in psycholinguistics. The field-specific bias is, however, historical and cultural rather than caused by design specific advantages.

One difference between ANOVA and regression analysis is that regression analyses use a *t-statistic* rather than an *F-statistic* as ANOVAs do (but F is simply t squared ;-)). Both measure the ratio of explained to unexplained variance though.

The idea behind regression analyses can be best understood visually: imagine a line going through the data points in a coordinate system. What regression analyses do is to aim at finding that line going through the data points which has the smallest sum of *residuals*. Residuals are the difference between any observed data point and the predicted value for that data point. And the *sum of residuals* is the error in/of the regression model. The line of best fit line is then called the *regression line* (also called abline or the regression model). The slope of the line is called the *coefficient* and the point where the line crosses the y-axis is called the *intercept*.

The output in regression analyses provide a *standard error* (SE) which shows how much the coefficients (or estimates) differ across samples. A low standard error means that the coefficient will not differ much if we look at many samples while a high SE indicates that the coefficient will differ a lot between samples.

The example we are going to look at concerns the frequency of prepositions in the history of the English language. Using the *Penn Corpora of His-*

torical English (aka the Penn Parsed Corpus; see <http://www.ling.upenn.edu/hist-corpora/>) which consist of 605 texts written between 1125 and 1900, I extracted all words with the part-of-speech tag "/P" (preposition), calculated the per-1,000-words frequency of prepositions per text. After some data processing, we are left with a table of two columns: one with the text.id and the other with the normalized relative frequency of prepositions in the text. The question I want to investigate is whether the frequency of prepositions has increased over time.

The r code below is used to a) extract and process the data; b) visualize the data; c) perform the regression analysis, and d) test if the assumptions, that a linear regression is based on, are met.

```

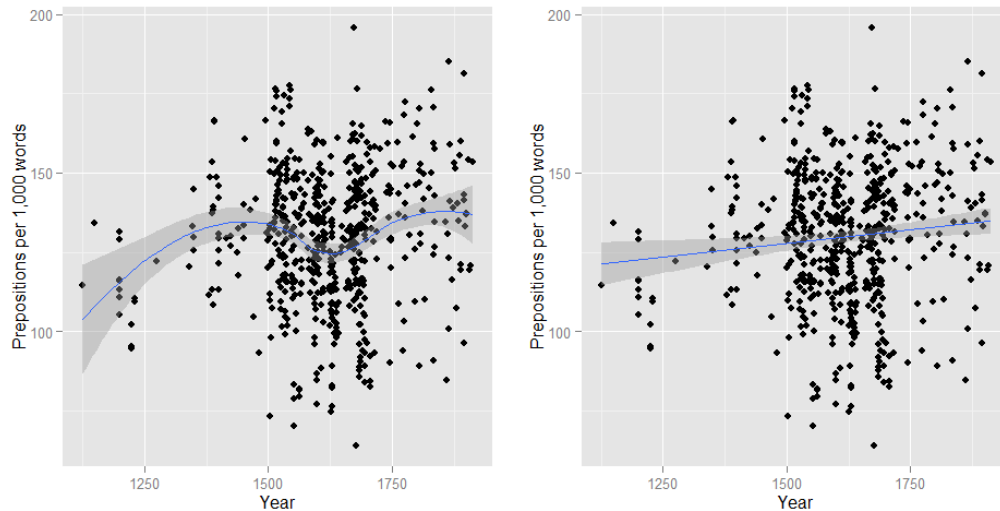
1  ### --- Prepare data
2  # Remove all lists from the current workspace
3  rm(list=ls(all=T))
4  # Install packages we need or which may be useful
5  # (to activate just delete the #)
6  #install.packages("QuantPsyc")
7  #install.packages("car")
8  # Initiate the packages
9  library(QuantPsyc)
10 library(car)
11 library(ggplot2)
12 source("C:\\R/multiplot_ggplot2.R")
13 # this call will only work, if you have stored this function
14   in
15 # the directory specified in the call
16 source("C:\\R/slr.summary.tb") # load modified summary
17   function for slr
18
19 ##
20 #####
21
22 ### Load and manipulate data
23 ##
24 #####
25
26 # WARNING: To use this script you need to set our own paths!
27 # Your path should be the path to the corpus on your own
28   computer!
29 # Remember to use double backslash instead of single
30   backslash, if
31 # you use Windows on your machine.
32 # Read in data
33 slr.data <- read.delim("C:\\MyProjects\\
34   SimpleLinearRegression/slr.data.txt", header = TRUE)
35 attach(slr.data) # we attach the data so we don't need to

```

```

27     specify the path all the time
28 # remove columns we do not need
29 slr.data <- as.data.frame(cbind(slr.data$datems, slr.data$P.
    ptw))
30 colnames(slr.data) <- c("year", "prep.ptw") # add column
    names
31 slr.data <- slr.data[!is.na(slr.data$year) == T, ] # delete
    incomplete cases
32 head(slr.data) # inspect data
33 # year prep.ptw
34 #1 1736 166.01
35 #2 1711 139.86
36 #3 1808 130.78
37 #4 1878 151.29
38 #5 1743 145.72
39 #6 1807 152.59
40 str(slr.data) # inspect data structure
41 #'data.frame': 603 obs. of 2 variables:
42 # $ year : num 1736 1711 1808 1878 1743 ...
43 # $ prep.ptw: num 166 140 131 151 146 ...
44 summary(slr.data) #inspect data values
45 # year prep.ptw
46 # Min. :1125 Min. : 63.97
47 # 1st Qu.:1545 1st Qu.:115.66
48 # Median :1615 Median :130.78
49 # Mean :1619 Mean :129.81
50 # 3rd Qu.:1687 3rd Qu.:144.08
51 # Max. :1913 Max. :195.86
52 #####
53 ### Visualize and eyeball data
54 # set up this plot from scratch
55 p2 <- ggplot(slr.data, aes(year, prep.ptw)) +
56   geom_point() +
57   labs(x = "Year") +
58   labs(y = "Prepositions per 1,000 words") +
59   geom_smooth()
60 # set up this plot from scratch
61 p3 <- ggplot(slr.data, aes(year, prep.ptw)) +
62   geom_point() +
63   labs(x = "Year") +
64   labs(y = "Prepositions per 1,000 words") +
65   geom_smooth(method = "lm") # with linear model smoothing!
66
67 multiplot(p2, p3, cols = 2)

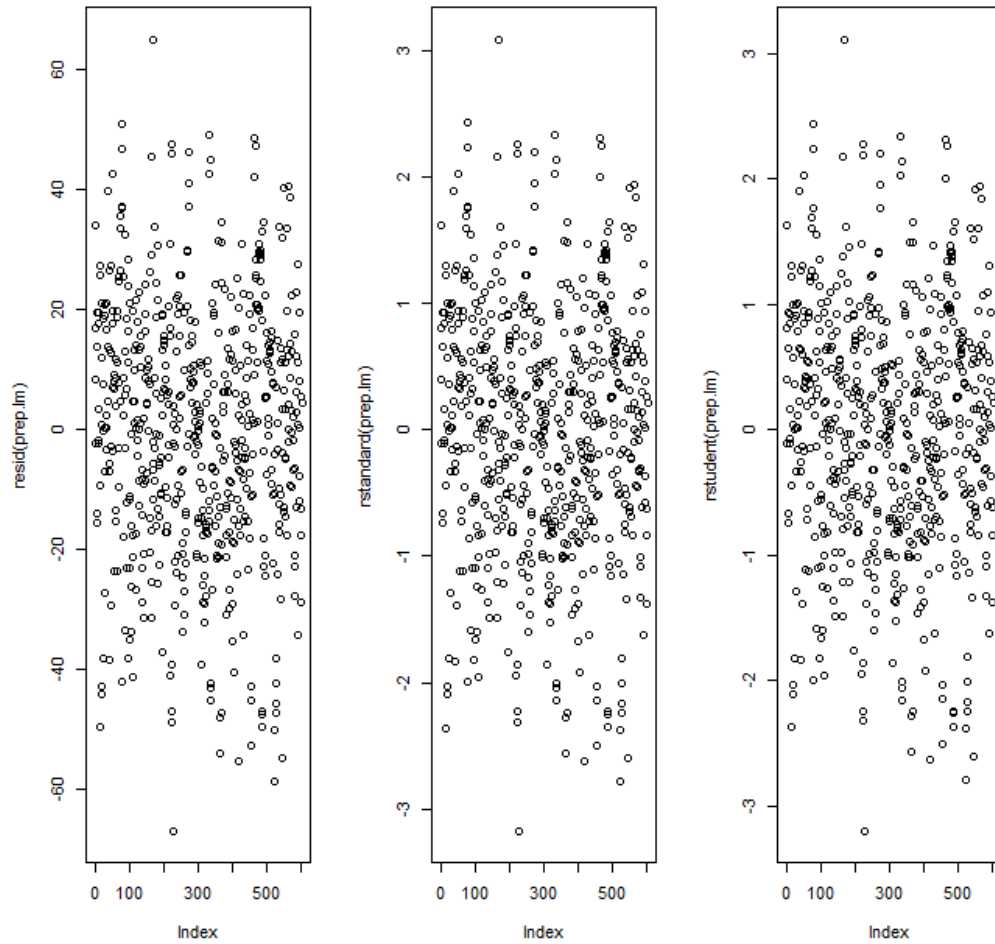
```



```

1 # set up Simple Linear egression model and inspect properties
  of the model
2 prep.lm <- lm(preop.ptw ~ year, data = slr.data)
3 summary(preop.lm)
4 #Call:
5 #lm(formula = prep.ptw ~ year, data = slr.data)
6 #Residuals:
7 # Min 1Q Median 3Q Max
8 #-66.842 -13.523  1.183 14.086 65.117
9 #Coefficients:
10 # Estimate Std. Error t value Pr(>|t|)
11 #(Intercept) 1.021e+02 1.086e+01 9.397 <2e-16 ***
12 #year 1.713e-02 6.691e-03 2.560 0.0107 *
13 #---
14 #Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '
  1
15 #Residual standard error: 21.11 on 601 degrees of freedom
16 # (2 observations deleted due to missingness)
17 #Multiple R-squared: 0.01079, Adjusted R-squared: 0.00914
18 #F-statistic: 6.553 on 1 and 601 DF, p-value: 0.01071
19 # use plots to check if there are problems with the model
20 # set graphic's parameters to display 3 plots in one row
21 par(mfrow = c(1, 3))
22 plot(resid(preop.lm))
23 plot(rstandard(preop.lm))
24 plot(rstudent(preop.lm))
25 par(mfrow = c(1, 1)) # restore original graphic's parameters

```

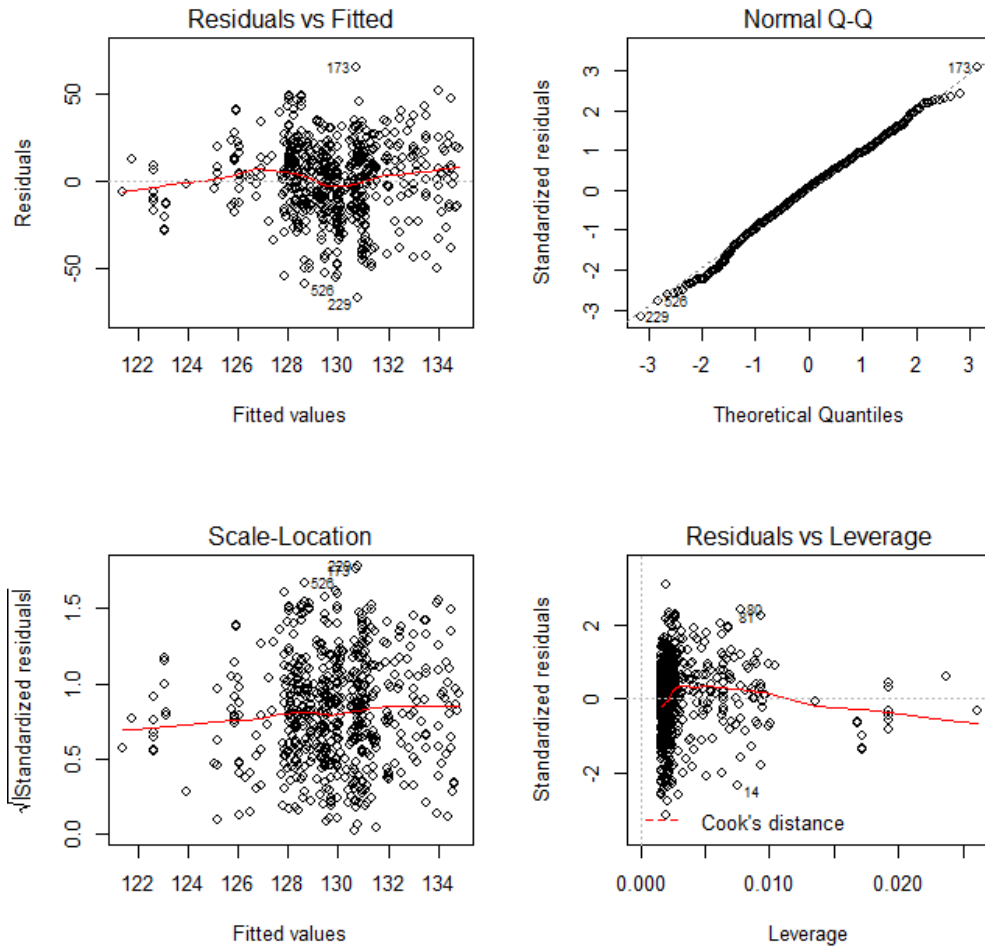


The left plot shows the residuals of the model (the difference between the observed and the predicted value). A shortcoming of this plot is that the residuals are not standardized, i.e. we cannot compare them to residuals from another regression model. To overcome this shortcoming, we have a look at standardized residuals which are residuals divided by their standard deviation (see the plot in the middle). This way, residuals are converted into standardized residuals, but more importantly, they now represent z-scores and we can use the z-distribution to determine which points are problematic. Here are three rules of thumb for diagnosing regression models based on residuals (see Field, Miles & Field (2012:268-269):

1. Points with extreme values, i.e. values above 3 (3.29 to be exact), should be excluded.
2. If more than 1% of data points in our sample has values above 2.5 (2.58 to be exact), then the error in our model is too high.
3. If more than 5% of data points in our sample has values above 2 (1.96 to be exact), then the error in our model is too high.

The plot on then right shows the *studentized residuals*, i.e. the *adjusted predicted value* of each point is divided by the standard error of the residuals. This way, we can use the Student's t-distribution to diagnose our model. Adjusted predicted values are also residuals, but a specific kind of residual: a model is calculated without one data point, then, this model is used to predict the data point. The difference between the observed data point and the point predicted by the model without this point is called the adjusted predicted value. In summary, studentized residuals are very helpful to detect influential data points. The plots show that there are two potentially problematic cases – the dots which are at the very top and at the very bottom. These two data points stand apart from the other points and are thus probably outliers. We will check later if they need to be removed.

```
1 # Create a 2x2 matrix of diagnostic plots
2 par(mfrow = c(2, 2))
3 plot(preplm)
4 par(mfrow = c(1, 1))
```



The diagnostic plots look very good - but what does this mean? (I will go more into detail in the post on Multiple Linear Regression)

Well, the upper left plot is useful for a) finding outliers, or b) finding correlations between residuals and fitted values: if there were a noticeable trend visible in the line or the data points (e.g. an upward, downward, or zig-zag trend in the line), we would have a problem and would have to delete data points.

The upper right plot is useful for checking if the residuals are normally distributed (a necessary condition for linear regressions). If the dots lie on the line then the residuals are indeed distributed following a normal distribution. If the dots diverge from the line at the top and bottom, then this indicates

poor model fit.

Let's consider the lower left plot: regression models are also based on the assumption of homoscedasticity, i.e. that the variance of the residuals do not change with x or - to put it differently - that the variance of the residuals do not correlate with the predictor or predictors in multiple regression). The assumption of homoscedasticity is met, then we should see a flat line without upward or downward trend. If there is a trend, you have a problem called heteroscedasticity.

The lower right plot is useful for detecting leverage, i.e. points that over-proportionally influence the regression model (this should not be the case as in simple linear regression all data points have the same weight). If there are points that have a high leverage, you should try a) to use a robust linear regression (which assigns different weights to the data points or b) delete the data points as outliers. This plot also shows Cook's distance which is a measure that captures how much the regression would change if the point in question was deleted. So Cook's distance assesses the influence of individual points on the model as a whole. Points which have a value for Cook's distance that is greater than 1 are problematic (Field, Miles & Field 2012:269). Leverage is also a measure that provides an estimate of how much a data point influences the accuracy of the model. Leverage values lie between 0 (no influence) to 1 (major influence: no good!). To check if a given point is still acceptable, you need to calculate a cut-off point (either $3(k + 1)/n$ or $2(k + 1)/n$).

We will now check, if we need to exclude data points and also calculate the standardized β .

```

1 # Extract standardized betas
2 lm.beta(preplm)
3 # year
4 #0.1038566
5 #####
6 ### Write a function for a neat output table
7 lm.summary <- function(x) {
8   p.nice <- function(z) {
9     as.vector(unlist(sapply(z, function(w) {
10      ifelse(w < .001, return("p < .001***"),
11      ifelse(w < .01, return("p < .01**"),
12      ifelse(w < .05, return("p < .05*"), return(w)))) } ))) }
13   intercept <- c(
14     round(summary(x)[[4]][1], 2),
15     "",
16     round(summary(x)[[4]][3], 2),
17     round(summary(x)[[4]][5], 2),
18     round(summary(x)[[4]][7], 4),

```



```

19 p.nice(summary(x)[[4]][7]))
20 predictor <- c(
21 round(summary(x)[[4]][2], 2),
22 round(lm.beta(x)[[1]], 4),
23 round(summary(x)[[4]][4], 2),
24 round(summary(x)[[4]][6], 2),
25 round(summary(x)[[4]][8], 4),
26 p.nice(summary(x)[[4]][8]))
27 mdl.statz <- c("", "", "", "", "", "Value")
28 nbcases <- c("", "", "", "", "", length(summary(x)[[3]]))
29 rse <- c("", "", "", "", "",
30 round(summary(x)[[6]], 2))
31 multR2 <- c("", "", "", "", "", round(summary(x)[[8]], 4))
32 adjR2 <- c("", "", "", "", "", round(summary(x)[[9]], 4))
33 F <- c("", "", "", "", "",
34 round(summary(x)[[10]][1], 2))
35 p <- c("", "", "", "", "", round(summary(x)[[4]][8], 4))
36 slrm.tb <- rbind(intercept, predictor, mdl.statz, nbcases,
37 rse, multR2, adjR2, F, p)
38 colnames(slrm.tb) <- c(colnames(summary(x)[[4]])[1],
39 "Std. Beta",
40 colnames(summary(x)[[4]])[c(2:4)],
41 "P-value sig.")
42 rownames(slrm.tb) <- c(
43 rownames(summary(x)[[4]])[1],
44 rownames(summary(x)[[4]])[2],
45 "Model statistics", "Number of cases in model",
46 paste("Residual standard error", paste("on", summary(x)
47 [[7]][2], "DF")),
48 "Multiple R-squared", "Adjusted R-squared",
49 paste("F-statistic",
50 paste("(", round(summary(x)[[10]][2], 0), ", ",
51 round(summary(x)[[10]][3], 0), ")"), sep = "", collapse = "")
52 ,
53 "Model p-value")
54 slrm.tb <- as.data.frame(slrm.tb)
55 return(slrm.tb)
56 }
57 lm.summary(preplm) # inspect the results
58 # Estimate Std. Beta Std. Error t value Pr(>|t|) P-value sig.
59 #(Intercept) 102.09 10.86 9.4 0 p < .001***
60 #year 0.02 0.1039 0.01 2.56 0.0107 p < .05*
61 #Model statistics Value
62 #Number of cases in model 603
63 #Residual standard error on 601 DF 21.11
64 #Multiple R-squared 0.0108
65 #Adjusted R-squared 0.0091
66 #F-statistic (1, 601) 6.55
67 #Model p-value 0.0107

```

Typically, the results of regression models are reported using a table like the one we just created (see again below).

	Estimate	Std. Beta	Std. Error	t-value	Pr(> t)	P-value sig.
(Intercept)	102.09		10.86	9.4	0	p < .001***
year	0.02	0.1039	0.01	2.56	0.0107	p < .05*
Model statistics						Value
Number of cases in model						603
Residual standard error on 601 DF						21.11
Multiple R-squared						0.0108
Adjusted R-squared						0.0091
F-statistic (1, 601)						6.55
Model p-value						0.0107

However, you can also summarize the results of regression analyses in prose form which would be something like this:

A simple linear regression model was fitted to the data. Visual inspection of diagnostic plots and data driven methods did not suggest problems concerning outliers and model fit. The final linear regression model was based on 603 data points and correlated significantly with the data (R^2 : 0.0108, F-statistic (1, 601): 6.553, p-value: 0.0107*) and confirmed a significant positive correlation between the year the text was written in and the relative frequency of prepositions in the text (coefficient: .02, std. β : 0.1039, SE: 6.691e-03, t-value: 2.560, p-value: .0107 *).

References

Field, Andy, Jeremy Miles & Zoe Field. 2012². *Discovering statistics using R*. London, Thousand oaks, CA, New Delhi, Singapore: SAGE.