

Business Analytics & Machine Learning

Tutorial sheet 2: Regression – Solution

Prof. Dr. Martin Bichler, Prof. Dr. Jalal Etesami
Julius Durmann, Markus Ewert, Johannes Knörr, Yutong Chao

Exercise T2.1 Gross national product

The following table displays the per capita gross national product (X in \$1,000) and the percentage of literate people among the population (Y).

Country	X	Y
Nepal	0.5	5
Uganda	0.6	28
Thailand	1.0	68
South Korea	1.4	77
Peru	1.8	48
Lebanon	3.6	48
Ireland	5.7	98
France	6.4	96
New Zealand	13.0	99

For convenience, we have precomputed: $\sum x_i = 34$, $\sum x_i^2 = 262.22$, $\sum y_i = 567$, $\sum x_i y_i = 2,914.3$, $\bar{x} = 3.78$, and $\bar{y} = 63$.

- a) Calculate the coefficients $\hat{\beta}_0$ and $\hat{\beta}_1$ of the simple linear regression model using *ordinary least squares*. Find the regression line using the formulas below:

$$\begin{aligned}\hat{\beta}_1 &= \frac{\text{cov}(x, y)}{\text{var}(x)} = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\left(\frac{1}{n} \sum_{i=1}^n x_i y_i\right) - \bar{x} \bar{y}}{\left(\frac{1}{n} \sum_{i=1}^n x_i^2\right) - \bar{x}^2}\end{aligned}$$

and

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

- b) Interpret the coefficients calculated in exercise a).
- c) Test the zero hypothesis $H_0 : \beta_1 \leq 0$ with significance level of $\alpha = 0.05$. Use the following *t-test* with a residual sum of squares of $\text{RSS} = 4,411.4$ and $\sum_{i=1}^n (x_i - \bar{x})^2 = 133.77$:

$$t_0 = \frac{\hat{\beta}_1}{\text{SE}(\hat{\beta}_1)} \sim t_{n-2}, \quad \text{SE}(\hat{\beta}_1) = \sqrt{\frac{\text{RSS}}{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \frac{1}{n-2}}.$$

- d) Now the above linear regression model will be used to estimate the percentage of literates among a country with known gross national product. Which problems might occur? Briefly explain your concerns using an example.

Solution

- a) The formulas for the coefficients are:

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}}{\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2} = \frac{\frac{1}{9}(2914.3) - (3.78)(63)}{\frac{1}{9}(262.22) - (3.78)(3.78)} \approx 5.77$$

and

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \approx 63 - (5.77)(3.78) \approx 41.19$$

Now, the regression line is given as $\hat{y}(x) = 41.19 + 5.77 \cdot x$.

- b) $\hat{\beta}_0$: A country with a capita gross national product of 0 has (approximately) 41.19 % of literate people among the population.

$\hat{\beta}_1$: With each increase of \$1,000, the percentage of literate people among the population increases by (approximately) 5.77.

- c) Use the “test manual” from last week’s tutorial to solve the exercise.

- 1) Not needed, because only one test is used in regression analysis.
- 2) $H_0 : \beta_1 \leq 0$ vs. $H_1 : \beta_1 > 0$.
- 3) Simple t -test for coefficient. We have

$$SE(\hat{\beta}_1) = \sqrt{\frac{RSS}{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \frac{1}{n-2}} \approx 2.17 \quad \text{with} \quad RSS \approx 4,411.4, \quad \sum_{i=1}^n (x_i - \bar{x})^2 \approx 133.77$$

and thus $t_0 \approx 2.66$.

- 4) $\alpha = 0.05$.
 - 5) For $df = n - 2 = 7$ degrees of freedom, we find a critical value of $t_{7,0.95}^c = 1.895$ in the t -table.
 - 6) We reject H_0 because $t_0 > t^c$ and conclude that $\hat{\beta}_1$ is statistically significant.
- d) A prediction for countries with per capita gross national product outside the sample range is problematic. For example, any country with a gross national product smaller than Nepal or greater than New Zealand may have illogical predicted percentage of literate people among population ($< 0\%$ or $> 100\%$).

Exercise T2.2 Testing Gauss-Markov assumptions

Please use the provided Jupyter notebook to solve this task.

You might need to install `statsmodels` (see [here](#) for the documentation) by running

```
pip install statsmodels
```

(Make sure that your virtual Python environment is active!)

You are given the data set in *gauss-markov.csv*. It contains values for three variables X_1 , X_2 , X_3 and values for a target variable Y . Our goal is to predict the target variable based on the three input variables.

- a) We start by using the simple linear regression model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

Using statsmodels, compute optimal values for the parameters. Let the model predict the values of $\hat{y} \approx y$.

Note: You may want to use `statsmodels.api.add_constant()` to add constant values for the intercept.

- b) Compute the residuals $e = \hat{y} - y$ of the resulting model. Plot the residuals over the input variables x_1 and x_2 . What do you observe?

Using a White test, show that we can reject the hypothesis of homoscedastic residuals at an α level of 0.01.

- c) Consider the alternative model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_1^2$$

Compute the optimal parameter values. You should observe that the R^2 value improves drastically over the previous model.

Although this model gives a very good fit of the data, there is another problem: Multicollinearity. Use the Variance inflation factor to check whether the variables are dependent.

- d) Consider a third model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1^2$$

and compute the optimal parameters.

Check if the model has multicollinear input variables using the VIF criterion.

Check if the model satisfies the homoscedasticity assumption using the White test and an α level of 0.01.

Solution

Please also refer to the [provided solution notebook](#).

- a) Parameter values:

β_0	31.7239
β_1	-14.0285
β_2	14.3497
β_3	2.5067

$$R^2 = 0.663.$$

b) Result of White test ($LM \sim \chi_3^2$): $LM = 73.744$, $p\text{-value} = 2.8 \times 10^{-12} \ll 0.01$

Since the p-value is very low, we can reject the null hypothesis (homoscedasticity).
 \Rightarrow Heteroscedastic residuals are very likely.

c) Parameter values:

β_0	3.9063
β_1	-3.8780
β_2	-0.4725
β_3	-0.8776
β_4	0.9977

$R^2 \approx 1$.

VIF values:

X_1	355.35
X_2	752.47
X_3	1080.27
X_1^2	1.01

The VIF values are very high ($\gg 10$). Consequently, we may assume that model has multicollinear input variables.

(Bonus) Result of White test ($LM \sim \chi_4^2$): $LM = 16.367$, $p\text{-value} = 0.23 > 0.01$

d) Parameter values:

β_0	2.3017
β_1	-6.5107
β_2	3.3908
β_3	0.9949

$R^2 = 0.99$.

VIF values:

X_1	1.00
X_2	1.00
X_1^2	1.00

The VIF values are low (< 10). Consequently, the model inputs are likely independent.

Result of White test ($LM \sim \chi_3^2$): $LM = 4.315$, $p\text{-value} = 0.83 > 0.01$

The p-value for the null hypothesis (homoscedastic residuals) is high. Therefore, the null hypothesis needn't be rejected.

Exercise T2.3 *Derivation of closed-form solution*

In this exercise, we will derive the closed-form solution of the regression problem

$$\beta^* = (X^T X)^{-1} X^T y \quad (1)$$

where

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

summarize our input variables $x_i \in \mathbb{R}^d$ and target variables $y_i \in \mathbb{R}$.
The model which we use is

$$\hat{y}_i = \beta^T x_i$$

You may assume that the inputs x_i already contain an entry equal to 1 which allows to include the intercept of our model without further consideration. Also, you may assume that the vectors in X are not colinear.

a) Formulate the sum of squared errors $e_i = \hat{y}_i - y_i$

$$\mathcal{L} = \sum_{i=1}^n e_i^2$$

which we seek to minimize in our analysis:

- (i) In terms of the individual elements x_i, y_i
- (ii) In terms of the matrix notation X, y

b) Calculate the derivatives (gradients):

$$\begin{aligned} \text{(i)} \quad \frac{\partial}{\partial \beta}(\beta^T a) &= \begin{pmatrix} \frac{\partial}{\partial \beta_1}(\beta^T a) \\ \vdots \\ \frac{\partial}{\partial \beta_d}(\beta^T a) \end{pmatrix} \text{ for } \beta, a \in \mathbb{R}^d \\ \text{(ii)} \quad \frac{\partial}{\partial \beta}(\beta^T A \beta) &= \begin{pmatrix} \frac{\partial}{\partial \beta_1}(\beta^T A \beta) \\ \vdots \\ \frac{\partial}{\partial \beta_d}(\beta^T A \beta) \end{pmatrix} \text{ for } \beta \in \mathbb{R}^d, A \in \mathbb{R}^{d \times d} \end{aligned}$$

Check your results with the matrix cookbook [1], chapter 2.4.

- c) Use these derivatives to compute the gradient of the loss: $\frac{\partial}{\partial \beta} \mathcal{L}(\beta)$.
- d) Set the derivative to zero (first order condition) to obtain (1).
- e) Why is there no need to check second-order derivatives to prove optimality?

Solution

a)

$$\begin{aligned} \mathcal{L} &= \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (\hat{y}_i - y_i)^2 \\ &= \sum_{i=1}^n (x_i^T \beta - y_i)^2 \quad (i) \\ &= (X\beta - y)^T (X\beta - y) = \beta^T X^T X \beta - 2\beta^T X^T y + y^T y \quad (ii) \end{aligned}$$

- b) (i) Derive the individual components $\frac{\partial \beta^T a}{\partial \beta_i}$ and aggregate to obtain the gradient:

$$\frac{\partial \beta^T a}{\partial \beta_i} = a_i \quad \Rightarrow \quad \frac{\partial}{\partial \beta}(\beta^T a) = a$$

- (ii) Reformulate the matrix-vector product to a summation of scalars:

$$\beta^T A \beta = \sum_{i=1}^d \sum_{j=1}^d \beta_i A_{ij} \beta_j$$

Derive the individual components $\frac{\partial \beta^T A \beta}{\partial \beta_k}$:

$$\begin{aligned} \frac{\partial \beta^T A \beta}{\partial \beta_k} &= \sum_{i=1}^d \sum_{j=1}^d \frac{\partial (\beta_i A_{ij} \beta_j)}{\partial \beta_k} = \sum_{i=1}^d \sum_{j=1}^d \mathbb{I}_{i=k} A_{ij} \beta_j + \beta_i A_{ij} \mathbb{I}_{j=k} \\ &= \sum_{j=1}^d A_{kj} \beta_j + \sum_{i=1}^d \beta_i A_{ik} = A_{k,:} \beta + \beta^T A_{:,k} = (A_{k,:} + (A_{:,k})^T) \beta \end{aligned}$$

where $A_{k,:}$ denotes the k^{th} row of A and $A_{:,k}$ denotes the k^{th} column of A .

Aggregate the gradient:

$$\frac{\partial \beta^T A \beta}{\partial \beta} = (A + A^T) \beta$$

c)

$$\begin{aligned} \frac{\partial \mathcal{L}(\beta)}{\partial \beta} &= \frac{\partial}{\partial \beta} (X\beta - y)^T (X\beta - y) = \frac{\partial}{\partial \beta} \beta^T X^T X \beta - \frac{\partial}{\partial \beta} 2\beta^T X^T y + \frac{\partial}{\partial \beta} y^T y \\ &= \underbrace{(X^T X + (X^T X)^T)}_{\text{b)(ii)}} \beta - \underbrace{2X^T y}_{\text{b)(i)}} = 2X^T X \beta - 2X^T y \end{aligned}$$

d)

$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta} = 0 \quad \Leftrightarrow \quad 2X^T X \beta - 2X^T y = 0 \quad \Leftrightarrow \quad 2X^T X \beta = 2X^T y \quad \Leftrightarrow \quad \beta = (X^T X)^{-1} X^T y$$

- e) The objective is fully convex in β (squared form). The solution from d) therefore is the unique global minimum.

—

- [1] Kaare Brandt Petersen and Michael Syskind Pedersen. *The Matrix Cookbook*. URL: <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>.