

# Business Analytics & Machine Learning

## Tutorial sheet 2: Regression

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.

## 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  $\alpha$  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  $\alpha$  level of 0.01.

### 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 [**MatrixCookbook**], 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?