

Introduction to R

Exam Summer Term 2018

Student Name, Student ID, StudentName@uni-muenster.de

Contents

1	Linear Equation	2
2	Functions	3
3	Passenger numbers	4
4	Porsche	5
5	Law of large numbers	6
6	Limits of maxima	7
7	Student teacher ratio	8
8	Asymptotic normality	9
9	Stochastic frontier analysis	10
10	Variance estimation in GARCH	11

-
- Answer **7** out of **10** of the following exercises in either German or English.
 - Hand in your solutions before Friday, 13 April 2018 at 10 am.
 - It is advised to regularly check the learnweb and your emails in case of urgent updates.
 - Please send your solutions files to Willi Mutschler. We will confirm the receipt of your work also by email.
 - The solution files should contain your executable and commented script file or preferably a R Notebook.
 - You may use *any* available R package you find fit to solve the exercise.
 - Please label your axes and title in your plots.
 - **All students must work on their own.**
-

1 Linear Equation

Solve the linear equation $A \cdot x = b$ with

$$A = \begin{pmatrix} 1 & 1 & 1 & 3 \\ 2 & 5 & 8 & 6 \\ 4 & 3 & 7 & 9 \\ 3 & 6 & 5 & 1 \end{pmatrix} \text{ and } b = \begin{pmatrix} 1 \\ 4 \\ 3 \\ 5 \end{pmatrix}$$

Solution:

and compute the inverse as well as Eigenvalues of A .

Solution:

2 Functions

1. Write a function `psum(n,a)` that computes

$$s_{n,a} := \sum_{k=0}^n \frac{k^a}{k^a + 1}$$

for any natural number $n \in \{1, 2, \dots\}$ and any $a > 0$.

Solution:

2. Write a function `mymatrix(n)` that returns a $n \times n$ matrix such that: the first and last row as well as the first and last column contain only ones, whereas the remaining values are zero, e.g.

$$\begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}$$

Solution:

3 Passenger numbers

The file **apass.csv** contains monthly data on passenger numbers of US airlines from January 1949 to December 1959.

1. Read the data into a data frame.

Solution:

2. Create a vector that contains the corresponding dates. Add this vector to your data frame.

Solution:

3. Plot the passenger numbers against the date vector.

Solution:

4. Calculate the mean passenger numbers for each month. Plot the means as a bar chart.

Solution:

4 Porsche

The file **Porsche911.csv** contains data on Porsche 911 cars, **name** is an id for the owner, **loc** is the location, **age** is the age of the car, **TKM** is the mileage in thousands kilometers and **price** the current price listed on an internet platform for used cars in thousand Euros.

1. Read in the data into a data frame and compute key descriptive statistics (mean, standard deviation, smallest and largest values, quartiles, covariance and correlation) for the variables **age**, **TKM** and **price**.

Solution:

2. Plot boxplots for **age**, **TKM** and **price**.

Solution:

3. Generate the empirical cumulative distribution function as well as histograms for each of the variables **age**, **TKM** and **price**.

Solution:

5 Law of large numbers

Let X_1, X_2, \dots be a sequence X_1, X_2, \dots from an $AR(1)$ process:

$$(X_i - \mu) = \rho (X_{i-1} - \mu) + \varepsilon_i$$

where ε_i is uniformly distributed on the interval $[-1, 1]$ and $|\rho| < 1$. Define the sequence of random variables

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

The weak law of large numbers states that the sample average \bar{X}_n converges in probability towards the expected value μ when n tends to infinity.

Show by simulation that the law of large numbers holds despite the intertemporal dependence in X . In particular, show the convergence by means of an appropriate plot. Hint: You may set the parameters to e.g. $\rho = 0.5$ and $\mu = 2$ (or any other value you find fit.)

Solution:

6 Limits of maxima

Let X_1, X_2, \dots be an i.i.d. sequence of random variables uniformly distributed on the interval $[0, 1]$. Define the random variables

$$M_n = \max_{i=1, \dots, n} X_i$$

and its normalized version $\overline{M}_n = (M_n - 1) \cdot n$. One can show that the limit distribution of \overline{M}_n is the Weibull distribution with density $\exp(x)$. Write an R program to illustrate that \overline{M}_n converges in distribution. To this end, set $n = 100$ and $R = 1000$ and consider the $R \times n$ matrix with uniformly distributed random variables X_i . Show the convergence by means of an appropriate plot.

Note: Try to avoid using loops (use e.g. `apply` instead). You will not get full points if you use a loop.

Solution:

7 Student teacher ratio

Load the dataset **caschool.csv** into the object **caschool**. This dataset is discussed in great detail in the textbook of Stock and Watson.

1. Make the following variables accessible:
 - test score **testscr**
 - student-teacher ratio **str**
 - percentage of English language learners **el_pct**
 - expenditures per student **expn_stu**

Solution:

2. Regress **testscr** on a constant and **str**. Assign the residuals of the regression into the variable **r1**. Now regress **testscr** on an intercept, **str**, **el_pct** and **expn_stu**. Put the residuals into the variable **r2**. Compute the sum of squared residuals for both regressions. Display **r1** and **r2** in one plot, where the points of **r2** are marked red.

Solution:

3. Consider the regression of **testscr** on a constant, **str**, **el_pct** and **expn_stu**. Predict the value of **testscr** for a school district with an average class size (**str**) of 25 students, a percentage of English learners (**el_pct**) of 60% and an average expenditures per student (**expn_stu**) of 4000\$.

Solution:

4. Reconsider the regression of **testscr** on a constant, **str**, **el_pct** and **expn_stu**. Compute the heteroscedastic robust standard errors.

Solution:

5. Test the null hypothesis that the coefficients on **str** and **expn_stu** both equal 0 and the coefficient on **el_pct** equals -0.7 . Hint: Use the **linearHypothesis** function of the **car** package.

Solution:

8 Asymptotic normality

Consider the multiple linear regression model $y = X\beta + u$. In R, generate the matrix X by executing the following commands:

```
library(MASS)
X <- cbind(1,mvrnorm(n=100,c(5,10),matrix(c(1,0.9,0.9,1),2,2)))
```

Assume that the true coefficient vector is

$$\beta = \begin{pmatrix} 3 \\ 2 \\ -1 \end{pmatrix}$$

and the error terms are i.i.d. uniformly distributed on the interval $[-1, 1]$. Hence, the assumption of normally distributed error terms is violated.

1. Write an R program that generates $R = 10000$ random samples of size $n = 100$ each. Generate an empty vector $Z \leftarrow \text{rep}(NA, R)$. For each sample $i = 1, \dots, R$, compute the OLS estimate $\hat{\beta}$ of β and store the second component of $\hat{\beta}$ in the i -th element of the vector Z .

Solution:

2. Plot the histogram of Z . Compute the mean m and standard deviation s of Z and add the density of $N(m, s)$ to the plot.

Solution:

9 Stochastic frontier analysis

Consider the Cobb-Douglas production function

$$y = Ax_1^\alpha x_2^\beta$$

By definition, the production function returns the maximal output for given inputs, and actual production cannot be larger than y . Due to inefficiencies, actual production could be modeled (in logs) as

$$\ln y = \ln A + \alpha \ln x_1 + \beta \ln x_2 - u$$

where u is a **non-negative** random variable. Since other disturbances (e.g. measurement errors) can enter the production function, it is more common to add another, **symmetrically distributed**, disturbance term v ,

$$\ln y = \ln A + \alpha \ln x_1 + \beta \ln x_2 - u + v$$

Assume that u is exponentially ($u \sim \text{Exp}(\lambda)$) and v normally ($v \sim N(0, \sigma^2)$) distributed. One can show that if u and v are independent then the density function of $\varepsilon = v - u$ is given by

$$f_\varepsilon(\varepsilon) = \lambda \exp\left(\lambda\varepsilon + \frac{1}{2}\lambda^2\sigma^2\right) \Phi\left(\frac{-\varepsilon}{\sigma} - \lambda\sigma\right)$$

where Φ is the distribution function (**pnorm**) of $N(0, 1)$ and \exp the exponential function (**exp**).

1. Load the dataset **sfa.csv**. This dataset is an abbreviated version of table F7.2 of Greene, 2008. The original data appeared in Zellner and Revankar, *Generalized Production Functions*, Review of Economic Studies, 36 (1969), 241-250.

Solution:

2. Write an R program to estimate the parameters A , α , β , λ and σ by maximum likelihood on this dataset.

Solution:

3. Compute the asymptotic standard errors.

Solution:

10 Variance estimation in GARCH

When one considers an iid sample X_1, \dots, X_n from $X \sim N(\mu, \sigma^2)$ then one usually estimates the variance σ^2 using

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

The distribution of the normalized estimator for the variance is given by:

$$\frac{(n-1)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-1}^2$$

where σ^2 is the true variance. Consider the case when the observations are not iid but are stochastically dependent over time. To this end, assume that X_1, \dots, X_n is a time series generated by a $GARCH(1,1)$ process

$$\begin{aligned} X_i &\sim N(0, \sigma_i^2) \\ \sigma_i^2 &= \omega + \alpha X_{i-1}^2 + \beta \sigma_{i-1}^2 \end{aligned}$$

with $\omega = 0.1$, $\alpha = 0.1$, $\beta = 0.85$ and sample size equal to $n = 2500$. Show by simulations that the distribution of the normalized estimator for the variance is not χ_{n-1}^2 -distributed. Hint: The true unconditional variance of this GARCH process is $\sigma^2 = 2$.

Solution: