

Exercises 3

Exercise 1

Use the Munich appartments dataset (`Miete2003.csv`) to estimate two different linear models with the net rent (`nettomiete`) as dependent variable.

- (a) Determine (with R) the predicted values \hat{y}_i and the residuals $\hat{\varepsilon}_i$. Display them in a graph by plotting \hat{y}_i against $\hat{\varepsilon}_i$.
- (b) Check with R that the following properties hold:

$$\bar{y} = \bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i \quad \text{and} \quad \bar{\hat{\varepsilon}} = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i = 0.$$

i.e., the mean of the \hat{y}_i 's equals \bar{y} and the residuals have mean 0.

- (c) Prove the equations in (b) for the simple linear regression case.

Exercise 2

Generate artificial regression data using the true model:

$$Y = 1 + 2X + X^2 + \varepsilon$$

where $X \sim N(0, 1)$ and $\varepsilon \sim N(0, 0.25)$. Choose a sufficiently large sample size, e.g. $n = 500$ and estimate the following regression models:

- simple linear regressin, i.e. regress Y on X ,
 - quadratic regression, i.e. regress Y on X and X^2 ,
 - cubic regression, i.e. regress Y on X , X^2 , and X^3 .
- (a) Determine the coefficients of determination (R^2) of your three regressions. Do the same for the RSS values. (See tables on the last page of this exercise sheet.)
 - (b) Analyse and compare the residuals for all 3 model fits (e.g. by residual plots or by boxplots). What do you conclude with respect to which of the 3 models seems to be appropriate?

Exercise 3

Use the data generating process and the models of Exercise 3 again, but now with a small sample size (e.g. $n = 10$).

- Construct (in R) the matrix \mathcal{X} for all 3 models (the matrix \mathcal{X} is called “design matrix”).
- Check, if \mathcal{X} and $\mathcal{X}^\top \mathcal{X}$ are of full rank. (Recall: What is the rank of a matrix? How could you determine this value using R?)
- Do also calculate the following matrix (the “hat matrix”) for each of the 3 models:

$$\mathbf{P} = \mathcal{X}(\mathcal{X}^\top \mathcal{X})^{-1} \mathcal{X}^\top.$$

Determine (for each model) the trace and the eigenvalues of \mathbf{P} . (Do you remember the relation between them?)

- By \mathbf{I} we denote the identity matrix (a matrix with diagonal elements 1 and 0 otherwise). Show with the help of R **and without R** that it holds:

$$\mathbf{P}^2 = \mathbf{P} \quad \text{and} \quad (\mathbf{I} - \mathbf{P})^2 = \mathbf{I} - \mathbf{P}$$

Exercise 4

We use the CPS1985 dataset which is included in an R package:

```
require(AER)
data(CPS1985)
```

Check the contents of the dataset (`?CPS1985`) and consider the variables $Y = \text{wage}$, $X_1 = \text{education}$ and $X_2 = \text{experience}$. Estimate the following models:

- regression of $\log(Y)$ on X_1 ,
- regression of $\log(Y)$ on X_1 and X_2 ,
- regression of $\log(Y)$ on X_1 , X_2 and X_2^2 .

By \log we denote the natural logarithm (denoted by \ln in mathematics, the function is `log` in R as in statistics we merely use this logarithm).

- Interpret the estimated coefficients. Do they make sense? (Plot a graph for the quadratic part of the 3rd model.)
- Analyse and compare the residuals for all 3 model fits (e.g. by comparing R^2 and by residual plots or by boxplots). What do you conclude with respect to which of the 3 models seems to be appropriate?

Components of a Linear Model Estimated in R

The following R functions can be used to extract components from a linear model.

Example: `model <- lm(y ~ x); summary(model)`

| Function | Description |
|------------------------|--|
| <code>summary</code> | summary output (see also <code>?summary.lm</code>) |
| <code>coef</code> | estimated coefficients |
| <code>residuals</code> | residuals $\hat{\varepsilon}_i$ |
| <code>fitted</code> | predicted values \hat{y}_i |
| <code>predict</code> | predicted values, useful for new data (see also <code>?predict.lm</code>) |
| <code>anova</code> | test for comparing two nested models |
| <code>plot</code> | some diagnostic plots |
| <code>confint</code> | confidence intervals for the coefficients |
| <code>deviance</code> | residual sum of squares RSS |
| <code>vcov</code> | estimated covariance matrix (of the coefficients) |
| <code>logLik</code> | log-likelihood (under normality assumption) |
| <code>AIC</code> | Akaike's information criterion (for model choice) |

Further terms can be extracted from `summary`. Example: `summary(model)$call`

| Function | Description |
|----------------------------|---|
| <code>call</code> | call of <code>lm</code> |
| <code>terms</code> | information on the explanatory variables |
| <code>residuals</code> | residuals $\hat{\varepsilon}_i$ |
| <code>coefficients</code> | table of coefficients, standard errors, t values and p values |
| <code>sigma</code> | estimated standard deviation $\hat{\sigma}$ |
| <code>df</code> | degrees of freedom |
| <code>r.squared</code> | coefficient of determination R^2 |
| <code>adj.r.squared</code> | adjusted coefficient of determination |
| <code>fstatistic</code> | F statistic with acc. degrees of freedom |
| <code>cov.unscaled</code> | unscaled covariance matrix (results in <code>vcov</code> when multiplied with $\hat{\sigma}^2$) |