

Series 7

1. **Collinearity and variable selection:** In a study about infection risk controlling in US hospitals a random sample from 113 hospitals contains the following variables:

id	randomly assigned ID of the hospital
length	average duration of hospital stay (in days)
age	average age of patients (in years)
inf	average infection risk (in percent)
cult	number of bacteriological tests per asymptomatic patient x 100
xray	number of X-rays per asymptomatic patient x 100
beds	number of beds
school	university hospital 1=yes 2=no
region	geographical region 1=NE 2=N 3=S 4=W
pat	average number of patients a day
nurs	number of full-employed, trained nurses
serv	percentage of available services from a fixed list of 35 references

Read in the data from: <http://stat.ethz.ch/Teaching/Datasets/senic.dat>. Since some observations span more than a single line, you have to use `scan()` to read the file into R:

```
senic <-scan("http://stat.ethz.ch/Teaching/Datasets/senic.dat",
  what=list(id=0,length=0,age=0,inf=0,cult=0,xray=0,beds=0,school=0,
  region=0,pat=0,nurs=0,serv=0))
```

Using `senic <- data.frame(senic); senic <- senic[, -1]` you turn the object into a user friendly data frame structure. Turn the variables `school` and `region` into so-called factor variables.

Using the variables `age`, `inf`, `region`, `beds`, `pat`, `nurs` as predictors and `length` as response variable, perform a linear regression analysis and find an optimal model by following the next instructions:

- Check the correlations between these variables. Which of them are problematic and why? Is there an intuitive explanation of this problem? Combine some of the predictors to improve the situation.
- Perform the necessary transformations on the predictors and response.
- Fit a linear regression using the transformed variables. Then, use this model as your starting equation to do backward elimination (using p-values).
- Perform a backward elimination using the AIC criterion. Use the function `step()`. Check the final model with the usual diagnostic plots.
- Now perform a forward selection using the AIC criterion. Thus, start with the empty model. Use the same function as before. Check also the diagnostic plots and comment on the differences with **c)** and **d)**.
- Optional:** Perform a stepwise selection. Start with the full model as well as with empty model and compare the results. Check the help file of `step()` on how to perform a stepwise selection.

- 2. Cross validation:** The goal of this exercise is to make you acquainted with the cross-validation technique. Use the data set `data(houseprices)` from the package `library(DAAG)`.

```
> head(houseprices)
```

```
   area bedrooms sale.price
9   694         4    192.0
10  905         4    215.0
11  802         4    215.0
12 1366         4    274.0
13  716         4    112.7
14  963         4    185.0
```

- a) Perform a leave-one-out cross validation for the model containing both predictors as main effects:
`sale.price ~ area + bedrooms`
 Is there a better model to predict the sale price? What other models are possible anyway? R hint: Use the R-function `CVlm()` from `library(DAAG)`.
- b) **Optional exercise for advanced users:** Instead of using the function `CVlm(data, formula, fold.number, ...)` you could also perform the cross validation “by hand” using a `for-loop`.

3. Logistic Regression for Binary Data

A car manufacturer instructed a market research company to analyze which families are going to buy a new car next year using a logistic regression model. Data stems from a random sample of 33 families from an agglomeration area. Assessed variables cover the yearly household income (in 1000 US \$) and the age of the oldest car in the family (in years). 12 months later, interviewers assessed which families had bought a new car in the meantime. The data is available in the file `car.dat` and can be read in with following command.

```
read.table("http://stat.ethz.ch/Teaching/Datasets/car.dat",header=T)
```

- a) Perform a logistic regression. Report the regression equation.
- b) Estimate $\exp \hat{\beta}_{income}$ and $\exp \hat{\beta}_{age}$ and interpret the values.
- c) How large is the estimated probability that a family with a yearly household income of 50 000 US \$ and whose oldest car is 3 years old will buy a new car?
- d) Do the residual plots show any abnormalities?
- e) Is the variable `age` required in the model?
- f) Is there a non-negligible interaction between `income` and `age`?

4. Logistic Regression for Binomial Data

In this task we analyze an example concerning hypertension. First, we need to enter the data. This is done as follows:

```
> no.yes <- c("No", "Yes")
> smoking <- gl(2,1,7, no.yes)
> obesity <- gl(2,2,7, no.yes)
> snoring <- gl(2,4,7, no.yes)
> n.total <- c(60, 17, 8, 187, 85, 51, 23)
> n.hyper <- c(5, 2, 1, 35, 13, 15, 8)
```

Here, the function `gl` creates a factor variable with the given levels. The factors `smoking`, `obesity` and `snoring` have an obvious meaning, `n.total` is the number of observations and `n.hyper` is the number of people with hypertension in each group.

- a) In order to fit a binomial logistic regression model construct a response matrix with two columns containing the number of people with and without hypertension, respectively.

- b) Fit a binomial regression model to the data. Does this model fit well? Assess the goodness-of-fit via the residual deviance.
- c) Which variables significantly influence the occurrence of hypertension?
- d) Try to find a suitable model using likelihood-ratio tests.
- e) Compare the observed and fitted proportions for hypertension using the model you found in d). Additionally, calculate the expected and observed counts.

Preliminary discussion: Monday, December 08.

Deadline: —.

Question hour: Thursday, January 15: 14:00 – 15:00, HG G 26.3 .