



Project 2: BMI survey II

Formalities, structure and expectations – second mandatory project

In this project, we'll continue to analyse data from the BMI survey. Here, we'll formulate and select a suitable multiple linear regression model for BMI. The assignment must be solved using the statistical software R. Some code suggestions are provided but, in addition, it's a good idea to take a look at the R code from project 1, as well as chapter 5 and 6 of the book.

The results of your analysis must be documented in a report with tables, figures, appropriate mathematical notation, and explanatory text. Relevant figures and tables must be included within the text, not in an appendix. Present the results of your analysis as you would when explaining them to one of your peers. Divide the report into subsections, one for each of the questions to be answered.

The report must be handed in as a pdf file. R code should not be included in the report itself, but must be handed in as an appendix (a .R file). The report and appendix must be handed in under Opgaver/Assignments on CampusNet at:

Assignments > Active Assignments > Obligatorisk opgave nr. 2:
BMI undersøgelse II > Answer > Answer Assignment

The report should not exceed 6 pages (excluding figures, tables, and the appendix). A page contains 2400 characters.

It's important that you describe and explain the R output in words – figures and tables cannot stand alone.

When you're asked to state a formula, insert numbers, and then perform certain computations, it's important to show that you've done this by including your intermediate

results. (In these cases, it's not enough to report results obtained directly from R). Furthermore, remember that when performing a hypothesis test, you must go through the following steps: State the hypothesis and significance level (α), compute the test statistic and state its distribution, compute the p -value, and summarize your findings.

Figures and tables are not included in the assessment of the length of the report. However, it's not in itself an advantage to include many figures, if they aren't relevant!

You may work in groups, but the report must be written individually. Questions may be addressed to the teaching assistants, see the guidelines on the *Projects* page of the course website.

Data

Read the dataset `bmi2_data.csv` into R. The following code may be used:

```
# Read the dataset 'bmi2_data.csv' into R
D <- read.table("bmi2_data.csv", header = TRUE, sep = ";")
```

The dataset for this project contains observations of the following variables:

- `id`: The respondent's id number (may be used for identification)
- `bmi`: The respondent's BMI (in kg/m^2)
- `age`: The respondent's age (in years)
- `fastfood`: No. of days per year that the respondent dines at fast-food restaurants

The `fastfood` variable was originally a categorical variable, but it was re-coded (see project 1), so that it may now be used as a numerical variable. In the following, we will refer to it as "*fast-food consumption*".

Before you proceed, add the following log-transformed BMI variable to the dataset:

```
# Add log-BMI to the dataset
D$logbmi <- log(D$bmi)
```

Statistical analysis

- a) Present a short descriptive analysis and summary of the data for the variables `logbmi`, `age`, and `fastfood`. Include scatter plots of the log-transformed BMI scores against the two other variables, as well as histograms and box plots of all three variables. Present a table containing summary statistics, which includes the number of observations, and the sample mean, standard deviation, median, and 0.25 and 0.75 quantiles for each variable.

In the following, the statistical model should only be fitted to the first 840 observations in the dataset. Later on, we'll use the last seven observations to evaluate the prediction capabilities of the final model. For example, the following code may be used to split the dataset into two parts, one for estimating the model (`D_model`), and the other for validating prediction accuracy (`D_test`):

```
# Subset containing the first 840 observations (for model estimation)
D_model <- subset(D, id <= 840)

# Subset containing the last 7 observations (for validation)
D_test <- subset(D, id >= 841)
```

- b) Formulate a multiple linear regression model with the log-transformed BMI scores as the dependent/outcome variable (Y_i), and age and fast-food consumption as the independent/explanatory variables ($x_{1,i}$ and $x_{2,i}$, respectively). Remember to state the model assumptions. (See Equation (6-1) and Example 6.1).
- c) Estimate the parameters of the model. These consist of the regression coefficients, which we denote by β_0 , β_1 , β_2 , and the variance of the residuals, σ^2 . You may use the following R code:

```
# Estimate multiple linear regression model
fit <- lm(logbmi ~ age + fastfood, data = D_model)

# Show parameter estimates etc.
summary(fit)
```

Give an interpretation of the estimates $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$, explaining what they tell us about the relation between the log-transformed BMI scores and the model's explanatory variables. (See Remark 6.14). Furthermore, present the estimated standard deviations of $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$, the degrees of freedom used for the estimated residual variance $\hat{\sigma}^2$, and the explained variation, R^2 .

- d) Perform model validation with the purpose of assessing whether the model assumptions hold. Use the plots, which can be made using the R code below, as a starting point for your assessment. (See section 6.4 on residual analysis).

```
# Plots for model validation

# Observations against fitted values
plot(fit$fitted.values, D_model$logbmi, xlab = "Fitted values",
     ylab = "log(BMI)")

# Residuals against each of the explanatory variables
plot(D_model$EXPLANATORY_VARIABLE, fit$residuals,
     xlab = "INSERT TEXT", ylab = "Residuals")

# Residuals against fitted values
plot(fit$fitted.values, fit$residuals, xlab = "Fitted values",
     ylab = "Residuals")

# Normal QQ-plot of the residuals
qqnorm(fit$residuals, ylab = "Residuals", xlab = "Z-scores",
      main = "")
qqline(fit$residuals)
```

- e) State the formula for a 95% confidence interval for the age coefficient, here denoted by β_1 . (See Method 6.5). Insert numbers into the formula, and compute the confidence interval. Use the R code below to check your result, and to determine confidence intervals for the two other regression coefficients.

```
# Confidence intervals for the model coefficients
confint(fit, level = 0.95)
```

- f) It is of interest whether β_1 might be 0.001. Formulate the corresponding hypothesis. Use the significance level $\alpha = 0.05$. State the formula for the relevant test statistic (see Method 6.4), insert numbers, and compute the test statistic. State the distribution of the test statistic (including the degrees of freedom), compute the p -value, and write a conclusion.
- g) Use backward selection to investigate whether the model can be reduced. (See Example 6.13). Remember to estimate the model again, if it can be reduced. State the final model, including estimates of its parameters.
- h) Use your final model from the previous question as a starting point. Determine predictions and 95% prediction intervals for the log-transformed BMI scores, for each of the seven observations in the validation set (`D_test`). See Example 6.8, Method 6.9 and the R code below. Compare the predictions to the observed log-BMI scores for the seven observations in the validation set and make an assessment of the prediction capabilities of the final model.

```
# Predictions and 95% prediction intervals
pred <- predict(FINAL_MODEL, newdata = D_test,
               interval = "prediction", level = 0.95)

# Observed values and predictions
cbind(id = D_test$id, logbmi = D_test$logbmi, pred)
```

Hence, don't write the formulas in the report, but instead refer to that the R function `predict` was used for the calculations. The formulas requires a matrix formulation, which are out of the curriculum (to derive the formulas use Equations (6-48) and (6-49) together with the derivations leading to Equations (5-57) and (5-58)).