

Written examination: 23. June 2022

Course name and number: **Introduction to Statistics (02323)**

Duration: 4 hours

Aids and facilities allowed: All

The questions were answered by

(student number)

(signature)

(table number)

This exam consists of 30 questions of the “multiple choice” type, which are divided between 8 exercises. To answer the questions, you need to fill in the “multiple choice” form (6 separate pages) on Digital Exam with the numbers of the answers that you believe to be correct.

5 points are given for a correct “multiple choice” answer, and -1 point is given for a wrong answer. ONLY the following 5 answer options are valid: 1, 2, 3, 4, or 5. If a question is left blank or an invalid answer is entered, 0 points are given for the question. Furthermore, if more than one answer option is selected for a single question, which is in fact technically possible in the online system, 0 points are given for the question. The number of points needed to obtain a specific mark or to pass the exam is ultimately determined during censoring.

The final answers should be given by filling in and submitting the form online via CampusNet. The table provided here is ONLY an emergency alternative. Remember to provide your student number if you do hand in on paper.

| | | | | | | | | | | |
|-----------------|-----|-----|-----|-----|-----|------|------|------|------|-------|
| Exercise | I.1 | I.2 | I.3 | I.4 | I.5 | II.1 | II.2 | II.3 | II.4 | III.1 |
| Question | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) | (9) | (10) |
| Answer | | | | | | | | | | |

| | | | | | | | | | | |
|-----------------|-------|-------|------|------|------|------|------|------|------|------|
| Exercise | III.2 | III.3 | IV.1 | IV.2 | IV.3 | IV.4 | IV.5 | V.1 | V.2 | VI.1 |
| Question | (11) | (12) | (13) | (14) | (15) | (16) | (17) | (18) | (19) | (20) |
| Answer | | | | | | | | | | |

| | | | | | | | | | | |
|-----------------|------|------|------|------|------|-------|-------|-------|--------|--------|
| Exercise | VI.2 | VI.3 | VI.4 | VI.5 | VI.6 | VII.1 | VII.2 | VII.3 | VIII.1 | VIII.2 |
| Question | (21) | (22) | (23) | (24) | (25) | (26) | (27) | (28) | (29) | (30) |
| Answer | | | | | | | | | | |

The exam paper contains 24 pages.

Continue on page 2

Multiple choice questions: *Note that in each question, one and only one of the answer options is correct. Furthermore, not all the suggested answers are necessarily meaningful. Always remember to round your own result to the number of decimals given in the answer options before you choose your answer. Also remember that there may be slight discrepancies between the result of the book's formulas and corresponding built-in functions in R.*

Exercise I

The birth weight of 50 newborn girls has been recorded in an unknown country, and the sample mean and standard deviation were found to be $\bar{x}_p = 3505.7$ g and $s_p = 467.9$ g.

Question I.1 (1)

What is the 95% confidence interval for the mean birth weight of girls, μ_p ?

- 1 [3328.3, 3683.0]
- 2 [3372.7, 3638.7]
- 3 [3371.4, 3640.0]
- 4 [3328.4, 3683.0]
- 5 [3499.6, 3511.8]

Question I.2 (2)

The mean birth weight of Danish girls is known to be 3449 g. We want to test if the birth weight from the unknown country differs significantly from the birth weight of Danish girls, so we test the hypothesis $H_0 : \mu_p = 3449$ g using the observed data.

The corresponding test statistic is computed as $t_{\text{obs}} = 0.857$. Which of the following statements is correct, when we use a significance level of $\alpha = 0.05$ (both p -value and conclusion must be correct)?

- 1 p -value = 0.198, and the hypothesis cannot be rejected.
- 2 p -value = 0.198, and the hypothesis is rejected.
- 3 p -value = 0.396, and the hypothesis cannot be rejected.
- 4 p -value = 1.604, and the hypothesis is rejected.
- 5 p -value = 0.396, and the hypothesis is rejected.

Continue on page 3

The birth weight of 50 newborn boys has also been recorded (in the same country). The sample mean and standard deviation were found to be $\bar{x}_d = 3619.4\text{g}$ and $s_d = 409.0\text{g}$. We want to test the hypothesis that girls and boys have the same mean birth weight against the alternative hypothesis that the means are different. We use a significance level of $\alpha = 0.05$ in the remainder of the questions.

Question I.3 (3)

Under the null-hypothesis of no difference between the mean birth weight of girls and boys, the (Welch) two-sample statistic, T , follows a t -distribution with ν degrees of freedom. What is ν in our case equal to?

- 1 100
- 2 98.24
- 3 49
- 4 98
- 5 96.28

Continue on page 4

Question I.4 (4)

Assume that the number of degrees of freedom, ν , is stored in R as v . Which command results in the correct critical value in the t -distribution, mentioned in the previous question, to be used for the hypothesis test of equal means?

1 `qt(0.975, v)`

2 `1-pt(0.975, v)`

3 `pt(0.95, v)`

4 `1-qt(0.95, v)`

5 `qt(0.95, v)`

Continue on page 5

Question I.5 (5)

The sampled birth weights of girls and boys are stored in `xp` and `xd`, respectively. Which of the commands below would generate the correct confidence interval for the difference in means?

1 `t.test(xp, xd, paired = TRUE)`

2 `t.test(xp, xd, paired = TRUE, conf.level = 0.90)`

3 `t.test(xp, xd, conf.level = 0.90)`

4 `t.test(xp, xd, paired = TRUE, conf.level = 0.95)`

5 `t.test(xp, xd)`

Continue on page 6

Exercise II

A Danish company wants to investigate whether the employees' professional training on a virtual reality (VR) platform affect their task quality score. 200 employees participated. The following count data provides an overview of task quality (poor, medium and good) versus VR training engagement level (below average, average, and above average).

| Task Quality Score \ VR engagement | Below Average | Average | Above Average | Row Total |
|------------------------------------|---------------|---------|---------------|-----------|
| Poor | 11 | 27 | 15 | 53 |
| Medium | 14 | 40 | 30 | 84 |
| Good | 5 | 23 | 35 | 63 |
| Column Total | 30 | 90 | 80 | 200 |

The null-hypothesis of independence between task quality score and VR training engagement score is to be tested by χ^2 -test.

Employees with "medium" and "good" task quality score are considered as "Efficient Employees".

Question II.1 (6)

What is the expected number of individuals with below average VR training engagement score and poor task quality score under H_0 (i.e assuming H_0 is true)?

- 1 7.95
- 2 21.83
- 3 25.2
- 4 19.43
- 5 9.45

Question II.2 (7)

What is the 95% confidence interval for the proportion of "Efficient Employees" based on the data given above?

- 1 [0.674, 0.796]
- 2 [0.621, 0.749]
- 3 [0.532, 0.668]
- 4 [0.426, 0.578]

5 [0.706, 0.824]

Question II.3 (8)

What is the 95% confidence interval for the difference in the proportion of good task quality scorers with VR training engagement score "above average" and "average" ($p_{\text{AbAvg, good}} - p_{\text{Avg, good}}$)?

1 [0.019, 0.361]

2 [0.043, 0.212]

3 [0.011, 0.313]

4 [0.041, 0.323]

5 [0.044, 0.091]

As a help for the next question the following R-code, where `training` is the table of counts, has been executed (some numbers have been replaced by letters)

```
chisq.test(training, correct = FALSE)
##
## Pearson's Chi-squared test
##
## data: training
## X-squared = 10.985, df = A, p-value = B
```

Question II.4 (9)

Considering the χ^2 -test statistic, what is the p -value and the correct conclusion using significance level $\alpha = 0.05$ (all parts of the answer must be correct)?

1 There is a significant dependence between VR training engagement and task quality, as $p\text{-value} = 0.027 < 0.05 = \alpha$

2 There is no evidence of significant dependence between VR training engagement and task quality, as $p\text{-value} = 0.027 < 0.05 = \alpha$

3 The Null-hypothesis cannot be rejected since $p\text{-value} = 0.50 > \alpha = 0.05$

4 There is a significant dependence between VR training engagement and task quality, as $p\text{-value} = 0.037 < 0.05 = \alpha$

5 There is no evidence of significant dependence between VR training engagement and task quality, as $p\text{-value} = 0.037 < 0.05 = \alpha$

Continue on page 8

Exercise III

In an experiment regarding poisoning of rats, survival time (days) for 24 rats were measured. Each rat received poison and was afterwards treated with one of four treatments, A, B, C, D. Let Y denote the logarithm of the survival time (`logt`) which is used for the analysis.

A one-way ANOVA model was fitted to the data:

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij} \text{ where } \varepsilon_{ij} \sim N(0, \sigma^2) \text{ and i.i.d.}$$

```
logt <- c(-1.02, -1.24, -0.92, -1.47, -0.08, -0.49, -0.71, 0.22, -0.82, -1.05,
          -1.17, -0.92, -0.58, 0.02, -0.34, -0.97, -1.51, -1.56, -1.2, -0.99,
          -1.47, -1.39, -1.2, -1.02)
treatments <-
  as.factor(c("A", "A", "A", "A", "B", "B", "B", "B", "C", "C", "C", "C",
             "D", "D", "D", "D", "A", "A", "B", "B", "C", "C", "D", "D"))
```

We are additionally informed that $SS(Tr) = 2.286$ and $SSE = 3.241$, and the values of the group means:

```
tapply(logt, treatments, mean)
##           A           B           C           D
## -1.2866667 -0.5416667 -1.1366667 -0.6816667
```

Question III.1 (10)

What is the estimate of the effect of treatment B, $\hat{\alpha}_B$?

- 1 -0.542
- 2 0.370
- 3 0.542
- 4 2.22
- 5 2.33

Continue on page 9

Question III.2 (11)

What is the value of the usual test statistic (F), for testing difference in treatments?

- 1 0.0121
- 2 0.7051
- 3 4.702
- 4 14.11
- 5 16.93

The researchers are particularly interested in comparing treatments B and D, as their knowledge in chemistry predicts that these treatments should be roughly equally good.

Question III.3 (12)

What is the conclusion on a 5% significance level regarding the post hoc difference in means between treatments B and D ($\alpha_D - \alpha_B$) (both argument and conclusion must be correct)?

- 1 The 95% confidence interval for the difference in means is $[-0.902, 0.622]$. Hence the treatments are not significantly different.
- 2 The 95% confidence interval for the difference in means is $[-0.659, 0.338]$. Hence the treatments are not significantly different.
- 3 The 95% confidence interval for the difference in means is $[-0.659, 0.338]$. Hence the treatments are significantly different.
- 4 The 95% confidence interval for the difference in means is $[-0.625, 0.345]$. Hence the treatments are not significantly different.
- 5 The 95% confidence interval for the difference in means is $[-0.206, -0.020]$. Hence the treatments are significantly different.

Continue on page 10

Exercise IV

In an office building, the duration of a room being available (i.e. empty) was measured during a period of approximately 8 months, there is one observation every time the room change from occupied to empty (this imply that there might be more than one observation pr. day). The duration of availability is measured as available hours during normal office hours, and the measurements were stored in the vector `time`. A summary of the measured duration of availability is given below.

```
summary(time)
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.2500  0.9375  2.5000  2.7017  3.7500 11.7500
```

Question IV.1 (13)

What is the Inter Quartile Range (IQR) for the presented data?

- 1 0.20
- 2 2.70
- 3 2.50
- 4 2.81
- 5 11.5

Assume that a random variable X follows an exponential distribution with expected value equal to the observed average of the time of availability.

Question IV.2 (14)

What is the median of X ?

- 1 2.50
- 2 0.169
- 3 1.87
- 4 2.70
- 5 0.741

Continue on page 11

Question IV.3 (15)

Still assuming the exponential distribution with mean equal to the observed average time of availability, which of the following pieces R-code calculate a 95% parametric bootstrap confidence interval for the expected value of the time of availability (in all cases $n = \text{length}(\text{time})$ and $k=10^4$)?

1

```
m <- mean(time)
X <- matrix(rexp(n * k, m), ncol=k)
quantile(apply(X, 2, mean), prob = c(0.025, 0.975))
```

2

```
m <- mean(time)
X <- matrix(rexp(n * k, 1/m), ncol=k)
quantile(apply(X, 2, mean), prob = c(0.025, 0.975))
```

3

```
X <- replicate(k, sample(time, replace = TRUE))
quantile(apply(X, 2, mean), prob = c(0.025, 0.975))
```

4

```
X <- replicate(n, sample(time, replace = TRUE, size = k))
quantile(apply(X, 2, mean), prob = c(0.025, 0.975))
```

5

```
m <- mean(time)
X <- matrix(rnorm(n * k, m, sd(time)), n)
quantile(apply(X, 2, mean), prob = c(0.025, 0.975))
```

Continue on page 12

It is of interest to examine if the coefficient of variation is equal 1, for that purpose the following R-code have been evaluated (including the results):

```
k <- 1e4
n <- length(time)

X <- replicate(k, sample(time, replace = TRUE))
quantile(apply(X, 2, sd) / apply(X, 2, mean),
         prob = c(0.025, 0.05, 0.95, 0.975))

##      2.5%      5%      95%      97.5%
## 0.7586455 0.7727351 0.9228017 0.9382224

quantile(apply(X, 2, var) / apply(X, 2, mean),
         prob = c(0.025, 0.05, 0.95, 0.975))

##      2.5%      5%      95%      97.5%
## 1.508020 1.570758 2.324900 2.400079

X2 <- replicate(k, rexp(n, m))
quantile( apply(X2, 2, sd) / apply(X2, 2, mean) ,
         prob = c(0.025, 0.05, 0.95, 0.975))

##      2.5%      5%      95%      97.5%
## 0.8767490 0.8933601 1.1094034 1.1368452

quantile(apply(X2, 2, var) / apply(X2, 2, mean),
         prob = c(0.025, 0.05, 0.95, 0.975))

##      2.5%      5%      95%      97.5%
## 0.2740744 0.2873999 0.4669260 0.4916159
```

Question IV.4 (16)

Based on the R-code above what can we conclude using significance level $\alpha = 0.05$, and not using any distribution assumption (both conclusion and argument should be correct)?

- 1 It cannot be rejected that the coefficient of variation is equal to 1, since $1 > 0.94$
- 2 The coefficient of variation is less than 0.7 since a 95% confidence interval is $[0.27, 0.49]$
- 3 It cannot be rejected that the coefficient of variation is equal to 1, since $1 \notin [1.51, 2.4]$
- 4 It cannot be rejected that the coefficient of variation is equal to 1, since $1 \in [0.88, 1.14]$
- 5 The coefficient of variation is not equal to 1, since $1 \notin [0.76, 0.94]$

Continue on page 13

A similar set of measurements was taken from another room, it is desired to compare the mean time of availability between the two rooms. The summary for the data from the second room is given below.

```
summary(time2)
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.250  1.500   2.500   2.551   3.250  14.250
```

Assuming independence between the rooms, it was decided to test if there is a significant difference in the expected times of availability between the rooms using a test that does not use any distribution assumptions.

Question IV.5 (17)

Which of the following pieces of R-code can be used to test the hypothesis that there is no difference between the mean time of availability for the two rooms.

1 `t.test(time, time2, paired = TRUE)`

2 `prop.test(sim1,sim2)`

3 `sim <- replicate(k, sample(time - time2, replace = TRUE))`
`quantile(apply(sim1, 2, mean), prob = c(0.025, 0.975))`

4 `sim1 <- replicate(k, sample(time, replace = TRUE))`
`sim2 <- replicate(k, sample(time2, replace = TRUE))`
`quantile(apply(sim1,2,mean) - apply(sim2,2,mean), prob = c(0.025,0.975))`

5 `t.test(time, time2)`

Continue on page 14

Exercise V

A sample was randomly taken from a population and entered into R by:

```
x <- c(8.3, 10.5, 10.6, 6.7, 10.9, 10.2, 6.3, 7.7)
```

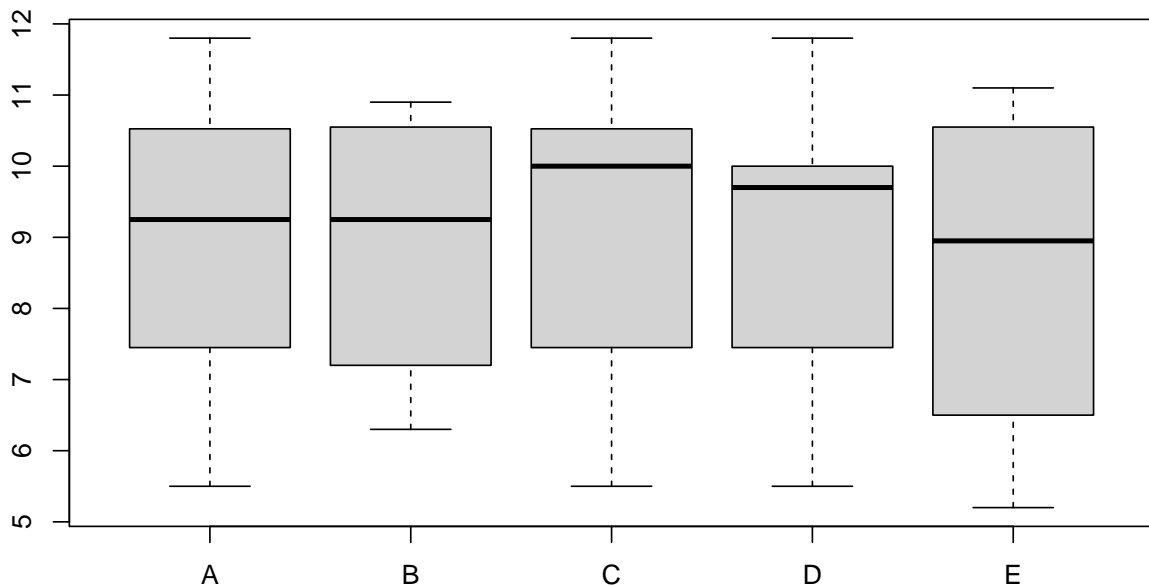
Question V.1 (18)

What is the sample mean?

- 1 1.87
- 2 3.51
- 3 7.38
- 4 8.90
- 5 9.25

Question V.2 (19)

The figure below contains boxplots of five samples:



Continue on page 15

Which of these boxplots correspond to the sample given in the exercise, hence of 'x'?

1 A

2 B

3 C

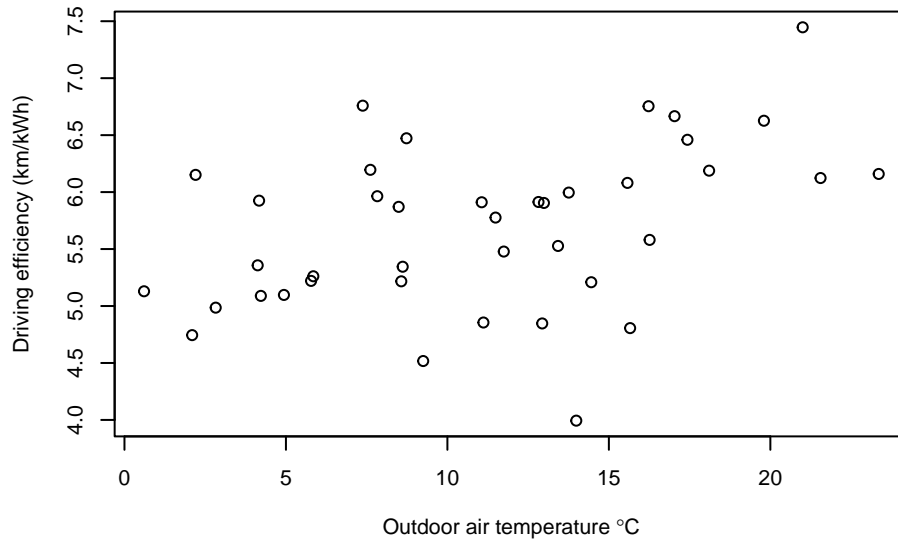
4 D

5 E

Continue on page 16

Exercise VI

The owner of an electrical car wanted to find out what effect the ambient temperature has on the driving range. So she collected the driving efficiency (trip length per unit of energy), as well as the outdoor air temperature, on every trip she made during a period. The data can be seen in the scatter plot below:



A simple linear regression model with the driving efficiency as model output, and outdoor air temperature as model input, was fitted. The results were:

```
##  
## Call:  
## lm(formula = Efficiency ~ Toutdoor)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -1.84694 -0.27181  0.01402  0.43993  1.26562   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  5.10493    0.22551  22.638  <2e-16 ***   
## Toutdoor     0.05259    0.01799   2.924  0.0058 **    
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.6573 on 38 degrees of freedom  
## Multiple R-squared:  0.1836, Adjusted R-squared:  0.1622   
## F-statistic: 8.548 on 1 and 38 DF,  p-value: 0.0058
```


Question VI.1 (20)

Which of the following statements is correct (both the conclusion and the argument must be correct)?

- 1 At significance level of $\alpha = 0.01$ a significant correlation between the driving efficiency and the outdoor temperature could be detected, since $0.0058 < 0.01$.
- 2 At significance level of $\alpha = 0.01$ a significant correlation between the driving efficiency and the outdoor temperature could not be detected, since $0.05259 > 0.01$.
- 3 At significance level of $\alpha = 0.05$ a significant correlation between the driving efficiency and the outdoor temperature could be detected, since $0.01799 < 0.05$.
- 4 At significance level of $\alpha = 0.05$ a significant correlation between the driving efficiency and the outdoor temperature could not be detected, since $0.01799 < 0.05$.
- 5 At significance level of $\alpha = 0.05$ a significant correlation between the driving efficiency and the outdoor temperature could not be detected, since $0.6573 > 0.05$.

Question VI.2 (21)

The battery size was 54 kWh. How long is the predicted mean driving range at a temperature level of 5 °C according to the model and the estimated parameters?

- 1 250 km
- 2 260 km
- 3 270 km
- 4 280 km
- 5 290 km

Continue on page 18

Question VI.3 (22)

The $i = 5$ data point had the observation of temperature at 2.096 °C and at a driving efficiency of 4.744 km/kWh. What is the residual (i.e. the realized error) for this data point?

- 1 -0.471
- 2 0.226
- 3 0.657
- 4 0.843
- 5 1.634

Question VI.4 (23)

The car owner wanted to investigate the effect of humidity on the driving range and therefore got hold of observations of air humidity from a nearby weather station and matched them with her observations.

She wanted to fit a multiple linear regression model with both the temperature and the humidity as inputs, but before she did some considerations. Which of the following statements about fitting a multiple linear regression model is not correct?

- 1 It's most often a good idea to investigate scatter plots of all possible pairs of the variables (a pairs plot in R).
- 2 It's important to carry out a model selection.
- 3 The level of correlation between the inputs cannot impact the results.
- 4 The number of observations impacts the results.
- 5 It's important to carry out a model validation with the selected model.

Continue on page 19

Question VI.5 (24)

She fitted a multiple linear regression model with both the outdoor air temperature and the humidity. The obtained result was:

```
##
## Call:
## lm(formula = Efficiency ~ Toutdoor + Humidity)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.76776 -0.34382 -0.01327  0.38670  1.34920
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.96398    0.24226  20.490 < 2e-16 ***
## Toutdoor     0.06454    0.01952   3.306  0.00211 **
## Humidity    -0.16622    0.11379  -1.461  0.15250
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6477 on 37 degrees of freedom
## Multiple R-squared:  0.2282, Adjusted R-squared:  0.1864
## F-statistic: 5.469 on 2 and 37 DF,  p-value: 0.008302
```

What is the conclusion of a backward selection step on a 5% significance level for the fitted multiple linear regression model from this result (both the conclusion and the argument must be correct)?

- 1 None of the two inputs should be removed from the model, since $0.2282 > 0.05$.
- 2 The outdoor temperature should be removed from the model, since $0.00211 < 0.05$.
- 3 The outdoor temperature should be removed from the model, since $0.06454 > 0.05$.
- 4 The humidity should be removed from the model, since $0.15250 > 0.05$.
- 5 The humidity should be removed from the model, since $0.16622 > 0.05$.

Continue on page 20

Question VI.6 (25)

In the model summary from the previous question it is stated that Multiple R-squared: 0.2282. Which of the following statements is correct?

- 1 There is a positive correlation between Efficiency and Toutdoor because the multiple R-squared value is positive.
- 2 0.2282% of the variance in Efficiency is explained by the model.
- 3 22.82% of the variance in Efficiency is explained by the model.
- 4 There is a positive correlation between Humidity and Efficiency because the multiple R-squared value is positive.
- 5 $(1 - 0.2282) \cdot 100 = 77.18\%$ of the variance in Efficiency is explained by the model.

Continue on page 21

Exercise VII

A family is playing an old game called Mouse. In this game 10 pieces of candy are put on a plate. One family member is the player and looks away, while the others point to 2 pieces of candy, which are then called the “mice”. The player now selects one piece at a time. If the selected one is a mouse, the turn is over and the player keeps all the pieces picked up before selecting the “mouse”. It can be assumed that the player selects the pieces completely at random.

Question VII.1 (26)

What is the probability that the player gets all 8 possible pieces?

- 1 2.2%
- 2 3.6%
- 3 5.8%
- 4 6.4%
- 5 9.2%

Question VII.2 (27)

If the player already got 5 pieces without picking the mouse, what is then the probability that the player will get all 8 possible pieces?

- 1 5%
- 2 8%
- 3 10%
- 4 16%
- 5 20%

Continue on page 22

Question VII.3 (28)

Sometimes when buying candy, it is actually passed the expiration date. For products in a particular shop it is known that there is a 20% probability that a product is expired. If 10 products are randomly selected, what is the probability that at least 2 products are expired?

1 32.2%

2 62.4%

3 72.5%

4 89.3%

5 93.1%

Continue on page 23

Exercise VIII

Question VIII.1 (29)

At a pharmaceutical company one is interested in designing an experiment which is able to detect a effect size $\mu_0 - \mu_1$. In a typical one-sample scenario we would like to detect $\mu_0 - \mu_1 = 0.3$. Assume $\alpha = 0.05$ and a population standard deviation $\sigma = 1.5$. How many observations should at least be measured in order to facilitate a statistical power of at least 80%?

- 1 We can't answer this question without additional specification of the confidence interval width.
- 2 197
- 3 393
- 4 412
- 5 257

Question VIII.2 (30)

Which statement regarding the R-output below is correct?

```
##  
##      Two-sample t test power calculation  
##  
##              n = 50  
##            delta = 2  
##             sd = 4.6  
##    sig.level = 0.05  
##      power = 0.576369  
## alternative = two.sided  
##  
## NOTE: n is number in *each* group
```

Continue on page 24

- 1 The statistical power is 58%. It is the probability of making a type II error.
- 2 The statistical power is 58%. Hence, the probability of making a type II error is approx. 42%.
- 3 If we would measure 60 observations instead of 50 in each group we would decrease the power of the statistical analysis.
- 4 The statistical power is 58%. Hence, the probability of making a type I error is approx. 42%.
- 5 The statistical power is 58%. It is the probability of making a type I error.

The exam is finished. Enjoy the summer!