

Written examination: 16 December 2018

Course name and number: **Introduction to Statistics (02402)**

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 14 exercises. To answer the questions, you need to fill in the “multiple choice” form (6 separate pages) on CampusNet 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	II.1	III.1	III.2	III.3	IV.1	IV.2	V.1	V.2	V.3
Question	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Answer	1	2	1	2	3	5	5	4	3	2

Exercise	VI.1	VI.2	VI.3	VI.4	VI.5	VII.1	VIII.1	IX.1	X.1	X.2
Question	(11)	(12)	(13)	(14)	(15)	(16)	(17)	(18)	(19)	(20)
Answer	5	4	2	5	3	1	3	5	5	1

Exercise	X.3	X.4	XI.1	XI.2	XII.1	XII.2	XIII.1	XIV.1	XIV.2	XIV.3
Question	(21)	(22)	(23)	(24)	(25)	(26)	(27)	(28)	(29)	(30)
Answer	1	4	2	5	1	4	3	3	1	2

The exam paper contains 31 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.

Exercise I

In the analysis of a single sample, 10 measurements are assumed to be independent and sampled from a normal distribution with mean μ and variance σ^2 . The sample mean is $\bar{x} = 0.57$, while the sample standard deviation is $s = 0.32$.

Question I.1 (1)

Which of the following is a standard 99% confidence interval for the theoretical standard deviation σ ?

1* [0.20, 0.73]

2 $0.57 \pm 1.96 \cdot 0.32$

3 [0.22, 0.58]

4 [0.05, 0.34]

5 [0.03, 0.53]

----- FACIT-BEGIN -----

See Method 3.19. Here, $n = 10$ and $\alpha = 0.01$, so the left and right endpoints are, respectively,

```
sqrt( (10-1)*0.32^2/qchisq(0.995, df = 9) )  
## [1] 0.1976575  
sqrt( (10-1)*0.32^2/qchisq(0.005, df = 9) )  
## [1] 0.7288361
```

which result in the interval [0.20, 0.73] when rounded to two decimals.

----- FACIT-END -----

Exercise II

We would like to determine the median of X_1/X_2 , when X_1 and X_2 are independent stochastic variables, which are both normal distributed with mean 1 and variance 1. The distribution of the ratio is not trivial; therefore we resort to simulation to determine an estimate and a confidence interval for the median of the distribution of X_1/X_2 .

Question II.1 (2)

First, 10000 medians are simulated, each being the median of 10000 ratios. We store these in R in the vector `medians`:

```
ratio <- replicate(10000, rnorm(10000, mean = 1)/rnorm(10000, mean = 1))
medians <- apply(ratio, 2, median)
```

Subsequently, the sample mean and a series of percentiles are calculated for these 10000 medians:

```
mean(medians)
## [1] 0.6193

quantile(medians, c(0.005, 0.025, 0.05, 0.5, 0.95, 0.975, 0.995), type = 2)
##   0.5%   2.5%    5%   50%   95%  97.5%  99.5%
## 0.5873 0.5949 0.5989 0.6193 0.6402 0.6443 0.6515
```

Which of the following choices yields an estimate for the median of X_1/X_2 and a 95% confidence interval for this median?

- 1 Estimate: 1.
95% confidence interval: $[1 - 1.96 \cdot 0.6193, 1 + 1.96 \cdot 0.6193]$.
- 2* Estimate: 0.6193.
95% confidence interval: $[0.5949, 0.6443]$.
- 3 Estimate: 1.
95% confidence interval: $[1 - 0.5949, 1 + 0.6443]$.
- 4 Estimate: 0.6193.
95% confidence interval: $[0.5873, 0.6515]$.
- 5 Estimate: 0.6193.
95% confidence interval: $[0.6193 - 0.5949, 0.6193 + 0.5949]$.

----- FACIT-BEGIN -----

The estimate is the average of the simulated medians, i.e. the estimated median is 0.6193. The left and right endpoints of the 95% confidence interval are, respectively, the 2.5% and 97.5% quantiles of the simulated medians, so the confidence interval becomes [0.5949, 0.6443].

----- FACIT-END -----

Exercise III

A normal distributed population has mean $\mu = 100$ and standard deviation $\sigma = 15$.

Question III.1 (3)

In a random draw, what is the probability of obtaining an observation below 90?

- 1* 0.252
- 2 0.482
- 3 0.518
- 4 0.631
- 5 0.748

----- FACIT-BEGIN -----

Let $X \sim N(100, 15^2)$. Then, we may find $P(X < 90) = P(X \leq 90)$ as:

```
pnorm(90, mean = 100, sd = 15)
## [1] 0.2524925
```

----- FACIT-END -----

Question III.2 (4)

If a random sample of $n = 10$ independent observations is drawn from the population, what is the probability that the sample mean is below 90?

- 1 0.000783
- 2* 0.0175

3 0.146

4 0.252

5 0.482

----- FACIT-BEGIN -----

Let \bar{X} denote the sample mean, i.e. the average of 10 independent random variables X_1, \dots, X_{10} , each with the same distribution as X . Use the mean and variance identities for linear combinations of independent random variables (Theorem 2.54) to compute the mean

$$E(\bar{X}) = E\left(\frac{1}{10} \sum_{i=1}^{10} X_i\right) = \frac{1}{10} \sum_{i=1}^{10} E(X_i) = \frac{1}{10} \cdot 10 \cdot \mu = \mu = 100$$

and the variance

$$V(\bar{X}) = V\left(\frac{1}{10} \sum_{i=1}^{10} X_i\right) = \frac{1}{10^2} \sum_{i=1}^{10} V(X_i) = \frac{1}{10^2} \cdot 10 \cdot \sigma^2 = \frac{1}{10} 15^2 = 22.5.$$

Now, $P(\bar{X} < 90) = P(\bar{X} \leq 90)$ may be found as

```
pnorm(90, mean = 100, sd = sqrt(22.5))
```

```
## [1] 0.01750749
```

----- FACIT-END -----

Question III.3 (5)

Suppose that a random sample of n independent observations is repeatedly drawn from the population, and that the sample variance S^2 is calculated in each repetition. What holds true for S^2 ?

1 $n^2 S^2$ is F -distributed with $n - 1$ and $n - 2$ degrees of freedom.

2 S^2 is χ^2 -distributed with $n - 1$ degrees of freedom.

3* $(n - 1)S^2/\sigma^2$ is χ^2 -distributed with $n - 1$ degrees of freedom.

4 S^2 is normal distributed with mean μ and variance σ^2/n^2 .

5 S^2 has the same distribution as $(Z - \sigma^2)/n$, where Z is standard normal distributed.

----- FACIT-BEGIN -----

See Section 3.1.6 from beginning and after one page you find the result, that the sampling distribution of the sample variance transformed by multiplying with $(n - 1)$ and dividing with σ^2 is χ^2 -distributed with $n - 1$ degrees of freedom. It is stated in Equation 3-17.

----- FACIT-END -----

Exercise IV

10 people have had their daily energy intake measured (in kJ). The measurements in the sample are shown in the table below:

Energy intake (kJ):	8230	5470	7515	5260	6390	6180	6515	6805	7515	5640
---------------------	------	------	------	------	------	------	------	------	------	------

Question IV.1 (6)

What is the median of the sample?

- 1 6390
- 2 6515
- 3 $(8230+5260)/2$
- 4 $(6390+6180)/2$
- 5* $(6390+6515)/2$

----- FACIT-BEGIN -----

You can rather easily copy from the pdf into an R script and add some commas between the values, and then:

```
# Read data into R and sort:
x <- sort(c(8230, 5470, 7515, 5260, 6390, 6180, 6515, 6805, 7515, 5640))
x
## [1] 5260 5470 5640 6180 6390 6515 6805 7515 7515 8230
```

As there are 10 observations, the median is computed as the mean of observations number 5 and 6 after sorting:

$$\frac{(6390 + 6515)}{2} = 6452.5$$

The result may be verified using R:

```
median(x)
## [1] 6452.5
```

or:

```
quantile(x, prob=0.5, type=2)
##      50%
## 6452.5
```

----- FACIT-END -----

Question IV.2 (7)

The sample mean is $\bar{x} = 6552$, while the sample standard deviation is $s = 975.94$. It is assumed that the daily energy intake may be modelled by a normal distribution, and that the observations are independent and identically distributed. What is the p -value for the t -test that tests the hypothesis that the mean daily energy intake is 7725 kJ?

- 1 0.4
- 2 0.06
- 3 0.04
- 4 0.006
- 5* 0.004

----- FACIT-BEGIN -----

See Method 3.23. With $\bar{x} = 6552$, $s = 975.94$ and $n = 10$, the observed t -test statistic is calculated as

$$t_{\text{obs}} = \frac{6552 - 7725}{975.94/\sqrt{10}} = -3.8007989$$

and the p -value is thus found as

$$2P(T \leq t_{\text{obs}}) = 2 \cdot 0.0021061 = 0.004.$$

$P(T \leq t_{\text{obs}})$ is found in R as:

```
pt(-3.8007989, 10-1)
## [1] 0.002106097
```

The in-built function could also be used:

```
t.test(x, mu=7725)
##
## One Sample t-test
##
## data: x
## t = -3.8008, df = 9, p-value = 0.004212
## alternative hypothesis: true mean is not equal to 7725
## 95 percent confidence interval:
## 5853.853 7250.147
## sample estimates:
## mean of x
## 6552
```

----- FACIT-END -----

Exercise V

A married couple visits the same restaurant several times a month. Typically, they order a glass of red wine with their food. One day, they decide to complain to the owner. They believe that one of the waiters pours less wine into the glass than what they pay for. Consequently, the owner launches an experiment with three of the restaurant's waiters in order to investigate how much they pour into wine glasses, when they pour using a rule of thumb. Each of the three waiters (here anonymized by A, B, and C) were asked to pour red wine into 20 wine glasses, after which the content in each glass was measured. The data were read into R in two variables: `waiter`, indicating which waiter poured the wine, and `wine`, indicating the amount of wine in the glass (in mL).

The following code was run in R to analyze the data:

```
anova(lm(wine ~ waiter))
## Analysis of Variance Table
##
## Response: wine
##           Df Sum Sq Mean Sq F value    Pr(>F)
## waiter     2 1043.4   521.71   6.9594 0.001976 **
```



```
## Residuals 57 4273.0 74.97
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Question V.1 (8)

What may be concluded from the R output above, when a significance level of 5% is used (both the reasoning and conclusion must be correct)?

- 1 As the observed F -test statistic is larger than the 0.95 quantile of the $F(57, 2)$ -distribution, there is a significant difference in the expected amount of wine in glasses poured by the three different waiters.
- 2 As the p -value is larger than 5%, there is no significant difference in the expected amount of wine in glasses poured by the three different waiters.
- 3 As the sum of squared errors, SSE , is more than four times the size of the treatment sum of squares, $SS(Tr)$, there is too much noise in the data for it to be meaningful to perform a one-way analysis of variance.
- 4* As the observed F -test statistic is larger than the 0.95 quantile of the $F(2, 57)$ -distribution, there is a significant difference in the expected amount of wine in glasses poured by the three different waiters.
- 5 As the p -value is less than 5%, there is no significant difference in the expected amount of wine in glasses poured by the three different waiters.

----- FACIT-BEGIN -----

See Theorem 8.6. From the R-output, it is seen that the relevant F -test statistic is $F_{\text{obs}} = 6.9594$. The 0.95 quantile of the $F(2, 57)$ distribution may be found using R:

```
qf(0.95, df1 = 2, df2 = 57)
## [1] 3.158843
```

----- FACIT-END -----

Question V.2 (9)

Among other things, the owner would like to make a comparison between waiter A (the young waiter, whom the couple complained about) and waiter B (an older waiter with many years of experience in the business). On average, waiter A poured 127 mL of wine into each glass, while

waiter B poured 135 mL. Compute the t -test statistic for the post hoc pairwise hypothesis test which compares the expected amount of wine in glasses poured by waiter A and waiter B.

- 1 $t_{obs} = -0.92$
- 2 $t_{obs} = -4.13$
- 3* $t_{obs} = -2.92$
- 4 $t_{obs} = -1.07$
- 5 $t_{obs} = -0.11$

----- FACIT-BEGIN -----

See Method 8.10. The relevant post hoc t -test statistic is computed as follows:

$$t_{obs} = \frac{(127 - 135)}{\sqrt{74.97 \left(\frac{1}{20} + \frac{1}{20}\right)}} = -2.92$$

So its like the two-sample t.test, except that the estimate of the error variance is taken from the model fitted to all the data $\hat{\sigma}^2 = MSE = \frac{SSE}{n-k}$, i.e. the pooled variance estimate.

----- FACIT-END -----

Question V.3 (10)

In addition to the information in the previous question, it is given that, on average, waiter C poured 136 mL into each glass. Compute the Bonferroni corrected LSD (“least significant difference”) value used to perform all possible pairwise comparisons between the three waiters, and determine where there are significant differences (both the LSD value and the conclusion must be correct). Use the significance level $\alpha = 5\%$.

- 1 $LSD_{0.05/3} = 7$ mL, so there is a significant difference between the expected amount of wine in glasses poured by waiters B and C, but no significant difference between waiters A and B or between waiters A and C.
- 2* $LSD_{0.05/3} = 7$ mL, so there is a significant difference between the expected amount of wine in glasses poured by waiters A and B as well as between waiters A and C, but no significant difference between waiters B and C.
- 3 $LSD_{0.05/3} = 4$ mL, so there is a significant difference between the expected amount of wine in glasses poured by waiters A and B and between waiters A and C, but no significant difference between waiters B and C.

- 4 \square $LSD_{0.05/3} = 17$ mL, so there is a significant difference between the expected amount of wine in glasses poured by waiters A and B and between waiters A and C, but no significant difference between waiters B and C.
- 5 \square $LSD_{0.05/3} = 4$ mL, so there is a significant difference between the expected amount of wine in glasses poured by waiters B and C, but no significant difference between waiters A and B or between waiters A and C.

----- FACIT-BEGIN -----

See Remark 8.13.

$$LSD_{0.05/3} = t_{1-(0.05/3)/2} \sqrt{2 \cdot MSE/20} = 2.466687 \cdot \sqrt{2 \cdot 74.97/20} = 6.8,$$

where $t_{1-(0.05/3)/2} = t_{5.95/6}$ is the $5.95/6 = 0.9916667$ quantile of the t -distribution with $60 - 3 = 57$ degrees of freedom, found in R as follows:

```
qt(5.95/6, df = 60-3)
## [1] 2.466687
```

So we can use that to determine which of the three waiters will be tested significantly different in two-sample post hoc comparisons. We have information about the average for each waiter:

$$\begin{aligned}\bar{x}_A &= 127 \text{ mL} \\ \bar{x}_B &= 135 \text{ mL} \\ \bar{x}_C &= 136 \text{ mL}\end{aligned}$$

from which we can see that A is significantly different from B and C, since their differences are higher than 7 mL, and that there is no significant difference between B and C.

----- FACIT-END -----

Exercise VI

A spring is characterized by its spring constant, k . When a spring is stretched, Hooke's law states that

$$F = -k \cdot x,$$

where x is the length (in meters) by which the spring is extended, and F is the applied force (in Newtons). The following six observations were made for a given spring:

	1	2	3	4	5	6
x	0.22	0.24	0.26	0.28	0.30	0.32
F	-0.51	-0.85	-0.89	-1.59	-1.97	-2.06

The observations were read into two vectors in R, x (length) and F (force), respectively, after which the following model was estimated:

```
model1 <- lm(F ~ x)
```

The output from `summary(model1)` is shown below, where some numbers are replaced by letters:

```
##
## Call:
## lm(formula = F ~ x)
##
## Residuals:
##      1      2      3      4      5      6
## -0.04484 -0.04146  0.25365 -0.10667 -0.15758  0.09690
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.2433     0.5483      A      C    **
## x            -16.8663     2.0148      B      D    **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1686 on 4 degrees of freedom
## Multiple R-squared:  0.946, Adjusted R-squared:  0.9325
## F-statistic: 70.08 on 1 and 4 DF,  p-value: 0.001114
```

Question VI.1 (11)

How may the statistical model corresponding to `model1` be described?

- 1 $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, where Y_i represents the length by which the spring is extended when the force x_i is applied, and $\varepsilon_1, \dots, \varepsilon_6$ are assumed to be independent and identically $N(0, \sigma^2)$ -distributed.
- 2 $Y_i = \beta_1 x_i + \varepsilon_i$, where Y_i represents the force used to extend the spring by the length x_i , and $\varepsilon_1, \dots, \varepsilon_6$ are assumed to be independent and identically $N(0, 1)$ -distributed.
- 3 $Y_i = \beta_1 x_i + \varepsilon_i$, where Y_i represents the length by which the spring is extended when the force x_i is applied, and $\varepsilon_1, \dots, \varepsilon_6$ are assumed to be independent and identically $N(0, 1)$ -distributed.
- 4 $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, where Y_i represents the length by which the spring is extended when the force x_i is applied, and $\varepsilon_1, \dots, \varepsilon_6$ are assumed to be independent and identically $N(0, 1)$ -distributed.

5* $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, where Y_i represents the force used to extend the spring by the length x_i , and $\varepsilon_1, \dots, \varepsilon_6$ are assumed to be independent and identically $N(0, \sigma^2)$ -distributed.

----- FACIT-BEGIN -----

See Chapter 5. The R-code `lm(F ~ x)` fits a simple linear regression model, in which F (The force) is the dependent variable and x (distance) is the explanatory variable. The model is defined in Equation 5-16 and some more information is given in Remark 5.6.

There are $n = 6$ observations in the sample, hence since $i = 1, \dots, n$, there are: six stochastic variables Y_i , six variables x_i and six stochastic variables ε_i . The i.i.d. assumption is that the six errors:

- all come from the same population, which is normal distributed $N(0, \sigma^2)$
- are drawn independently from the population

The assumption of independence of the errors is actually not trivial! It can be summed up in that the conditions, which leads to “unmodelled” variance in Y_i , must be varied randomly. As an example think of: if another variable (e.g. temperature) actually affected the dependent variables Y_i , and this variable is not measured and thus not included in the model, then the experiment should actually be carried out such that this variable is varied randomly. If not then the sample will be biased and eventually (some of) the conclusions drawn can be affected (estimates, p -values, ...). This basically means, that one should be very careful when designing experiments and making sure that the studied phenomena is not affected by some unmeasured non-random conditions during the experiment...

----- FACIT-END -----

Question VI.2 (12)

Based on the estimated slope in `model1`, give an estimate of the spring constant, k :

- 1 0.5483
2 3.2433
3 2.0148
4* 16.8663
5 5.2004

----- FACIT-BEGIN -----

According to Hooke's law given above, the spring constant corresponds to the estimated slope, but with the opposite sign, i.e. $\hat{k} = -\hat{\beta}_1$.

----- FACIT-END -----

Question VI.3 (13)

It is of interest to test whether the model's intercept is significantly different from zero. Give the relevant test statistic:

- 1 -8.371
- 2* 5.915
- 3 0.004
- 4 0.548
- 5 0.169

----- FACIT-BEGIN -----

The null hypothesis is $H_0 : \beta_0 = 0$, so $\beta_{0,0} = 0$ in Theorem 5.12. Using the R output (and standard notation from the book), the observed t -test statistic may be computed as

$$t_{\text{obs}} = \frac{\hat{\beta}_0}{\hat{\sigma}_{\beta_0}} = \frac{3.2433}{0.5483} = 5.915.$$

----- FACIT-END -----

Question VI.4 (14)

What is the distribution of the test statistic used to test whether the model's slope can be assumed to be zero?

- 1 A t -distribution with 6 degrees of freedom.
- 2 A standard normal distribution.
- 3 An F -distribution with 6 degrees of freedom.
- 4 A normal distribution with mean zero and standard deviation 0.1686.

5* A t -distribution with 4 degrees of freedom.

----- FACIT-BEGIN -----

See again Theorem 5.12. Here $n = 6$, so degrees of freedom is $n - 2 = 4$.

----- FACIT-END -----

Question VI.5 (15)

In a simple linear regression like the above, the estimators of the intercept and slope parameters are often correlated. When is this correlation zero?

- 1 When the standard deviation of the dependent variable is 1.
- 2 When the slope is estimated as zero.
- 3* When the average of the explanatory variable is zero.
- 4 When the standard deviation of the explanatory variable is 1.
- 5 When the average of the dependent variable is zero.

----- FACIT-BEGIN -----

According to Theorem 5.8 Equation 5-29, the covariance, and hence the correlation, between the intercept and the slope is zero if the sample mean of the explanatory variable, \bar{x} , is zero.

----- FACIT-END -----

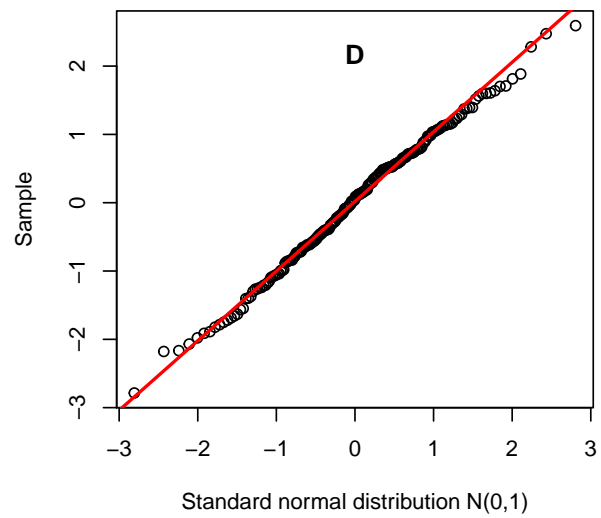
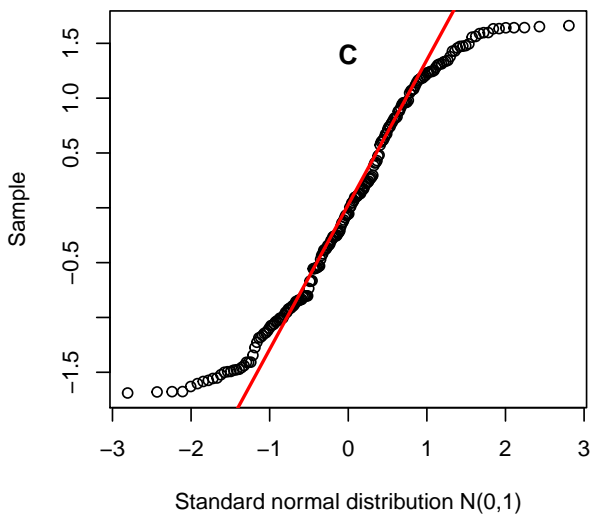
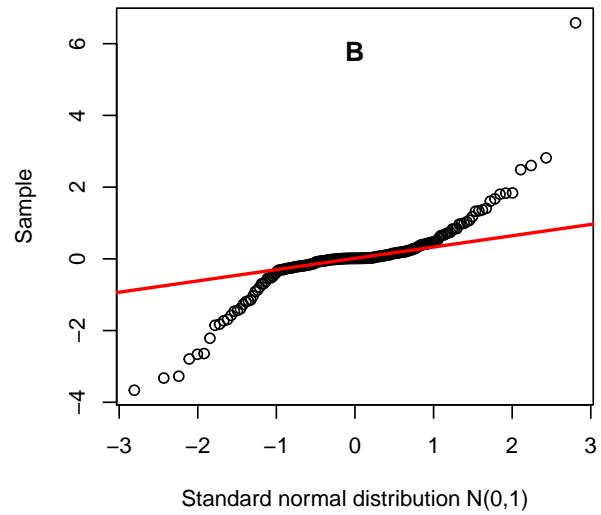
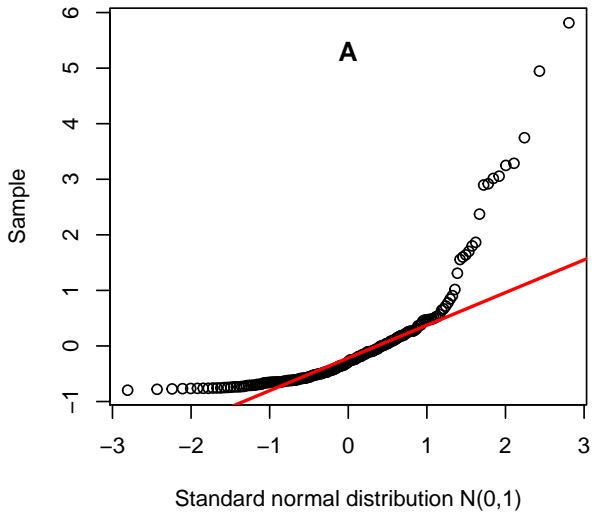
Exercise VII

In order to investigate whether data from a single sample is log-normal distributed, one could compare the data to a normal distribution using a qq-plot. If the data is log-normal distributed there will (typically) be fewer small values and more large values in the data, compared to a normal distribution with the same mean and variance as the sample.

Question VII.1 (16)

Below, four qq-plots are shown in which four different samples with mean 0 and variance 1 are each compared to a standard normal distribution. Let $z_{0.25}$ and $z_{0.75}$ denote the first and third quartile of the standard normal distribution, respectively, while $q_{0.25}$ and $q_{0.75}$ denote the first

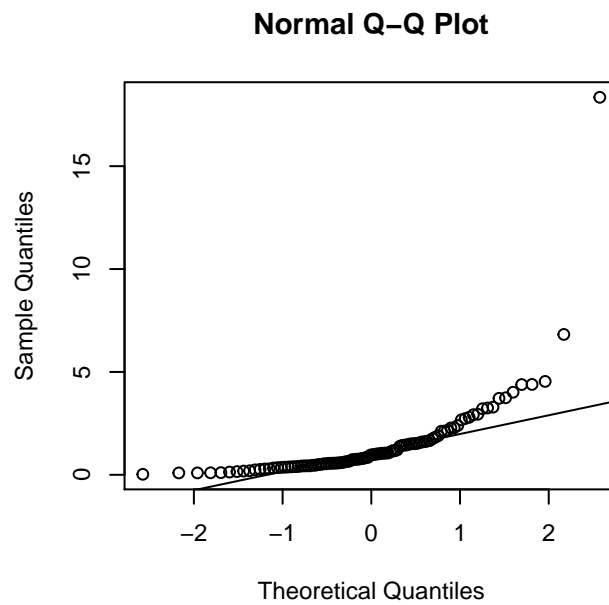
and third quartile of the sample. The red line is drawn through the points $(z_{0.25}, q_{0.25})$ and $(z_{0.75}, q_{0.75})$. Which sample fulfills the above description of log-normal distributed data?



- 1* A
- 2 B
- 3 C
- 4 D
- 5 None of the above.

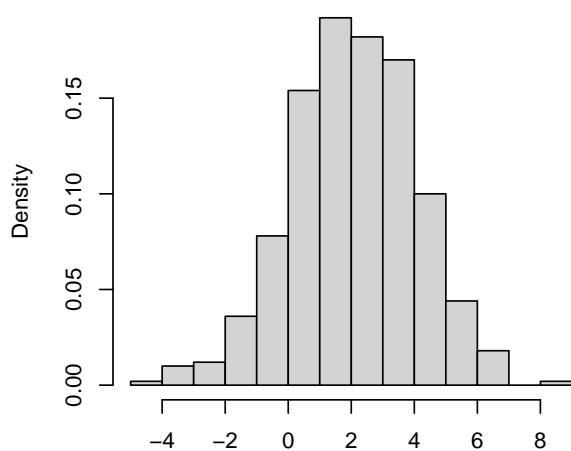
The answer is A. In B, the sample has more small values as well as more large values. The sample in C has fewer small values and fewer large values. The sample in D seems to be normal distributed. Verify the shape of a qq-plot of a log-normal distribution vs. a normal distribution in R:

```
x <- rlnorm(100)
qqnorm(x)
qqline(x)
```

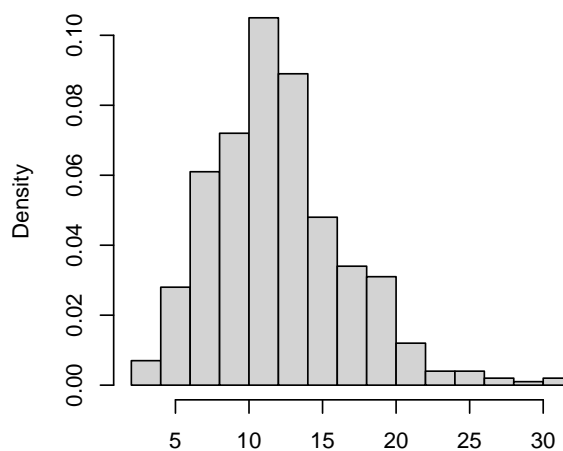


Exercise VIII

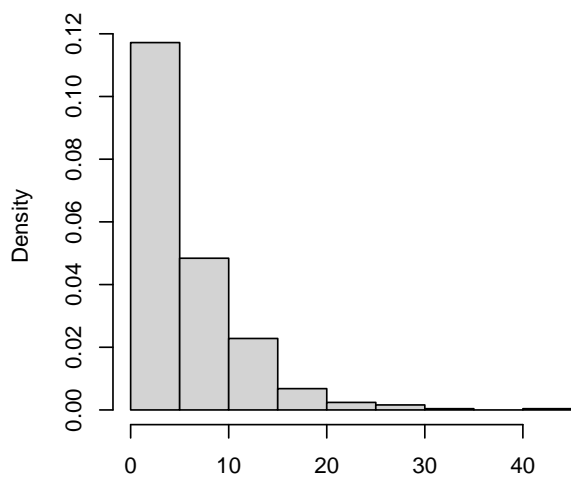
Histogram 1



Histogram 2



Histogram 3



Question VIII.1 (17)

Which distributions are simulated above? ($N(\mu, \sigma^2)$ refers to the normal distribution with mean μ and variance σ^2 , χ_a^2 to the χ^2 distribution with a degrees of freedom, and $Exp(\beta)$ to the exponential distribution with rate β).

1 1: $N(0, 4)$, 2: χ_{10}^2 , 3: $Exp(1/5)$

2 1: χ_4^2 , 2: $N(2, 4)$, 3: χ_1^2 .

3* 1: $N(2, 4)$, 2: χ_{12}^2 , 3: $Exp(1/5)$

4 1: $N(2, 4)$, 2: $Exp(5)$, 3: χ_1^2

5 1: $N(2, 4)$, 2: χ_1^2 , 3: $Exp(1/5)$

----- FACIT-BEGIN -----

The distribution in Histogram 1 takes negative values (which the χ^2 distribution doesn't), and it seems symmetric around 2, so based on the available choices, it can only be the $N(2, 4)$ distribution. The χ_1^2 distribution has mean 1, and the $Exp(5)$ distribution has mean $1/5$. Thus, based on Histogram 2 (where there isn't even any values in the data which are as small as 1), option 3 is the only possible choice. (This may be further verified by considering the means of some of the other distributions as well).

----- FACIT-END -----

Exercise IX

Two groups of rats are put on a diet while growing up, and their weight gain between day 28 and day 84 is recorded. 10 rats are put on a diet with a high protein content, while 7 rats are put on a diet with a low protein content. The collected data (weight gain in grams) is shown in the table below, with the total weight gain in each group given in the last row:

	High protein content	Low protein content
	134	70
	146	118
	104	101
	119	85
	124	107
	161	132
	107	94
	83	
	113	
	129	
Total	1220	707

Using the numbers in the table, the sample variances in the two groups are calculated to be $s_H^2 = 495$ and $s_L^2 = 425$, where H and L indicate the groups with high and low protein content, respectively. It is further given that the usual test, for whether the expected weight gain is the same for rats on high and low protein diets, has 13.7 degrees of freedom.

Question IX.1 (18)

Which of the following choices is correct (both statements need to be correct)?

- 1 Rats in the low protein diet group gain more weight than rats in the high protein diet group. However, the difference is not statistically significant at the significance level $\alpha = 0.05$.
- 2 Rats in the high protein diet group gain more weight than rats in the low protein diet group. The difference is statistically significant at the significance level $\alpha = 0.05$.
- 3 Rats in the high protein diet group gain more weight than rats in the low protein diet group. The difference is statistically significant at the significance level $\alpha = 0.01$.
- 4 Rats in the low protein diet group gain more weight than rats in the high protein diet group. The difference is statistically significant at the significance level $\alpha = 0.05$.
- 5* Rats in the high protein diet group gain more weight than rats in the low protein diet group. However, the difference is not statistically significant at the significance level $\alpha = 0.05$.

----- FACIT-BEGIN -----

The estimated difference in expected weight increase is

$$\frac{1220}{10} - \frac{707}{7} = 21,$$

that is, the increase is larger in the high protein group. However, the observed t -test statistic is

$$t_{\text{obs}} = \frac{21}{\sqrt{495/10 + 425/7}} = 2.000324$$

and as the 0.975 percentile of the t -distribution with 13.7 degrees of freedom is

```
qt(0.975, df = 13.7)
```

```
## [1] 2.149201
```

we conclude that the difference is not significant. This could also be concluded by writing in the values in R and using the `t.test()` function.

----- FACIT-END -----

Exercise X

Statistics Denmark provides data related to Denmark at www.statistikbanken.dk, among it data on traffic accidents. The following count data is taken from there:

Year Type Zone	2010				2017			
	All		Alcohol		All		Alcohol	
	City	Rural	City	Rural	City	Rural	City	Rural
Single-vehicle accidents	240	491	107	178	174	340	55	96
Others	1779	988	161	84	1456	819	106	48

Values under “all” count all accidents (including drunk-driving accidents) while numbers under “alcohol” include only drunk-driving accidents.

Question X.1 (19)

Give a 99% confidence interval for the total proportion of drunk driving accidents in 2010, where you use the relevant normal distribution approximation.

1 $0.848 \pm 2.58\sqrt{\frac{0.848}{3498}}$

2 $0.152 \pm 2.58\sqrt{\frac{0.848}{3498}}$

3 $0.848 \pm 1.96\sqrt{\frac{0.152 \cdot 0.848}{3498}}$

4 $0.848 \pm 2.58\sqrt{\frac{0.152}{3498}}$

5* $0.152 \pm 2.58\sqrt{\frac{0.152 \cdot 0.848}{3498}}$

----- FACIT-BEGIN -----

See Method 7.3. Here, $x = 107 + 178 + 161 + 84 = 530$ and $n = 240 + 491 + 1779 + 988 = 3498$, so

$$\hat{p} = \frac{530}{3498} = 0.152, \quad 1 - \hat{p} = 0.848,$$

and $z_{0.995} = 2.58$ is the 0.995 quantile of the standard normal distribution.

----- FACIT-END -----

Question X.2 (20)

Assume that the proportion of drunk-driving accidents in the “single-vehicle accidents” category is representative of the total proportion of drunk-driving. (Thus, data from the “others” category should *not* be used in this question).

Then, using the numbers from the table above and the wording from Table 3.1 of the book, what may be concluded about the difference in drunk driving between the years 2010 and 2017?

- 1* There is very strong evidence of a decrease in the proportion of drunk-driving accidents.
- 2 There is weak evidence of a decrease in the proportion of drunk-driving accidents.
- 3 There is little or no evidence of a difference in the proportion of drunk-driving accidents.
- 4 There is weak evidence of an increase in the proportion of drunk-driving accidents.
- 5 There is very strong evidence of an increase in the proportion of drunk-driving accidents.

----- FACIT-BEGIN -----

See Method 7.18. Here, $x_1 = 107 + 178 = 285$, $n_1 = 240 + 491 = 731$, $x_2 = 55 + 96 = 151$, $n_2 = 174 + 340 = 514$. The test for equality of two proportions can be tested in R using the following:

```
x1 <- 285
n1 <- 731
x2 <- 151
n2 <- 514

prop.test(c(x1,x2), c(n1,n2), correct = FALSE)

##
## 2-sample test for equality of proportions without continuity
## correction
##
## data:  c(x1, x2) out of c(n1, n2)
## X-squared = 12.249, df = 1, p-value = 0.0004656
## alternative hypothesis: two.sided
## 95 percent confidence interval:
##  0.04318181 0.14902331
## sample estimates:
##   prop 1    prop 2
## 0.3898769 0.2937743
```

The estimated proportion of alcohol-related accidents is smaller in 2017 than in 2010, and the small p -value indicates very strong evidence against the hypothesis of no change from 2010 to 2017.

----- FACIT-END -----

Question X.3 (21)

From the same source, there is also data available on the speed limits for the road stretches where the accidents occurred. The following data, describing the number of rural zone accidents at different speed limits in the years 2010 and 2017, were extracted:

	2010	2017
0 to 50 km/h	54	58
50 to 100 km/h	1280	966
100 to 130 km/h	144	135

What is the result of the usual test for no change in the distribution of accidents in the speed limit intervals between the two years (both your conclusion and reasoning must be correct)? Use the significance level $\alpha = 1\%$.

- 1* No significant difference is found in the distribution of speed limits between the two years, as the p -value is larger than the significance level.
- 2 A significant difference is found in the distribution of speed limits between the two years, as the p -value is larger than the significance level.
- 3 A significant difference is found in the distribution of speed limits between the two years, as the p -value is smaller than the significance level.
- 4 No significant difference is found in the distribution of speed limits between the two years, as the p -value is smaller than the significance level.
- 5 None of the above statements are true.

----- FACIT-BEGIN -----

Data is read into R and a χ^2 -test is carried out:

```
data <- matrix(c(54, 1280, 144, 58, 966, 135), ncol = 2)
chisq.test(data)

##
## Pearson's Chi-squared test
##
## data:  data
## X-squared = 5.8273, df = 2, p-value = 0.05428
```

It shows that the p -value for the test is above 1%, and hence a significant difference is found.

Question X.4 (22)

In connection with the usual test for whether the distribution of speed limits is the same in the two years, the following question is asked: What is the estimated proportion of accidents on roads with speed limits from 50 to 100 km/h in 2017 under the null hypothesis?

- 1 $(58 + 966 + 135)/(54 + 1280 + 144 + 58 + 966 + 135) = 0.440$
- 2 $(966)/(54 + 1280 + 144 + 58 + 966 + 135) = 0.366$
- 3 $(966)/(58 + 966 + 135) = 0.833$
- 4* $(1280 + 966)/(54 + 1280 + 144 + 58 + 966 + 135) = 0.852$
- 5 $(54 + 58 + 144 + 135)/(54 + 1280 + 144 + 58 + 966 + 135) = 0.148$

Under the null hypothesis, the proportion of accidents that happen at 50 to 100 km/h roads does not depend on the accident year, then the proportion is estimated using the data from both years: $x_{50\text{to}100\text{km/h}} = 1280 + 966 = 2246$ and $n = 54 + 1280 + 144 + 58 + 966 + 135 = 2637$. Then

$$\hat{p}_{50\text{to}100\text{km/h}} = \frac{x_{50\text{to}100\text{km/h}}}{n} = 0.852$$

Exercise XI

Below is a sample of 20 independent observations, read into R in the vector **x**:

```
x <-  
c(13, 12, 9, 7, 12, 15, 12, 10, 6, 13, 7, 13, 19, 12, 6, 4, 15, 16, 11, 18)
```

The data do not originate from a known distribution, but we are interested in the population mean and the uncertainty of its estimate.

Question XI.1 (23)

What is the sample mean \bar{x} and variance s^2 (both quantities must be correct)?

- 1 $\bar{x} = 11.2$ and $s^2 = 16.7$.
- 2* $\bar{x} = 11.5$ and $s^2 = 16.7$.
- 3 $\bar{x} = 11.2$ and $s^2 = 4.1$.
- 4 $\bar{x} = 11.5$ and $s^2 = 4.1$.
- 5 $\bar{x} = 11.5$ and $s^2 = 16.7^2$.

----- FACIT-BEGIN -----

Data can be read into R, and sample mean and variance calculated:

```
x <-
  c(13, 12, 9, 7, 12, 15, 12, 10, 6, 13, 7, 13, 19, 12, 6, 4, 15, 16, 11, 18)
mean(x)

## [1] 11.5

var(x)

## [1] 16.68421
```

----- FACIT-END -----

Question XI.2 (24)

Now, we perform a resampling of \mathbf{x} to get an idea of the uncertainty of the sample mean. We draw 200 resamples with replacement from the 20 observations in \mathbf{x} , each with sample size 20. Subsequently, the mean of each of the 200 resamples is calculated. The R code for this operation is:

```
apply(replicate(200, sample(x, replace = TRUE)), 2, mean)
```

Below, the 10 largest and 10 smallest sample means of the 200 resamples are shown.

smallest	9.00	9.65	9.65	9.80	9.90	9.95	10.00	10.00	10.00	10.05
largest	12.95	12.95	12.95	13.00	13.05	13.10	13.10	13.10	13.15	13.40

Using the results above and the book's definition of percentiles ("type = 2" in R), which of the following is a 95% bootstrap confidence interval for the population mean?

- 1 [10.05, 12.95]

2 [9.00, 13.40]

3 [9.80, 13.10]

4 [9.65, 13.10]

5* [9.925, 13.075]

----- FACIT-BEGIN -----

See Definition 1.7. For $n = 200$, and $p_1 = 0.025$, $p_2 = 0.975$, it holds that $p_1n = 5$ and $p_2n = 195$. Then, the relevant 0.25 quantile $q_{0.025}$ is the average of the 5th and 6th ordered averages, while the 0.975 quantile $q_{0.975}$ is average of the 195th and 196th ordered averages:

$$q_{0.025} = \frac{9.90 + 9.95}{2} = 9.925 \quad \text{and} \quad q_{0.975} = \frac{13.05 + 13.10}{2} = 13.075$$

----- FACIT-END -----

Exercise XII

During the preparation for a small festival, the toilet facilities are taken under consideration. Mobile toilets need to be ordered such that the capacity is sufficient, but not too high, since will lead to more cleaning and higher costs.

It is assumed that, on average, 150 guests need to use the toilets every hour, and that their arrival follows a Poisson distribution. In addition, it is assumed that each toilet can serve 20 guests per hour.

Question XII.1 (25)

Suppose that 10 toilets are ordered. What is then the probability that, in a randomly selected hour, the number of guests who arrive at the toilets exceeds the capacity?

1* 0.0042%

2 2.3%

3 11%

4 24%

5 99%

----- FACIT-BEGIN -----

Let X represent the number of guests arriving at the toilets in a randomly selected hour, then $X \sim Pois(150)$. The capacity is $10 \cdot 20 = 200$ per hour, hence we need to calculate $P(X > 200) = 1 - P(X \leq 200)$:

```
1 - ppois(200, lambda=150)
## [1] 4.205886e-05
```

----- FACIT-END -----

Question XII.2 (26)

A group of DTU students have decided to help small festivals optimize their logistical conditions. Among other things, the students have collected data on the use of toilets at small festivals. An examination of these data shows that a better model can be made to represent the number of guests who need to use the facilities in a randomly selected hour. This number can be modelled by an exponential distribution with mean $\frac{\text{“number of guests”}}{10}$, where “number of guests” is the total number of guests at the festival. In this question, this new model must be used.

A festival with 1500 guests is now considered. How many toilets should, at least, be ordered to ensure that the probability that not everyone can use the facilities is less than 2% in a randomly selected hour (given as a call in R)? It is still assumed that each toilet can serve 20 guests per hour.

- 1 `ppois(20, lambda = 15) * 20`
- 2 `qpois(0.98, lambda = 1500/10) / 20`
- 3 `qexp(0.98, rate = 10/15)`
- 4* `qexp(0.98, rate = 10/1500) / 20`
- 5 `qexp(0.98, rate = 10/1500) * 20`

----- FACIT-BEGIN -----

Let Y represent the number of guests arriving at the toilets in a randomly selected hour. Then $Y \sim Exp(10/1500)$. We need to find y such that

$$P(Y \leq 20y) \geq 0.98.$$

We can solve $P(Y \leq 20y) = 0.98$ for y by computing

```
qexp(0.98, rate = 10/1500) / 20
```

```
## [1] 29.34017
```

Thus, ordering 30 toilets or more ensures that the probability in focus stays below 2%.

----- FACIT-END -----

Exercise XIII

Below, there's a small sample with 5 independent observations:

Observations:	11.8071067	-1.7913888	-9.1872410	-4.4860901	-0.2324924
---------------	------------	------------	------------	------------	------------

Question XIII.1 (27)

Which of the following answer options is the only one that can possibly be correct?

- 1 It is impossible that the observations were sampled from a normal distribution with mean 0 and variance 10^2 .
- 2 It is possible that the observations were sampled from a uniform distribution with parameters -9 and 12.
- 3* It is possible that the observations were sampled from a t -distribution with 1 degree of freedom.
- 4 It is possible that the observations were sampled from an F -distribution with 1 and 2 degrees of freedom.
- 5 It is possible that the observations were sampled from an exponential distribution with rate 1.

----- FACIT-BEGIN -----

The F distributions and exponential distributions don't give rise to negative observations, which eliminates options 4 and 5. Likewise, the uniform distribution with parameters -9 and 12 wouldn't yield the observation -9.1872410 , which eliminates option 2. There is no reason why the observations could not come from the $N(0, 10^2)$ distribution, which eliminates option 1. We're then left with option 3: There's no reason why these observations couldn't come from a t -distribution with 1 degree of freedom (like normal distributions, t distributions take both positive and negative values).

Exercise XIV

As the share of wind power in the energy production in Denmark and the rest of Europe increases, accurate predictions become more crucial. The prediction of wind energy naturally depends on the weather forecast, but also the estimation method and model structure have an influence. The table below shows a sample of a data set giving the average weekly predicted wind production for a wind farm (measured as a percentage of installed power) for different prediction models (m1, . . . , m5).

Week	m1	m2	m3	m4	m5
1	0.6039	0.6232	0.6083	0.5751	0.6232
2	0.5143	0.5301	0.5049	0.4644	0.4850
3	0.5551	0.5603	0.5415	0.5091	0.5219
4	0.5396	0.5393	0.5766	0.4697	0.5245
⋮	⋮	⋮	⋮	⋮	⋮

Below, the data is loaded into R. The vector `prediction` contains the average weekly predictions, `model` indicates which of the five prediction models was used, and `week` indicates the week number.

One wants to investigate whether the 5 different models (m1, . . . , m5) can be assumed to give the same expected predictions (on a weekly basis) or whether there is a significant difference.

To investigate the hypothesis, the following model has been formulated

$$Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$$

where α_i describes the effect of model i , and β_j describes the effect of week j .

Question XIV.1 (28)

Under the usual assumptions, which of the following statements is correct?

- 1 $\alpha_i + \beta_j = 0$ for all combinations of i and j .
- 2 ϵ_{ij} are independent and normal distributed with mean 0 and a variance which depends on α_i and β_j .
- 3* Y_{ij} are independent and normal distributed with the same variance for all combinations of i and j .
- 4 $\sum_i \alpha_i = \sum_j \beta_j = \mu$.

5 Y_{ij} are independent and identically distributed for all combinations of i and j .

----- FACIT-BEGIN -----

The usual statistical model assumes that the Y_{ij} 's are all independent and normal distributed with the same variance, usually denoted σ^2 , but that the mean of each Y_{ij} depends on the model i and week j , i.e. $E(Y_{ij}) = \mu + \alpha_i + \beta_j$.

----- FACIT-END -----

Question XIV.2 (29)

To investigate the hypothesis that there is no difference in the expected predictions, the following R-code was run. Note that some of the results have been removed, and some numbers have been replaced by letters.

```
anova(lm(prediction ~ model+factor(week),data=dat))

## Analysis of Variance Table
##
## Response: pred
##           Df Sum Sq Mean Sq F value Pr(>F)
## model      4 0.01056 0.0026391  7.3754 5.389e-05 ***
## factor(week) 17 0.33946 0.0199684 55.8051 < 2.2e-16 ***
## Residuals   68 0.02433 0.0003578
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

How many weeks are there in the data set? (Due to rounding in the R-output, the relevant calculation will result in a decimal number which must be rounded correctly to the nearest integer).

1* 18

2 68

3 56

4 17

5 55

----- FACIT-BEGIN -----

In the row `factor(week)`, we find $MS(Week) = \frac{SS(Week)}{x-1} = 0.0199684$ and $SS(Week) = 0.33946$, where x denotes the number of weeks in the data set. Thus, isolating x gives

$$x = \frac{SS(Week)}{MS(Week)} + 1,$$

i.e.:

```
( x <- 0.33946/0.0199684 + 1 )  
## [1] 17.99986
```

which must be rounded to 18.

----- FACIT-END -----

Question XIV.3 (30)

Consider again the R-output from the previous question, and use the significance level $\alpha = 0.01$. Is there a significant difference between the five models m_1, \dots, m_5 , when the statistical model takes into account the difference between weeks (both the conclusion and argument must be correct)?

- 1 No, as $0.0106 > 0.01$.
- 2* Yes, as $5.389 \cdot 10^{-5} < 0.01$.
- 3 No, as $0.020 > 0.01$.
- 4 Yes, as $0.0026 < 0.01$.
- 5 Yes, as $0.024 > 0.01$.

----- FACIT-BEGIN -----

Yes, there is a significant difference. The relevant p value $5.389 \cdot 10^{-5}$ may be found in the `model` row of the R output, and it is smaller than the significance level 0.01.

----- FACIT-END -----

The exam paper is finished. Have a great Christmas vacation!