

Written examination: 22. June 2023

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 11 exercises. To answer the questions, you need to fill in the “multiple choice” form on exam.dtu.dk.

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. 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	II.1	II.2	III.1	III.2	IV.1	IV.2	IV.3	IV.4
Question	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Answer	2	3	1	4	5	5	2	5	1	2

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

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

The exam paper contains 40 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

A researcher is interested in comparing the average weight gain (in grams) of three different groups of mice fed with three different diets. The data is provided below.

```
group1 <- c(27, 22, 18, 26, 24)
group2 <- c(32, 22, 32, 25, 25)
group3 <- c(29, 25, 30, 30, 24)
```

Question I.1 (1)

Perform a one-way ANOVA and test the usual null hypothesis of equal treatment means at significance level $\alpha = 0.05$. Is there a significant difference in weight gain among the three groups?

- 1 The p -value is 0.03. The difference between group means is not significant because the p -value is less than 0.05.
- 2* The p -value is 0.1879. The difference between group means is not significant because the p -value is greater than 0.05.
- 3 The p -value is 0.1879. The difference between group means is significant because the p -value is greater than 0.05.
- 4 The p -value is 0.3758. The difference between group means is significant because the p -value is greater than 0.05.
- 5 The p -value is 0.03. The difference between group means is significant because the p -value is less than 0.05.

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

Data can be read into an ANOVA-suitable data frame using:

```
data <- data.frame(group = factor(rep(1:3, each = 5)),
                  weight_gain = c(group1, group2, group3))
```

The ANOVA table is

```
model <- lm(weight_gain ~ group, data = data)
anova(model)

## Analysis of Variance Table
##
## Response: weight_gain
##           Df  Sum Sq Mean Sq F value Pr(>F)
## group      2   53.733  26.867   1.9282 0.1879
## Residuals 12  167.200  13.933
```

from where we see that the p -value is 0.1879. This is larger than the significance level, therefore the difference is not significant.

One can also find this p -value manually using the formulas in Chapter 8.2.

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

Continue on page 4

Question I.2 (2)

The experiment described above was repeated (same number of mice) by a second researcher who collected a different data set. Again, one-way ANOVA was used to test for significant difference between treatment means. The following ANOVA table was obtained. Please note that some elements have been replaced by question marks.

```
## Analysis of Variance Table

## Response: weight_gain
##           Df Sum Sq Mean Sq F value Pr(>F)
## group      2  78.53  39.267   1.069 0.3739
## Residuals  ? 440.80      ?
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Which of the following statements is correct?

- 1 $Df(\text{Residuals}) = 14$ and $\text{Mean Sq}(\text{Residuals}) = 31.486$.
- 2 $Df(\text{Residuals}) = 15$ and $\text{Mean Sq}(\text{Residuals}) = 29.387$.
- 3* $Df(\text{Residuals}) = 12$ and $\text{Mean Sq}(\text{Residuals}) = 36.733$.
- 4 $Df(\text{Residuals}) = 14$ and $\text{Mean Sq}(\text{Residuals}) = 2.805$.
- 5 $Df(\text{Residuals}) = 13$ and $\text{Mean Sq}(\text{Residuals}) = 3.021$.

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

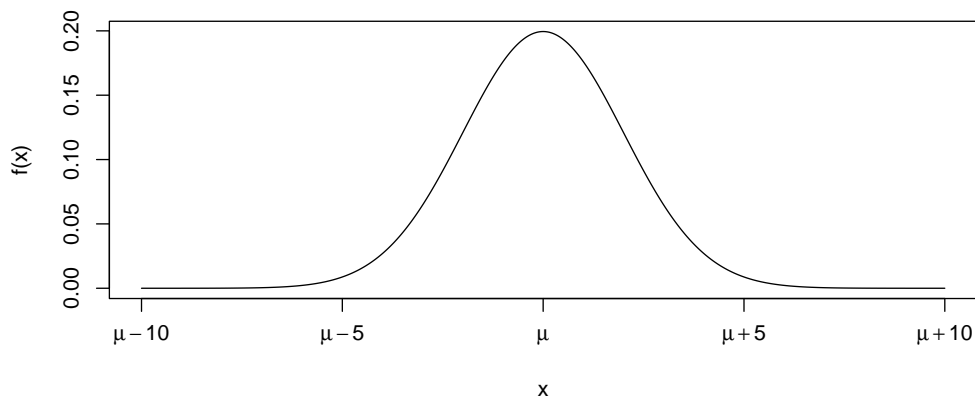
We use the formulas for the ANOVA table. n =no. of observations= 15 and k =no. of groups= 3.
Hence $Df(\text{Residuals}) = 15 - 3 = 12$ and $\text{Mean Sq}(\text{Residuals}) = 440.80/12 = 36.733$.

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

Continue on page 5

Exercise II

Let the random variable X be normal distributed with mean μ and standard deviation $\sigma = 2$, i.e. $X \sim N(\mu, 2^2)$, hence its' pdf is:



Question II.1 (3)

Let another random variable be defined by the function

$$Y_1 = a_1 + b_1 \cdot X + b_2 \cdot X$$

What is the mean and variance of Y_1 ?

- 1* $E(Y_1) = a_1 + (b_1 + b_2)\mu$ and $V(Y_1) = (b_1 + b_2)^2 \cdot 4$
2 $E(Y_1) = a_1 + b_1 + b_2$ and $V(Y_1) = b_1^2 + b_2^2$
3 $E(Y_1) = a_1 + b_1 + b_2$ and $V(Y_1) = b_1 + b_2$
4 $E(Y_1) = 0$ and $V(Y_1) = b_1^2 + b_2^2$
5 $E(Y_1) = 0$ and $V(Y_1) = b_1 + b_2$

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

We must use the rules for identities for the mean and variance in Section 2.7.

$$\begin{aligned} E(Y_1) &= E(a_1 + b_1 \cdot X + b_2 \cdot X) \\ E(Y_1) &= E(a_1) + E(b_1 \cdot X + b_2 \cdot X) \\ E(Y_1) &= E(a_1) + E((b_1 + b_2) \cdot X) \\ E(Y_1) &= E(a_1) + (b_1 + b_2) E(X) \\ E(Y_1) &= a_1 + (b_1 + b_2)\mu \end{aligned}$$

and similarly for the variance

$$\begin{aligned}V(Y_1) &= V(a_1 + b_1 \cdot X + b_2 \cdot X) \\V(Y_1) &= V(a_1) + V(b_1 \cdot X + b_2 \cdot X) \\V(Y_1) &= V(a_1) + V((b_1 + b_2) \cdot X) \\V(Y_1) &= V(a_1) + (b_1 + b_2)^2 V(X) \\V(Y_1) &= (b_1 + b_2)^2 \cdot 4\end{aligned}$$

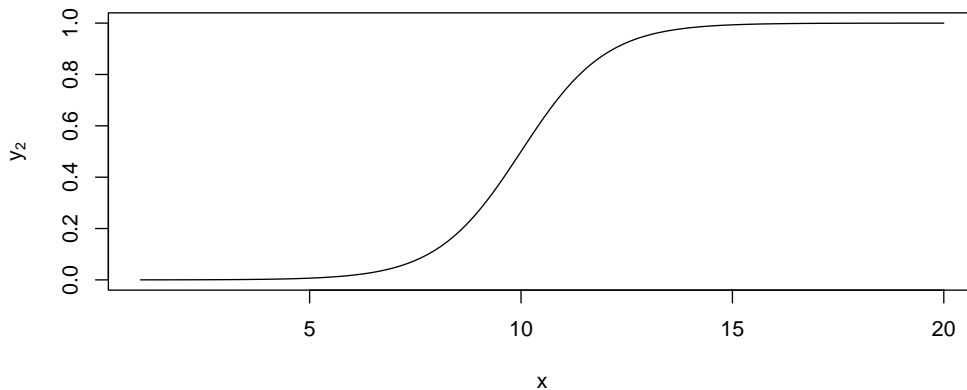
----- FACIT-END -----

Question II.2 (4)

Let another random variable be defined by the function

$$Y_2 = \frac{1}{1 + \exp(a_2 + b_2 \cdot X)}$$

where $a_2 = 10$ and $b_2 = -1$. A plot of this function is:



This function is called the logistic function (or Sigmoid function).

The notation $V(Y_2|\mu = \mu_0)$ means the variance of Y_2 when μ is equal to μ_0 . E.g. $V(Y_2|\mu = 0)$ is the variance of Y_2 when μ is equal to 0.

Which one of the following statements is correct?

- 1 $V(Y_2|\mu = 0) < V(Y_2|\mu = 10) < V(Y_2|\mu = 20)$
- 2 $V(Y_2|\mu = 0) < V(Y_2|\mu = 20) < V(Y_2|\mu = 10)$
- 3 $V(Y_2|\mu = 0) = V(Y_2|\mu = 20) = V(Y_2|\mu = 10)$

$$4^* \square V(Y_2|\mu = 0) = V(Y_2|\mu = 20) < V(Y_2|\mu = 10)$$

$$5 \square V(Y_2|\mu = 20) < V(Y_2|\mu = 10) < V(Y_2|\mu = 0)$$

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

When we want the variance of a non-linear function we must look into the techniques in Section 4.1.3. Here we can either use simulation or the approximate error propagation rule.

Using the error propagation rule in Method 4.3 we realise that the variance is scaled with the partial derivative of the function, in this case only one variable (x), so we can see the on the plot the derivative (how steep is the function at the value of x). When $x = 0$ and $x = 20$ the function is almost flat, so the variance is much lower through the function, than for $x = 10$.

One can think of the normal distribution bell on the x -axis centred at $x = 0$ and then mapped through the function to the y -axis: It will be mapped almost same values on the y -axis. For $x = 10$ the bell will be wider mapped through to the y -axis.

One can also try simulation:

```
a <- 10
b <- -1
k <- 100000
x1 <- rnorm(k, mean=0)
x2 <- rnorm(k, mean=10)
x3 <- rnorm(k, mean=20)
var(1/(1+exp(a+b*x1)))

## [1] 1.062608e-08

var(1/(1+exp(a+b*x2)))

## [1] 0.04344787

var(1/(1+exp(a+b*x3)))

## [1] 9.658104e-09
```

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

Continue on page 8

Exercise III

In a statistics class with 589 students, all students have taken an exam with two parts each lasting 2 hours. The instructors of the class are interested in evaluating whether the two exams parts have been equally difficult by comparing the mean scores in the two parts.

Question III.1 (5)

What test should the instructors apply in order to evaluate whether the mean scores of the two exam parts are equal?

- 1 A one-way ANOVA
- 2 An F -test with 2 and 589 degrees of freedom
- 3 A two-sample t -test assuming equal variances in the two groups
- 4 A two-sample t -test with a pooled variance
- 5* A paired t -test

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

Since the two samples are taken from the same students, this is a paired experiment, and the appropriate test is therefore a paired t -test. Hence, answer 5 is correct.

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

Question III.2 (6)

One of the instructors gives the same course (and the same exam) at a different university, where 240 students are enrolled in the course and subsequently take the exam. Some summary statistics concerning the exam results are given in the below table:

	University A	University B
Students	589	240
Average score	736.4	769.9
Variance of score	169.1	402.7

When calculating the 90% confidence interval, not assuming equal variances in the two groups, for the difference in mean scores, the instructor has to use a quantile from a t -distribution. The instructor has to use which quantile of the t -distribution with how many degrees of freedom?

- 1 The 10% quantile of the t -distribution with 323.93 degrees of freedom
- 2 The 90% quantile of the t -distribution with 829 degrees of freedom
- 3 The 90% quantile of the t -distribution with 323.93 degrees of freedom
- 4 The 95% quantile of the t -distribution with 829 degrees of freedom
- 5* The 95% quantile of the t -distribution with 323.93 degrees of freedom

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

The instructor applies method 3.47 and begins by calculating the degrees of freedom of the t -distribution

$$\nu = \frac{\left(\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}\right)^2}{\frac{(s_A^2/n_A)^2}{n_A-1} + \frac{(s_B^2/n_B)^2}{n_B-1}} = \frac{\left(\frac{169.1}{589} + \frac{402.7}{240}\right)^2}{\frac{(169.1/589)^2}{589-1} + \frac{(402.7/240)^2}{240-1}} = 323.9297.$$

Next, the instructor notices that a 90% confidence interval corresponds to using a significance level of $\alpha = 0.1$. Therefore, it is the $1 - \alpha/2 = 95\%$ quantile that must be used in the calculations. Thus, answer 5 is correct.

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

Continue on page 10

Exercise IV

On April 14 1912 the passenger ship Titanic hit an iceberg and sank the following day. The table below shows the number of survivors and total number of passengers distributed on different passenger categories.

Class	1st	2nd	3rd	Crew	Total
Survived	202	117	178	212	709
Total	325	285	706	885	2201

Question IV.1 (7)

Based on the table above, what is a 95% confidence interval for the probability of survival (regardless of passenger category) given the data?

- 1 [0.66, 0.70]
- 2* [0.30, 0.34]
- 3 [0.45, 0.50]
- 4 [0.46, 0.49]
- 5 [0.66, 0.69]

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

The confidence interval can be calculated by

$$\hat{p} = x/N \quad (1)$$

where x is the number of successes and N is the total, and the confidence interval can be calculated by

$$\hat{p} \pm z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{N}} \quad (2)$$

This can be calculated in R, either using the formula or by `prop.test`

```
phat <- 709/2201
phat + c(-1, 1) * qnorm(0.975) * sqrt(phat * (1 - phat) / 2201)

## [1] 0.3026043 0.3416484

prop.test(709, 2201, correct=FALSE)
```

```
##
## 1-sample proportions test without continuity correction
##
## data: 709 out of 2201, null probability 0.5
## X-squared = 278.55, df = 1, p-value < 2.2e-16
## alternative hypothesis: true p is not equal to 0.5
## 95 percent confidence interval:
## 0.3029287 0.3419437
## sample estimates:
##      p
## 0.3221263
```

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

Question IV.2 (8)

Is there a statistically significant difference in the survival probability between the crew and the 3rd class passengers, using a 5% significance level (both the argument and the conclusion should be correct)?

- 1 Yes, since the test statistics for the relevant test is -1.67
- 2 No, since the test statistics for the relevant test is 0.66
- 3 Yes, since the test statistics for the relevant test is 1.67
- 4 No, since the p -value for the relevant test is 0.41
- 5* No, since the p -value for the relevant test is 0.56

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

The test statistics for the difference can be calculated as

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1 - \hat{p})/N}} \quad (3)$$

where \hat{p} is the estimate for the probability under the null hypothesis. And the p -value is calculated based on the standard normal density, the problem can also be solved using `prop.test` or `chisq.test`;

```
s <- c(178,212)
n <- c(706,885)
```

```

p <- s/n

ph <- sum(s)/sum(n)

z <- diff(p)/sqrt(ph*(1-ph)*sum(1/n))
## p-values 2 ways
2*(1-pnorm(abs(z)))

## [1] 0.5623276

1-pchisq(z^2,1)

## [1] 0.5623276

## or chisq.test
chisq.test(cbind(s,n-s),correct=FALSE)

##
## Pearson's Chi-squared test
##
## data:  cbind(s, n - s)
## X-squared = 0.33569, df = 1, p-value = 0.5623

## or prop.test
prop.test(s,n,correct=FALSE)

##
## 2-sample test for equality of proportions without continuity correction
##
## data:  s out of n
## X-squared = 0.33569, df = 1, p-value = 0.5623
## alternative hypothesis: two.sided
## 95 percent confidence interval:
## -0.03004594  0.05519919
## sample estimates:
##   prop 1    prop 2
## 0.2521246 0.2395480

```

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

Continue on page 13

Question IV.3 (9)

Considering the entire table, what is the relevant observed test statistics (q), critical value (CV), and conclusion for a test of the hypothesis that the survival probability is the same across all classes, using significance level $\alpha = 0.05$?

- 1* $q=187.1$, $CV=7.8$, hence there is a significant difference
- 2 $q=84.37$, $CV=15.5$, hence there is a significant difference
- 3 $q=84.37$, $CV=7.8$, hence there is a significant difference
- 4 $q=187.1$, $CV=15.5$, hence there is not a significant difference
- 5 $q=84.37$, $CV=7.8$, hence there is not a significant difference

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

This is most easily solved using `chisq.test` and hence inputting the entire table, and the critical value is calculated based on the χ^2 -distribution with 3 degrees of freedom;

```
## Test statistics
s <- c(202, 117, 178,212)
n <- c(325,285, 706, 885)
tab <- rbind(s,n-s)
chisq.test(tab)

##
## Pearson's Chi-squared test
##
## data:  tab
## X-squared = 187.11, df = 3, p-value < 2.2e-16

## Critical value
qchisq(0.95,df=3)

## [1] 7.814728
```

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

Question IV.4 (10)

We wish to test if the probability of survival of 1st class passengers differs by more than 20 percentage points compared to the average of all other passengers, which of the following statements regarding that is correct (using significance level $\alpha = 0.05$)?

- 1 Since $\hat{p}_{1st} - \hat{p}_{rest} = 0.35$ there is a significant difference and it is greater than 0.2
- 2* The relevant confidence interval is $[0.29, 0.41]$, and hence the survival probability of 1st class passengers is at least 20 percentage point higher than the survival probability of other passengers
- 3 0.2 is not included in the relevant confidence interval, which is $[0.29, 0.41]$, and hence there is not a significant difference
- 4 The relevant confidence interval is $[0.33, 0.37]$, and hence the survival probability of 1st class passengers is at least 20 percentage points higher than the survival probability of other passengers
- 5 0.2 is not included in the relevant confidence interval, which is $[0.33, 0.37]$, and hence there is not a significant difference

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

All options are based on confidence intervals, and hence we will need the confidence interval for the difference, this is given by

$$\hat{p}_1 - \hat{p}_2 \pm z_{1-\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \quad (4)$$

I can also be calculated using `prop.test`;

```
s <- c(202, 117, 178, 212)
n <- c(325, 285, 706, 885)

p1 <- s[1]/n[1]
p2 <- sum(s[-1])/sum(n[-1])
n1 <- sum(n[1])
n2 <- sum(n[-1])
ph <- sum(s)/sum(n)

p1-p2+c(-1, 1)*qnorm(0.975)*sqrt(p1*(1-p1)/n1+p2*(1-p2)/n2)

## [1] 0.2948538 0.4077114

n_rest<-117+178+212
tot_rest<-285+706+885
prop.test(c(202,n_rest),c(325,tot_rest),correct=FALSE)

##
## 2-sample test for equality of proportions without continuity correction
##
## data:  c(202, n_rest) out of c(325, tot_rest)
```

```
## X-squared = 156.54, df = 1, p-value < 2.2e-16
## alternative hypothesis: two.sided
## 95 percent confidence interval:
##  0.2948538 0.4077114
## sample estimates:
##   prop 1   prop 2
## 0.6215385 0.2702559
```

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

Continue on page 16

Exercise V

A school class with 20 children are collecting trash on a beach, it is assumed that the mean value of the collected trash is 1kg/child with a standard deviation of 0.2 kg/child.

Question V.1 (11)

If the amount of trash collected by each child is assumed independent, what is the standard deviation (σ) of all the collected trash then?

- 1 $\sigma = 4.0$ kg
- 2 $\sigma = 0.8$ kg
- 3 $\sigma = 0.18$ kg
- 4 $\sigma = 2.0$ kg
- 5* $\sigma = 0.89$ kg

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

Using the independence assumption the variance of the total can be calculated as

$$V(Tot) = \sum_{i=1}^{20} \sigma^2 = 20\sigma^2 \quad (5)$$

and hence the standard deviation of the total is $\sqrt{20} \cdot 0.2$;

```
sqrt(20)*0.2
```

```
## [1] 0.8944272
```

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

Question V.2 (12)

After they returned, one of the children had collected 21 items, of which 6 were made of plastic. She is now asked to pick 5 items at random to be discussed. What is the probability that 3 of those are made of plastic?

- 1* 0.103
- 2 0.247

3 0.119

4 0.023

5 0.052

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

Since the total number of pieces and the number of plastic pieces is known this will be a hyper-geometric distribution, and the probability can be calculated by

```
dhyper(3,6,15,5)
```

```
## [1] 0.1031992
```

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

Continue on page 18

Question V.3 (13)

On average, 32% of the trash found is made of plastic. Another child collected 18 items, what is the chance that 3 of those items are made of plastic?

1* 0.082

2 0.100

3 0.124

4 0.876

5 0.958

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

```
dbinom(3, 18, 0.32)
```

```
## [1] 0.08218145
```

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

Continue on page 19

Exercise VI

The quality assurance department at a candy factory has taken a random sample of 26 chocolate bars of a certain brand. Each chocolate bar in the sample is weighted, and it is found that the average weight is 200.3 grams and the observed standard deviation is 0.75 grams.

Question VI.1 (14)

What is the 95% confidence interval for the standard deviation?

- 1 [0.346, 1.072]
- 2* [0.588, 1.035]
- 3 [0.611, 0.981]
- 4 [0.447, 1.053]
- 5 [0.462, 1.038]

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

The formula for the confidence interval is found in eq. (3-19) of the textbook:

$$\left[\sqrt{\frac{(n-1)s^2}{\chi_{1-\alpha/2}^2}}, \sqrt{\frac{(n-1)s^2}{\chi_{\alpha/2}^2}} \right] = \left[\sqrt{\frac{(26-1)0.75^2}{\chi_{1-0.05/2}^2}}, \sqrt{\frac{(26-1)0.75^2}{\chi_{0.05/2}^2}} \right] = [0.588, 1.035],$$

where the χ^2 -distribution has $\nu = n - 1$ degrees of freedom. The calculations in R are as follow:

```
sqrt((26-1)*0.75^2/qchisq(1-0.05/2,26-1))  
## [1] 0.588193  
sqrt((26-1)*0.75^2/qchisq(0.05/2,26-1))  
## [1] 1.035307
```

Thus, answer 2 is correct.

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

Question VI.2 (15)

The candy factory wants to test the null-hypothesis $\mathcal{H}_0 : \mu = 200$ grams (against a two-sided alternative) using a t -test. Which of the following statements is correct based on the hypothesis test? (Both the argument and the conclusion must be correct)

- 1 Using a significance level of 5%, the null-hypothesis is rejected since the test statistic is greater than $t_{0.975}(26)$
- 2 Using a significance level of 5%, the null-hypothesis is accepted since the test statistic is greater than $t_{0.975}(26)$
- 3 Using a significance level of 5%, the null-hypothesis is rejected since the test statistic is greater than $t_{0.975}(25)$
- 4 Using a significance level of 10%, the null-hypothesis is accepted since the test statistic is greater than $t_{0.95}(25)$
- 5* Using a significance level of 10%, the null-hypothesis is rejected since the test statistic is greater than $t_{0.95}(25)$

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

We apply method 3.36 and calculate the test statistics as

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{200.3 - 200}{0.75/\sqrt{26}} = 2.04.$$

This should be tested against a critical value of $t_{1-\alpha/2}(n-1)$ (the $(1-\alpha/2)$ quantile of a t -distribution with $n-1$ degrees of freedom). We find that $t_{0.95}(25) = 1.71$ and $t_{0.975}(25) = 2.06$, cf.

```
qt(0.95,df=25)
## [1] 1.708141
qt(0.975,df=25)
## [1] 2.059539
```

Since the test statistic is greater than $t_{0.95}(25)$, the null-hypothesis is rejected using a significance level of 10%. Thus, answer 5 is correct.

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

Continue on page 21

Question VI.3 (16)

To further investigate the mean weight of the chocolate bars, the candy factory is also planning another experiment. The quality assurance department wants to detect a difference in mean weight of 0.3 grams (against a two-sided alternative hypothesis) while using 0.75 grams as a guess of the standard deviation. Furthermore, the quality assurance department wants to keep both the Type I and the Type II error rates at (or below) 5%. What is the minimum number of chocolate bars to be included in the experiment in order to meet the criteria set by the department?

- 1 10 or 12 depending on whether you apply the normal approximation
- 2 68 or 70 depending on whether you apply the normal approximation
- 3* 82 or 84 depending on whether you apply the normal approximation
- 4 97 or 98 depending on whether you apply the normal approximation
- 5 162 or 164 depending on whether you apply the normal approximation

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

A sample size formula based on the normal approximation is found in method 3.65. The quantities in the formula are given as $\sigma = 0.75$, $\alpha = 0.05$, $\beta = 0.05$, and $\mu_0 - \mu_1 = 0.3$, which yields

$$n \geq \left(\sigma \frac{z_{1-\beta} + z_{1-\alpha/2}}{\mu_0 - \mu_1} \right)^2 = \left(0.75 \frac{1.645 + 1.960}{0.3} \right)^2 = 81.22,$$

which means the experiment should include at least 82 chocolate bars. The same calculations using R, which does not invoke the normal approximation, give a necessary sample size of 84 chocolate bars. The sample size determination in R is obtained as:

```
power.t.test(delta=0.3,sd=0.75,sig.level=0.05,power=0.95,type="one.sample")  
  
##  
##      One-sample t test power calculation  
##  
##          n = 83.16425  
##          delta = 0.3  
##          sd = 0.75  
##          sig.level = 0.05  
##          power = 0.95  
##          alternative = two.sided
```

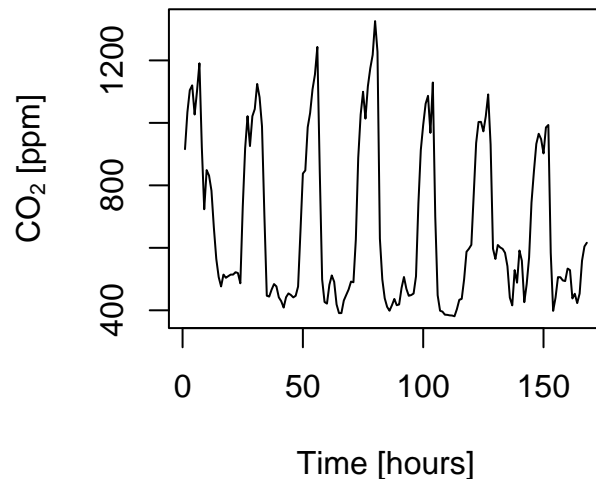
Thus, answer 3 is correct.

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

Continue on page 22

Exercise VII

CO₂ concentration is an important factor for well-being in the indoor environment, the figure below shows hourly CO₂ concentration [ppm] during a one week period in one room of a dwelling. The variance of the natural logarithm of the CO₂-concentration is 0.137.



As an initial analysis the CO₂ concentration is modeled as a function of time of day using the model

$$Y_i = \beta_0 + x_{1,i}\beta_1 + x_{2,i}\beta_2 + \epsilon_i,$$

where Y_i is the natural logarithm of CO₂ concentration at time i , $\epsilon_i \sim N(0, \sigma^2)$ and iid., and

$$x_{1,i} = \sin\left(2\pi\frac{h_i}{24}\right)$$
$$x_{2,i} = \cos\left(2\pi\frac{h_i}{24}\right),$$

where h_i is the hour of day for observation i .

The model is fitted and the result is reported below (some numbers are replaced by characters);

```
Call:
```

```
lm(formula = y ~ x1 + x2)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-0.59619 -0.09527  0.03135  0.12898  0.42424
```

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  6.43959    0.01468   t1      pv1
x1           0.40303    0.02076   t2      pv2
x2           0.20019    0.02076   t3      pv3
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: Sig on 165 degrees of freedom
Multiple R-squared:  R2, Adjusted R-squared:  0.7369
F-statistic: 234.9 on 2 and 165 DF, p-value: < 2.2e-16

```

Question VII.1 (17)

What is the total number of observations used for the estimation?

- 1 165
- 2 166
- 3 164
- 4 167
- 5* 168

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

The total number of obs. is the degrees of freedom plus 3, i.e. 168.

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

Question VII.2 (18)

What is the order of the p -values ($pv1$, $pv2$, and $pv3$) in the R-summary above?

- 1 $pv2 < pv3 < pv1$
- 2* $pv1 < pv2 < pv3$
- 3 $pv1 < pv3 < pv2$
- 4 $pv3 < pv1 < pv2$

5 $pv1 < pv2 = pv3$

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

The test statistics can be calculated as

```
c(6.4,0.4,0.2)/c(0.014,0.02,0.02)
## [1] 457.1429 20.0000 10.0000
```

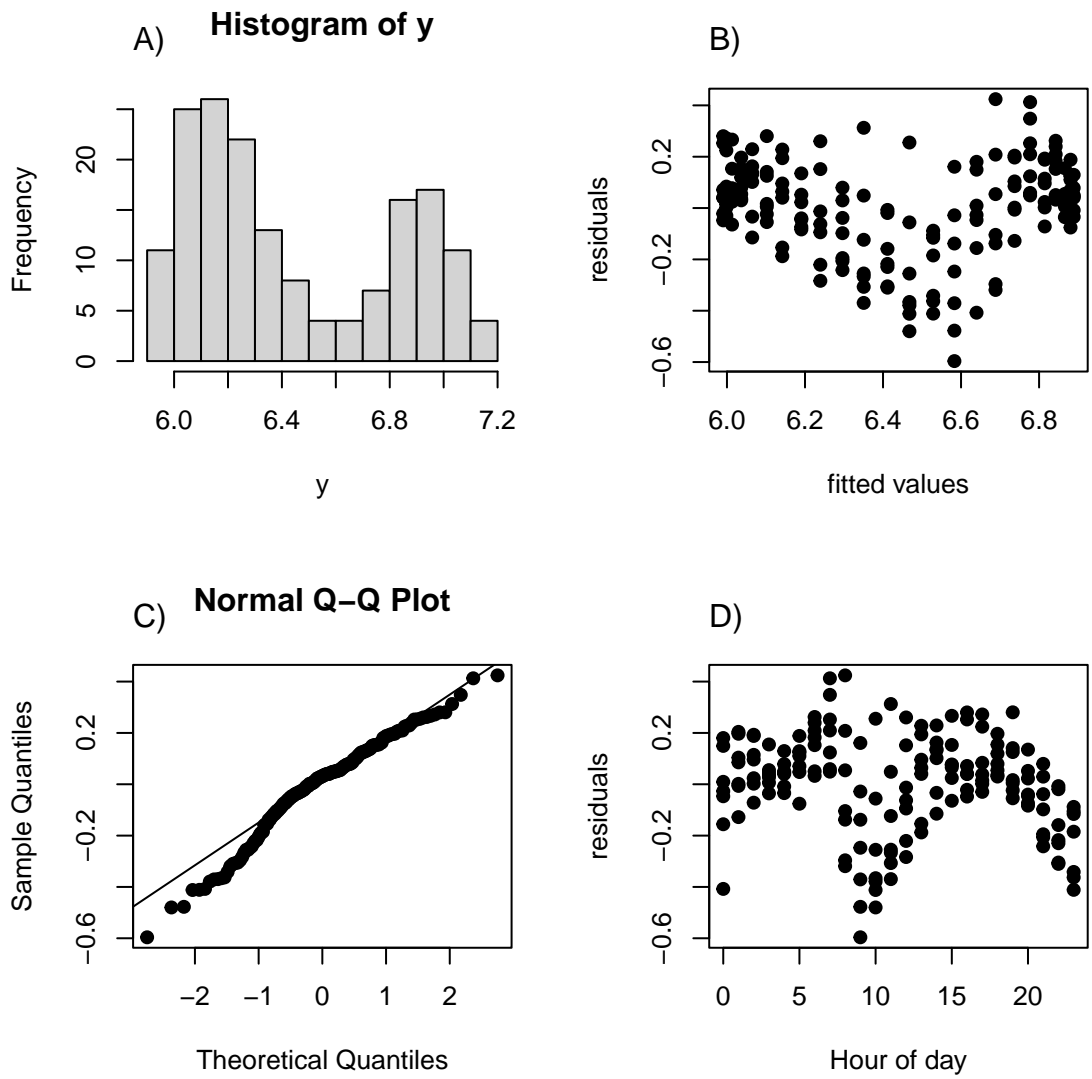
under the null hypothesis these all follow the same t-distribution, and hence the higher the test statistic the lower the p-value, hence $pv1 < pv2 < pv3$.

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

Continue on page 26

As part of the model validation the figure below is created. The plot show

- A) Histogram of y (log-CO₂ concentration)
- B) Residuals as a function of the fitted values using the model
- C) Normal quantile-quantile plot of the residuals from the model
- D) Residuals from the model as a function of hour of day



Continue on page 27

Question VII.3 (19)

Based on the plots in the figure, which of the following statements is correct (both the statement and figure reference should be correct)?

- 1 Based on figure A we should consider log-transforming the outcome
- 2 The residuals seems to be independent (figure C)
- 3 The normality assumption is clearly violated (figure A)
- 4 The residuals seems to be normally distributed (figure B)
- 5* There are still systematic effects related to time of day (figure D)

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

We will go through the statements.

- 1. Figure A is a histogram of the original data, not the residuals, hence it does not give information on the assumption of the model.
- 2. Figure C does not hold information on independence, hence the statement is false.
- 3. Figure A is a histogram of the original data, not the residuals, the plot can therefore not be used to validate the assumptions.
- 4. Figure B is not concerned with the normality assumption, hence it is false
- 5. There are clear systematic effect in the residuals a function of time of day (as seen in figure D), hence the statement is true.

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

Question VII.4 (20)

If x_1 and x_2 was removed from the model (so a constant mean model), what would the standard error related to the estimate of β_0 then be (hint: the variance of the outcomes is given above)?

- 1 0.0147
- 2 0.00990
- 3 0.00734

4 0.0208

5* 0.0106

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

This is calculated directly from the variance given in the initial statements of the exercise i.e.

```
0.137 / sqrt(168)
```

```
## [1] 0.01056978
```

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

Continue on page 29

Exercise VIII

14 days of whole-sale electricity prices and wind power forecasts have been collected in order to assess the effect of the wind production on electricity prices in some electricity market. Assume data follows a linear regression with normally distributed errors:

$$Y_i = \beta_0 + \beta_1 \cdot x_i + \varepsilon_i, \text{ where } \varepsilon_i \sim N(0, \sigma^2) \text{ and iid}$$

```
wind <- c(1063, 1450, 879, 1980, 406, 1542, 1212,
          1157, 1730, 1105, 775, 856, 802, 851)
elpris <- c(26.84, 24.87, 21.65, 13.26, 24.49, 21.90, 23.29,
            22.47, 19.26, 27.86, 27.96, 20.85, 21.83, 34.04)
```

Question VIII.1 (21)

What is the 99% confidence interval for the effect of wind power forecast on electricity price (β_1)?

- 1* [-0.0148, 0.0017]
- 2 [-0.0124, -0.0007]
- 3 [-0.0066, 0.0027]
- 5 [21.16, 40.87]
- 4 We have insufficient information to determine this.

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

```
L <- lm(elpris ~ wind)
confint(L, level = 0.99)

##              0.5 %          99.5 %
## (Intercept) 21.15572424 40.874663581
## wind        -0.01476638  0.001653692
```

or use Method 5.15

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

Question VIII.2 (22)

What is the 95% prediction interval for the electricity price when the wind power forecast is 1000 MWh?

- 1 [0.70, 12.41]
- 2 [11.41, 37.50]
- 3* [15.15, 33.76]
- 4 [21.95, 26.97]
- 5 [23.98, 38.05]

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

```
L <- lm(elpris ~ wind)
predict(L, level = 0.99, newdata = data.frame(wind = 1000), interval = "prediction")

##          fit          lwr          upr
## 1 24.45885 11.41344 37.50426
```

or use Method 5.18.

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

Continue on page 31

Exercise IX

A researcher is interested in investigating the effects of fertilizer and watering frequency on plant growth. A two-way ANOVA model for this data is:

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

where, Y_{ij} is the plant growth when applying the i 'th fertilizer and j 'th watering frequency ("Daily", "Twice a week" or "Weekly").

Question IX.1 (23)

Considering the statistical model above, which of the following statements regarding α_i is correct (note the statements are on the underlying model not on statistical tests)?

- 1 α_i denotes the effect size for watering frequency. $\alpha_i \neq 0$ implies that expected plant growth depends on watering frequency.
- 2 α_i denotes the effect size for watering frequency. This term should be omitted when plant growth depends on fertilizer type.
- 3* α_i denotes the effect size for fertilizer. $\alpha_i \neq 0$ implies that expected plant growth depends on fertilizer type.
- 4 α_i denotes the effect size for fertilizer. This term should be omitted when plant growth depends on fertilizer type.
- 5 α_i denotes the mean of the i -th fertilizer.

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

i refers to fertilizer, so α_i relates to the effect of fertilizer. More specifically, it refers to effect size (not the mean) for fertilizer i , and $\alpha_i \neq 0$ implies that fertilizer type has an effect on plant growth.

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

Question IX.2 (24)

A two-way ANOVA was carried out. The resulting ANOVA table is shown below. Please note that the p -values have been replaced by question marks.

```
## Analysis of Variance Table
```

```
## Response: Plant_Growth
```

##		Df	Sum Sq	Mean Sq	F value	Pr(>F)
##	Fertilizer	1	8.4017	8.4017	78.766	?
##	Watering_Frequency	2	4.0133	2.0067	18.812	?
##	Residuals	2	0.2133	0.1067		

Calculate the critical F-value for fertilizer and test the hypothesis of equal plant growth among fertilizers ($\alpha = 0.05$). Which of the following statements is the correct one?

- 1 $F_{crit} = 38.51$. We reject the null hypothesis of equal plant growth among fertilizers because $F_{obs} > F_{crit}$
- 2 $F_{crit} = 19$. We accept the null hypothesis of equal plant growth among fertilizers because $F_{obs} < F_{crit}$
- 3 $F_{crit} = 18.51$. We accept the null hypothesis of equal plant growth among fertilizers because $F_{obs} > F_{crit}$
- 4 $F_{crit} = 19$. We reject the null hypothesis of equal plant growth among fertilizers because $F_{obs} > F_{crit}$
- 5* $F_{crit} = 18.51$. We reject the null hypothesis of equal plant growth among fertilizers because $F_{obs} > F_{crit}$

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

We should use the 95% quantile from an F distribution with the right degrees of freedom. We read these from the ANOVA table (**Fertilizer** and **Residuals**). Hence we should use

```
qf(0.95,df1=1,df2=2)
## [1] 18.51282
```

$F_{obs} = 78.766$ is larger than F_{crit} , thus we reject the null hypothesis.

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

Question IX.3 (25)

Which of the following commands can be used to assess if the assumption of normality is fulfilled?

1*

```
lm1 <- lm(Plant_Growth~Fertilizer+Watering_Frequency, data)
qqnorm(lm1$residuals)
qqline(lm1$residuals)
```

2

```
lm1 <- lm(Plant_Growth~Fertilizer+Watering_Frequency, data)
lm1 <- anova(lm1)
qqnorm(lm1$residuals)
qqline(lm1$residuals)
```

3

```
qqnorm(data$Plant_Growth)
qqline(data$Plant_Growth)
```

4

```
qqnorm(data$Plant_Growth[data$Fertilizer=="Type 1"])
qqline(data$Plant_Growth[data$Fertilizer=="Type 1"])
```

5

```
qqnorm(rnorm(length(data)))
qqline(rnorm(length(data)))
```

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

We can always assess the normality assumption using the residuals of the model. Hence option 1 is correct. Option 2 converts `lm1` into an ANOVA table and would give an error when calling `qqnorm`.

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

Continue on page 34

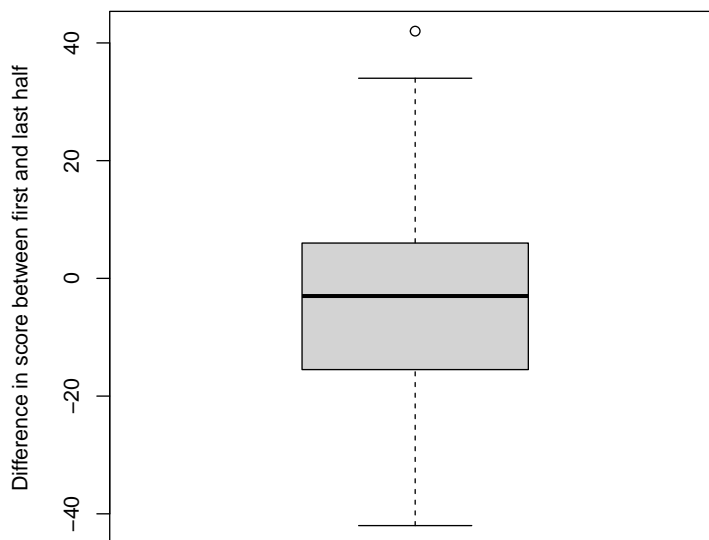
Exercise X

After a multiple choice exam in the introductory statistics course at DTU the teachers wanted to investigate the scores of different groups.

One question they would like to answer were: Were the students better at answering the first half of the exam (i.e. Question 1 to 15) than the last half (Question 16 to 30).

Let `xfirst` be a vector with students' scores in the first half of the exam and similarly `xlast` the students' scores for the second half of the exam. The observed differences in score between the last and the first half for all passed students is calculated and showed with a boxplot by:

```
x <- xlast - xfirst
boxplot(x, ylab="Difference in score between first and last half")
```



Continue on page 35

Question X.1 (26)

Which one of the following conclusions is wrong based on the information presented by the box-plot?

- 1 More than half of the students in the sample had a negative difference in scores.
- 2 More than 20% of the students in the sample had a positive difference in scores.
- 3 At least one student in the sample had a difference higher than 40 points in scores.
- 4* 60% of the students in the sample had a positive difference in scores.
- 5 No student in the sample had a difference in scores higher than 50 points.

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

We can not determine if 60% of the observations (i.e. the 60% quantile) in the sample is at 0, since we can only see that the median is below zero (i.e. the 50% quantile, the thick line in the middle of the box) and the 75% quantile (upper of the box) is above zero, so the 60% quantile can be anywhere in between.

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

Question X.2 (27)

The teachers want to test the null hypothesis

$$H_0 : \mu = 0,$$

where μ is the mean of the difference in scores between first and last part. They want to test without making any assumption of the distribution of the population where the sample was taken from.

The following code was run:

```
k <- 10000
simsamples <- replicate(k, sample(x, replace = TRUE))

quantile(apply(simsamples, 2, mean), c(0.05, 0.95))

##      5%      95%
## -6.54 -1.21

quantile(apply(simsamples, 2, mean), c(0.025, 0.975))
```

```
## 2.5% 97.5%
## -7.05 -0.77

quantile(apply(simsamples, 2, mean), c(0.005, 0.995))

## 0.5% 99.5%
## -8.09 0.23
```

Which one of the following answers is correct?

- 1* On a significance level $\alpha = 0.1$ a significant difference in scores between the first and last half is detected.
- 2 On a significance level $\alpha = 0.025$ a significant difference in scores between the first and last half is detected.
- 3 On a significance level $\alpha = 0.01$ a significant difference in scores between the first and last half is detected.
- 4 No conclusion can be made, the calculations don't meet the requirements, since in the calculations a normal distribution is assumed.
- 5 None of the answers above are correct.

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

The 90% confidence interval doesn't contain zero, so with $\alpha = 0.1$ a significant difference is detected. It would also be with $\alpha = 0.05$, but that is there is no answer with that.

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

Question X.3 (28)

The teachers wanted to investigate if the difference in score between the first and the second part of the exam differs according to the total score for a student. In order to investigate this the students were divided into two groups: one group that had a low total score and another group that had a high total score.

The score differences for low scoring students were stored in `xlow` and for high scoring students in `xhigh`.

The following code was executed:

```

k <- 10000
sim.xlow.samples <- replicate(k, sample(xlow, replace = TRUE))
sim.xhigh.samples <- replicate(k, sample(xhigh, replace = TRUE))

sim.xlow.means <- apply(sim.xlow.samples, 2, mean)
sim.xhigh.means <- apply(sim.xhigh.samples, 2, mean)
sim.dif.means <- apply(sim.xhigh.samples, 2, mean) -
  apply(sim.xlow.samples, 2, mean)

quantile(sim.xlow.means, c(0.025, 0.975))

## 2.5% 97.5%
## -9.23 -2.94

quantile(sim.xhigh.means, c(0.025, 0.975))

## 2.5% 97.5%
## -3.11 1.92

quantile(sim.dif.means, c(0.025, 0.975))

## 2.5% 97.5%
## 1.41 9.56

```

Which of the following conclusions is correct about the difference in mean of the two groups at significance level $\alpha = 0.05$ (both conclusion and argument must be correct)?

- 1 A significant difference between the two groups is not detected, since their one-sample confidence intervals overlap.
- 2 A significant difference between the two groups is detected, since their one-sample confidence intervals overlap.
- 3 A significant difference between the two groups is detected, since the one-sample confidence interval of one group includes zero, but that of the other groups' does not.
- 4* A significant difference between the two groups is detected, since the confidence interval for the difference in mean doesn't include zero.
- 5 None of the above conclusions are correct.

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

Only if the one-sample confidence intervals don't overlap we can use them for concluding if there is a difference, so we have to use the difference in mean to conclude.

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

Continue on page 39

Exercise XI

Question XI.1 (29)

Karl eats muesli in the morning, however he is picky and don't like raisins. Assuming that raisins appear at random (ie. can be described by a Poisson process), and that Karl's muesli portion contains 4 raisins on average, what is the probability that Karl's portion contains no raisins?

1 0.001

2* 0.018

3 0.183

4 0.250

5 0.368

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

```
dpois(0,4)
```

```
## [1] 0.01831564
```

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

Question XI.2 (30)

Karl's sister Karoline loves raisins and eats a muesli portion double the size of Karl's. What is the probability that her portion contains five or more raisins?

1 0.092

2 0.099

3 0.191

4 0.809

5* 0.900

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

```
1-ppois(4,8)
```

```
## [1] 0.9003676
```

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

The exam is finished. Enjoy the summer!