Technical University of Denmark

Written examination: 17.05.2025

Course name and number: Introduction to Statistics (02323)

Duration: 4 hours

Aids and facilities allowed: All aids - no internet access

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 12 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 not predetermined.

The answers should be given by filling in and submitting the digital 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	I.3	II.1	II.2	III.1	III.2	IV.1	IV.2	V.1
Question	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Answer										

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

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

The exam paper contains 22 pages.

Multiple choice questions: Note that in each question, one and <u>only</u> 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 the corresponding built-in functions in Python.

Exercise I

This exercise covers fundamental statistical concepts and basic Python functionality.

Question I.1 (1)

Which function in the scipy.stats-package is used to find theoretical quantiles?

- $1 \square$ The cdf-function
- $2 \square$ The pdf-function
- $3 \square$ The ppf-function
- $4 \square$ The rvs-function
- 5 \Box The std-function

Question I.2 (2)

Consider the sample (1, 2, 3, 9, 10). Which of the following statements is false?

- 1 \Box The sample lower quartile (Q1) is 2.
- $2 \square$ The sample median is 3.
- $3 \square$ The sample mean is 5.
- $4 \square$ The sample range is 9.
- $5 \square$ The sample variance is 16.

Question I.3 (3)

Which of the following quantities is not a random variable

- \Box $\,$ The population mean μ
- \Box The sample mean \bar{X}
- \Box The sample variance S^2
- \Box The test statistic T

Exercise II

Let Y follow an exponential distribution with rate parameter 2, and let U follow a continuous uniform distribution on the interval [3, 6]. The two random variables are independent.

Question II.1 (4)

What is the probability that Y exceeds 2?

 $1 \Box 2^{-2} = 0.250$ $2 \Box 2^{-1} = 0.500$ $3 \Box 1 - e^{-4} = 0.982$ $4 \Box e^{-4} = 0.018$ $5 \Box 2e^{-4} = 0.037$

Question II.2 (5)

What is the standard deviation of U?

 $1 \square \frac{3}{2}$ $2 \square \frac{3}{4}$ $3 \square \frac{\sqrt{3}}{2}$ $4 \square \frac{\sqrt{3}}{4}$ $5 \square \frac{9}{2}$

Exercise III

Three experiments compare three medications (A, B, and C) designed to regulate appetite. For each individual in each experiment, the decrease in calorie intake was measured. The results are stored in Dataset 1, Dataset 2, and Dataset 3 and are visualized in the boxplots below:



Question III.1 (6)

Which of the following statements is false?

- $1 \square$ The average effects of the medications seem the same in Datasets 1 and 3.
- $2 \square$ The SS(TR) Treatment Sum of Squares is largest in Dataset 2.
- $3 \square$ The SSE Sum of Squared Errors is largest in Dataset 3.
- 4 Within each dataset, the variance in each group (A, B, and C) is approximately equal.
- 5 \Box In all three datasets, medication C could be a placebo (i.e., have no average effect).

Question III.2 (7)

The sample size of each group is the same in all three datasets (i.e., all nine groups are of the same size). For each dataset, an analysis of variance (ANOVA) is performed, and the observed F-statistic and p-value are computed for the usual null hypothesis. Which of the following statements is correct?

- 1 \Box The observed *F*-statistic is the same for all three datasets.
- 2 \Box The test performed is a Welch *t*-test.
- $3 \square$ The *p*-value is the same for all three datasets.
- 4 \Box The assumption of equal variances across groups appears violated in Dataset 1.
- 5 \Box The *p*-value for Dataset 1 is larger than those for Datasets 2 and 3.

Exercise IV

A fitness coach claims that the average weight loss for a new diet program is 4 kg. To test this claim, a random sample of 12 participants who completed the program was taken, and their weight losses (in kg) were recorded: 3.5, 4.2, 3.8, 4.4, 4.1, 3.9, 4.0, 4.3, 3.6, 3.7, 4.5, 4.1.

The researchers then perform a t-test to determine whether the average weight loss differs significantly from 4 kg at a significance level of 5%.

Question IV.1 (8)

Let t_{obs} be the observed test statistic for the *t*-test. What is the rejection rule?

- 1 \square Reject H_0 if $t_{obs} \leq t_{1-\alpha}$, where $t_{1-\alpha}$ is the $(1-\alpha)$ quantile of a *t*-distribution with 11 degrees of freedom.
- 2 \square Reject H_0 if $t_{obs} \leq t_{1-\alpha}$, where $t_{1-\alpha}$ is the $(1-\alpha)$ quantile of a *t*-distribution with 12 degrees of freedom.
- 3 \square Reject H_0 if $t_{obs} \ge t_{1-\alpha}$, where $t_{1-\alpha}$ is the $(1-\alpha)$ quantile of a *t*-distribution with 11 degrees of freedom.
- 4 \square Reject H_0 if $t_{obs} \leq -t_{\alpha/2}$ or if $t_{obs} \geq t_{1-\alpha/2}$, where $t_{\alpha/2}$ and $t_{1-\alpha/2}$ are the $\alpha/2$ and $(1-\alpha/2)$ quantiles of a t-distribution with 11 degrees of freedom, respectively.
- 5 \square Reject H_0 if $t_{obs} \leq -t_{\alpha/2}$ or if $t_{obs} \geq t_{1-\alpha/2}$, where $t_{\alpha/2}$ and $t_{1-\alpha/2}$ are the $\alpha/2$ and $(1-\alpha/2)$ quantiles of a t-distribution with 10 degrees of freedom, respectively.

Question IV.2 (9)

What is the value of the observed test statistic (t_{obs}) ?

- $1 \square t_{\rm obs} = 2.83$
- $2 \square t_{obs} = 1.09$
- $3 \square t_{\rm obs} = 1.03$
- $4 \square t_{\rm obs} = 0.32$
- 5 \Box $t_{\rm obs} = 0.09$

Exercise V

A computer engineer wants to test the efficiency of an algorithm. She records the execution time (time) for jobs of varying complexity (complex).

To analyze the data, she runs the code below, where dat contains the recorded data.

```
fit1 = smf.ols(formula = 'time ~ complex', data = dat).fit()
print(fit1.summary(slim=True))
```

OLS Regression Results								
Dep. Variab Model: No. Observa Covariance	nonro	time OLS 50 nonrobust		ared: R-squared: atistic: (F-statistic	·):	0.971 0.971 1631. 1.03e-38		
	coef	std err		t	P> t	[0.025	0.975]	
Intercept complex	-0.1732 0.0006	0.007 1.47e-05	23 4(3.620).391	0.000 0.000	-0.188 0.001	-0.158 0.001	

Question V.1 (10)

Consider the test statistic for the null hypothesis that the intercept is zero. Which distribution is the test statistic compared to in order to obtain the associated *p*-value?

- $1 \square$ A t(50)-distribution
- 2 \Box A t(48)-distribution
- $3 \square$ An F(1, 50)-distribution
- 4 \square A $\mathcal{N}(0, 1^2)$ -distribution
- 5 \square A $\mathcal{N}(0, 0.007^2)$ -distribution

Question V.2 (11)

Based on the output above, which statement about the validity of the model (i.e., the model assumptions) is correct?

- 1 \Box The validity of the model assumptions cannot be assessed based on the output.
- 2 \Box The assumptions must be satisfied because the R^2 -value is close to 1.
- 3 \Box The model should be extended with additional terms because the *p*-values are small.
- 4 \Box The assumptions must be satisfied because the *p*-values are close to zero.
- 5 \square The model should be extended with additional terms because the R^2 -value is close to 1.

For the following questions, the numbers below may be useful:

```
print(np.mean(dat["complex"]))
484.58
print(np.var(dat["complex"], ddof=1))
12603.187346938777
```

Question V.3 (12)

Let $\hat{\sigma}$ denote the estimated residual standard deviation. What is the 95% confidence interval for the mean (execution) time of a job with complexity 300 according to the model?

- $1 \square 0.18 \pm 3.32 \hat{\sigma}$
- $2 \square 0.0068 \pm 2.08 \hat{\sigma}$
- $3 \square 0.18 \pm 0.82 \hat{\sigma}$
- $4 \square 0.0068 \pm 1.27 \hat{\sigma}$
- $5 \square 0.0068 \pm 0.55 \hat{\sigma}$

Question V.4 (13)

What is the estimate of the residual standard deviation $(\hat{\sigma})$?

The engineer figures that the execution time is proportional to the third power of complexity and formulates the following model:

$$Y_i = \beta_0 x_i^{\beta_1} \varepsilon_i; \quad \varepsilon_i \sim LN(0, \sigma^2),$$

where Y_i and x_i are the time and complexity of job *i*, respectively. She hence tests the hypothesis $\mathcal{H}_0: \beta_1 = 3$ against a two-sided alternative by running the below code chunk:

```
ltime = np.log(dat['time'])
lcomplex = np.log(dat['complex'])
dat2 = pd.DataFrame({'ltime' : ltime, 'lcomplex' : lcomplex})
fit2 = smf.ols(formula = 'ltime ~ lcomplex',data=dat2).fit()
print(fit2.summary(slim=True))
```

		OLS R	legres	sion Re	esults			
Dep. Variab	1	time	R-squ	lared:		0.999		
Model:		OLS	Adj.	R-squared:		0.999		
No. Observa		50	F-sta	atistic:		8.032e+04		
Covariance	nonro	bust	Prob	(F-statist	ic):	4.87e-79		
	coef	std err		t	P> t	[0.025	0.975]	
Intercept lcomplex	-18.0015 2.5459	0.055 0.009	-32! 283	5.241 3.406	0.000 0.000	-18.113 2.528	-17.890 2.564	
			=====				============	

Question V.5 (14)

Based on the output above, what is the test statistic (t_{obs}) for the usual test of \mathcal{H}_0 , and what is the conclusion of the test at a significance level of 5%?

- \Box The hypothesis is rejected as $t_{\rm obs} = -5.645$.
- \Box The hypothesis is accepted as $t_{\rm obs} = 13.84$.
- \Box The hypothesis is rejected as $t_{\rm obs} = -50.46$.
- \Box The hypothesis is rejected as $t_{\rm obs} = 283.4$.
- \Box The hypothesis is accepted as $t_{\rm obs} = -5.645$.

Exercise VI

The two questions in this exercise must be solved using simulation.

The number of customers shopping in a supermarket on a randomly selected day follows a Poisson distribution with a mean of 5000. The amount a randomly selected customer spends is exponentially distributed, with an average of 200 DKK. It is assumed that customers act independently of each other and independently of the number of customers.

The following code snippet can be used to simulate k daily revenues, but you need to complete the missing parts of the code by specifying the parameters of the distributions:

```
# Set seed
np.random.seed(2025)
# Number of samples
k = 100000
# Number of customers (array of customer counts for the k days)
N = stats.poisson.rvs(mu=____,size=k)
# The revenues (array with revenues for the k days)
Y = np.array([np.sum(stats.expon.rvs(scale=____, size=n)) for n in N])
```

Let Y denote the supermarket's revenue (or turnover) on a randomly selected day.

Question VI.1 (15)

What is the expectation and the standard deviation of Y?

- 1 \square $\mathbb{E}[Y] = 1,500,000$ DKK and SD(Y) = 1,000,000 DKK
- 2 \square $\mathbb{E}[Y] = 1,500,000$ DKK and SD(Y) = 20,000 DKK
- $3 \square \mathbb{E}[Y] = 1,000,000 \text{ DKK} \text{ and } \text{SD}(Y) = 1,000,000 \text{ DKK}$
- 4 \square $\mathbb{E}[Y] = 1,000,000$ DKK and SD(Y) = 20,000 DKK
- 5 \square $\mathbb{E}[Y] = 500,000$ DKK and SD(Y) = 1,000,000 DKK

Question VI.2 (16)

What is the probability that the supermarket makes more than 1,500,000 DKK in revenue on a randomly selected day?

- $1 \square$ Almost 0%
- $2 \square$ Approximately 11%
- $3 \square$ Approximately 22%
- $4 \square$ Approximately 33%
- $5 \square$ Approximately 50%

Exercise VII

A shipping company has hired a sales representative to recruit new clients. The sales representative visits n (independent) clients each month, and historical data suggests that the probability of recruiting a client after a visit is p. Let X denote the number of recruited clients in a randomly selected month.

Question VII.1 (17)

Which one of the following models is the most appropriate?

 $1 \square X \sim U(0, n)$ $2 \square X \sim H(n, np, n)$ $3 \square X \sim \mathcal{N}(np, p^2)$ $4 \square X \sim \text{Pois}(np)$ $5 \square X \sim \text{Bin}(n, p)$

The sales representative earns a salary of 5000 USD per month and receives an additional bonus of 100 USD for each client recruited during the month. The sales representative's total compensation in a randomly selected month (Y) is therefore given by

Y = 100X + 5000,

where X denotes the number of clients recruited during that month. You may assume that the standard deviation of X is 2.

Question VII.2 (18)

What is the variance of Y?

- $1 \square 5200$
- $2\square$ 20000
- $3\square$ 25000
- 4 🗌 40000
- 5 🗌 45000

Question VII.3 (19)

Which one of the following statements about the correlation between X and Y is true?

- \square The correlation between X and Y is one: $\rho(X, Y) = 1$.
- \square The correlation between X and Y is between one and zero: $0 < \rho(X, Y) < 1$.
- \square The correlation between X and Y is zero: $\rho(X, Y) = 0$.
- \square The correlation between X and Y is negative: $-1 \le \rho(X, Y) < 0$.
- \Box The correlation cannot be determined without additional information.

Exercise VIII

A municipality tested a service that allows citizens to report damages to traffic infrastructure (e.g., roads and traffic lights) via an app. The reported issues were categorized into three groups: needs immediate repair, can be postponed, or can be ignored.

During a one-year trial, n = 165 independent issues were reported. The counts for each issue type were compiled by season and entered into the below table:

```
# Reading the data into Python
data = np.array([[44, 10, 32, 18], [20, 8, 4, 2], [11, 2, 10, 4]])
# Converting to a Pandas dataframe
data = pd.DataFrame(data, index=['Immediate', 'Postpone', 'Ignore'],
columns=['Spring', 'Summer', 'Fall', 'Winter'])
print(data)
           Spring
                   Summer
                            Fall
                                  Winter
Immediate
               44
                        10
                              32
                                      18
Postpone
               20
                         8
                               4
                                       2
                         2
                              10
                                       4
Ignore
               11
```

The municipality plans to test if the distribution of fault types depends on the season. Specifically, it tests the null hypothesis

$$H_0: p_{i1} = p_{i2} = p_{i3} = p_{i4} = p_i$$
 for $i = 1, 2, 3,$

where p_{ij} represents the proportion of observations in column j that fall into row i i.e., the proportion of fault type i out of all faults in season j, and p_i represents the proportion of all observations that are in row i.

Question VIII.1 (20)

Under the null hypothesis, what is the estimated proportion of faults requiring immediate repair?

- $1 \square 0.6303$
- $2 \square 0.4231$
- $3\square$ 0.2061
- $4 \square 0.7939$
- 5 \Box None of the above answers are correct.

Question VIII.2 (21)

At a significance level of $\alpha = 0.05$, what is the correct conclusion for the usual test of the null hypothesis? (Both argument and conclusion must be correct.)

- 1 \Box The null hypothesis is rejected, indicating a significant difference in distribution between the seasons, as the *p*-value is 0.04.
- 2 \Box The null hypothesis is rejected, indicating a significant difference in distribution between the seasons, as the *p*-value is 0.08.
- 3 \Box The null hypothesis is accepted, indicating no significant difference in distribution between the seasons, as the *p*-value is 0.04.
- 4 \Box The null hypothesis is accepted, indicating no significant difference in distribution between the seasons, as the *p*-value is 0.08.
- 5 \Box No conclusion can be drawn, since some expected cell counts under the null hypothesis are too low i.e., $e_{ij} < 5$ for some cells.

Another purpose was to test an AI system for automatically handling the reported faults, comparing its performance to manual evaluation, where a human categorized the faults. Based on expert knowledge, the number of discrepancies between the two methods is expected to be 21 (out of the 165 observed faults).

Question VIII.3 (22)

Let $X \sim Bin(n, p)$ be the count of how many faults were evaluated differently (You need to estimate p). What is the estimate of the standard deviation of X/n?

 $1 \Box \quad \hat{\sigma} = 0.3343$ $2 \Box \quad \hat{\sigma} = 0.1273$ $3 \Box \quad \hat{\sigma} = 0.1118$ $4 \Box \quad \hat{\sigma} = 0.0259$ $5 \Box \quad \hat{\sigma} = 0.0007$

Exercise IX

A researcher designs an experiment to estimate the mean time of a chemical reaction. The experiment should be able to detect a difference in mean of 0.5 seconds assuming a standard deviation of 1.2 seconds for the reaction time (previous studies suggest this is the best guess of σ). The researcher would like to use a 95% confidence level and a power of 80% for this experiment.

Question IX.1 (23)

What is the minimum sample size required to achieve the desired specifications? (Note that the facit is calculated with Python. Thus, facit can differ slightly from results obtained with formulas.)

- $1 \square 38$ $2 \square 47.17$
- $3\square 48$
- _
- 4 🗌 63
- $5\square$ 79

Exercise X

In an experiment, n = 85 individuals are tested for a binary response (positive or negative). Let the binomially distributed random variable

$$X \sim \operatorname{Bin}(n, p)$$

denote the number of positive outcomes.

The following null hypothesis is considered

$$H_0: p = 0.2.$$

Question X.1 (24)

Under the null hypothesis (i.e., assuming it is true), what is the probability of observing 25 or more positive outcomes?

- $1 \square 0.025$
- $2\square$ 0.030
- 3 🗌 0.061
- 4 🗌 0.132
- $5 \square 0.987$

Question X.2 (25)

If 25 positive outcomes are actually observed, what is the observed test statistic in the usual test of the null hypothesis?

- 1 🗌 8.000
- $2\square$ 2.169
- $3 \square 1.904$
- $4 \square 0.057$
- $5 \square 0.030$

Exercise XI

A study is conducted to compare the average scores of two groups of students in a math test. Group A consists of 12 students with a mean score of 78 and a standard deviation of 10, while Group B consists of 15 students with a mean score of 82 and a standard deviation of 8. Researchers want to perform a two-sample t-test to determine if there is a significant difference between the means of the two groups at a 5% significance level i.e., they want to test the null hypothesis

$$H_0: \mu_A - \mu_B = 0$$

against a two-sided alternative hypothesis.

Question XI.1 (26)

What is the observed test statistic in the appropriate test?

- 1 -3.422
- 2 🗌 -1.155
- 3 🗌 -1.127
- 4 🗌 -0.765
- 5 🗌 -0.317

Question XI.2 (27)

The test statistic follows a *t*-distribution with how many degrees of freedom?

- 1 🗌 20.853
- $2\square$ 22.382
- $3\square 24.998$
- $4\square 25$
- $5\square 26$

Exercise XII

Consider the random variable Y and the sample $\mathbf{y} = (6, 5, 4, 4, 1, 10, 7, 6, 5, 3)$.

Question XII.1 (28)

Using non-parametric bootstrapping with 100,000 replications (bootstrap samples), what is the 90% confidence interval for the mean of Y? (Note: The facit uses np.random.seed(2025) and your answer may vary according to the seed.)

- $1 \square [3.7, 6.6]$
- $2 \square [3.9, 6.3]$
- $3 \square [4.2, 6.0]$
- $4 \square [5.0, 5.2]$
- $5 \square [5.1, 5.1]$

Question XII.2 (29)

Assume that Y follows a Poisson distribution. What is the 90% confidence interval for the mean of Y using parametric bootstrapping with 100,000 replications (bootstrap samples)? (Note: The facit uses np.random.seed(2025) and your answer may vary according to the seed.)

- $1 \square [3.7, 6.6]$
- $2 \square [4.0, 6.3]$
- $3 \square [4.2, 6.0]$
- $4 \square [5.0, 5.2]$
- $5 \square [5.1, 5.1]$

The following chunk of code has been run:

```
# Set seed
np.random.seed(2025)
# Data
y = np.array([6, 5, 4, 4, 1, 10, 7, 6, 5, 3])
# Estimated mean
mu = y.mean()
# Simulate bootstrap samples
simsamples = stats.poisson.rvs(mu,size=(1000,len(y)))
simsamples = pd.DataFrame(simsamples)
# Compute medians
simmedians = simsamples.median(axis=1)
# Find confidence interval
print(np.quantile(simmedians,[0.05,0.95],
method='averaged_inverted_cdf'))
```

[3.5 6.5]

Question XII.3 (30)

Consider the null hypothesis \mathcal{H}_0 : median(Y) = 5. What is the correct conclusion about \mathcal{H}_0 based on the code chunk?

- 1 \square Based on the parametric bootstrapping procedure, \mathcal{H}_0 cannot be rejected on a significance level of 90%.
- 2 \square Based on the parametric bootstrapping procedure, \mathcal{H}_0 is rejected on a significance level of 10%.
- 3 \square Based on the parametric bootstrapping procedure, \mathcal{H}_0 cannot be rejected on a significance level of 10%.
- 4 \square Based on the non-parametric bootstrapping procedure, \mathcal{H}_0 is rejected on a significance level of 5%.
- 5 \square Based on the non-parametric bootstrapping procedure, \mathcal{H}_0 cannot be rejected on a significance level of 10%.

The exam is finished. Enjoy the vacation!