**ISEG – Lisbon School of Economics and Management**
**STATISTICS 2**
**Second Semester 2020/2021**
**Normal Exam**
**Monday 31 May 2021**

Duration: **2 hours**

**Name:**
**Student Number:**

| Question: | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Total |
|---|---|---|---|---|---|---|---|---|
| Points: | 20 | 10 | 30 | 10 | 50 | 40 | 40 | 200 |

**Justify** all your answers. A correct answer in a multiple choice question is worth 10 points; an incorrect one is worth $-2.5$ points. **You are allowed to use** the <u>formula sheet</u>, and the <u>statistical tables</u> that are provided on the webpage of the course, without ANY annotation. You can also use a <u>calculator</u> and <u>A4 empty scrap paper</u>. Consider a **5% significance level**, unless otherwise is stated.

1. The daily number of COVID-19 infections on the north portuguese municipalities are assumed to follow an exponential distribution with mean 35, whereas the daily number of COVID-19 infections on the south portuguese municipalities are assumed to follow an exponential distribution with mean 65. For now, assume that the COVID-19 infections, per municipality (both north and south), are independent.

(10)  (a) In a sample of 10 north portuguese municipalities, what is the probability that the daily maximum number of COVID-19 infections exceeds 100?

$$X_N \sim exp\left(\lambda = \tfrac{1}{35}\right) \quad E[X] = \bar{\lambda} = 35. \quad F_{X_N}(x) = 1 - e^{-\lambda x}$$
$$n = 10$$

$$P(X_{(n)} > 100) = 1 - P(X_{(n)} < 100)$$
$$= 1 - P(X_1 < 100, \ldots, X_n < 100)$$
$$= 1 - P(X_1 < 100) * \ldots * P(X_n < 100)$$
$$= 1 - \left\{P(X_N \leq 100)\right\}^n$$
$$= 1 - \left\{F_{X_N}(100)\right\}^n$$
$$= 1 - \left(1 - e^{-\tfrac{1}{35} * 100}\right)^{10}$$
$$\approx 0.4465$$

(10)     (b) Consider a sample of 100 north portuguese municipalities and 100 south portuguese municipalities. Approximate the probability that the average number of daily infections in the north municipalities is at least 31 lower than the average number of daily infections in the south municipalities?

$$\bar{X}_N \overset{appox}{\underset{CLT}{\sim}} N\left(\mu=35, \frac{35^2}{100}\right) \quad \bar{X}_S \overset{appox}{\underset{CLT}{\sim}} N\left(\mu_S=65, \frac{65^2}{100}\right)$$

$$P(\bar{X}_N < \bar{X}_S - 31) = P(\underbrace{\bar{X}_N - \bar{X}_S}_{N\left(-30, \frac{35^2 + 65^2}{100}\right)} < -31)$$

$$= P\left(\underbrace{\frac{\bar{X}_N - \bar{X}_S + 30}{\sqrt{(35^2 + 65^2)/100}}}_{N(0,1)} < \underbrace{\frac{-31 + 30}{\sqrt{(35^2 + 65^2)/100}}}_{-0.14}\right) \quad \longrightarrow 7.3824$$

$$= 1 - \Phi(0.14) \approx 0.444$$

$$\underbrace{}_{0.557 \ (tables)}$$

(10)   **2**. Which of the following statements is correct (only one statement is correct)?

    ○ A statistic is a function of the unknown population parameter(s).

    ○ A statistic can assume infinitely many values, assuming that we can draw infinitely many samples, with the same sample size, from the population.

    ○ The repeated sampling distribution of any statistic can be approximated by the Standard Normal distribution.

    ⊗ None of the above.

**Turn over**

**3**. Let $X$ be a random variable with probability distribution function given by:

$$f(x \mid \theta) = \theta (1-\theta)^{x-1}, \quad x = 1, 2, \ldots, \quad 0 < \theta < 1, \quad \mathrm{E}(X) = \frac{1}{\theta} \text{ and } \mathrm{Var}(X) = \frac{1-\theta}{\theta^2}$$

(15)      (a) Obtain the maximum likelihood estimator for $\theta$.

1) $f_X(x) = \theta (1-\theta)^{x-1}$

2) $\ell_i(\theta) = \log(\theta) + (x_i - 1) \log(1-\theta)$

   $\ell_n(\theta) = n \log(\theta) + \left( \sum_{i=1}^{n} x_i - n \right) \log(1-\theta)$

3) $\dfrac{d\ell_n(\theta)}{d\theta} = \dfrac{n}{\theta} - \dfrac{(\sum x_i - n)}{1-\theta} = 0$

   $\Leftrightarrow \dfrac{n(1-\theta) - \theta \sum x_i + n\theta}{\theta(1-\theta)} = 0 \quad \to \quad n - \cancel{n\theta} - \theta\sum x_i + \cancel{n\theta}$

   $\theta(1-\theta) > 0$

   $\Leftrightarrow n = \theta \sum_{i=1}^{n} x_i \quad \Leftrightarrow \quad \hat{\theta}_{ml} = \dfrac{1}{\bar{x}}$

(15)      (b) Obtain the maximum likelihood estimator for $\mathrm{E}(X)$. Make sure to explain your answer.

$\mathrm{E}[X] = \dfrac{1}{\theta}$

$\widehat{\mathrm{E}[X]}_{ml} = \dfrac{1}{\hat{\theta}_{ml}}$   by **invariance** of ml to (non linear) transformations

$= \dfrac{1}{(1/\bar{x})} = \bar{X}.$

(10) **4.** For any random variable $X$, with a given probability distribution function, $f(x \mid \boldsymbol{\theta})$, where $\boldsymbol{\theta}$ can be a vector of unknown parameters, which of the following statements is correct (only one statement is correct)?

     ◯ The Method of Moments estimator is always equal to the Maximum Likelihood estimator, for some unknown parameter of $\boldsymbol{\theta}$.

     ◯ There cannot be any estimator that takes a value below the Cramer-Rao lower bound.

     ◯ If $X \sim N(\mu, \sigma^2)$, then the Maximum Likelihood estimator for both $\mu$ and $\sigma^2$ is unbiased.

     ✗ None of the above.

**5.** Current evidence on COVID-19 shows that the time (in days) until an individual develops disease symptoms, after getting infected, follows a Normal distribution. More precisely, for male individuals, it follows a Normal distribution with unknown mean $\mu_M$ and unknown variance $\sigma_M^2$, whereas for female individuals, it follows a Normal distribution with unknown mean $\mu_F$ and unknown variance $\sigma_F^2$. The following sample results were obtained:

Male individuals:    $n_M = 13$,   $\bar{x}_M = 6.257$,   $s_M^2 = 0.709$   $\to S_M^2$

Female individuals:   $n_F = 12$,   $\bar{x}_F = 2.945$,   $s_F^2 = 0.518$   $\to S_F'^2$

*(handwritten annotations: 11, 13; 10, 12)*

(20)     (a) Based on a 99% confidence interval, comment on the statement: "in fact, the mean values for the distributions of the male and female individuals can be assumed to be the same".

step 1) Find pivotal quantity:

3.312   $\dfrac{(\bar{X}_M - \bar{X}_F) - (\mu_M - \mu_F)}{\sqrt{\dfrac{S_M'^2}{11} + \dfrac{S_F'^2}{10}}}$   approx $\sim N(0,1)$

0.341   0.929

step 2) Find quantiles:   $q_{0.995} = 2.576$   (tabelas)

step 3) Rewrite   $P(-2.576 < PQ < 2.576) = 0.99$:

$P\left(3.312 - 2.576 \times 0.341 < \mu_M - \mu_F < 3.312 + 2.576 \times 0.341\right) = 0.99$

2.434        4.190

statement incorrect.

(20)     (b) Test the hypothesis that the variances are equal, using $\alpha = 5\%$.

step 1) $\dfrac{S_m'^2/\sigma_m^2}{S_F'^2/\sigma_F^2} \;=\; \dfrac{S_m'^2}{S_F'^2} \sim F_{10,9}$

$\underset{H_0:\,\sigma_m^2=\sigma_F^2}{\Big\downarrow}$

step 2) Find quantiles:

$q_{0.975}^{F_{10,9}} = 3.96 \; ; \; q_{0.025}^{F_{10,9}} = \dfrac{1}{q_{0.975}^{F_{9,10}}} = \dfrac{1}{3.78} = 0.26$

step 3) $F_{obs} = \dfrac{0.709}{0.518} = 1.37 \;\in\; (0.26,\, 3.96)$

do $\underline{\underline{not}}$ reject $H_0$

(10)   (c) Which of the following statements is correct (only one statement is correct)?

- ◯ If, for a given test, the significance level is equal to $\alpha\%$, then the power of that test is equal to $(1-\alpha)\,\%$.
- ◯ The Type I error and the Type II error are equal for two-sided tests.
- ◯ A pivotal quantity is a test statistic, if its repeated sampling distribution is completely known when the alternative hypothesis is true.
- ☒ None of the above.

6. A survey to 1000 portuguese individuals, regarding the outcome of COVID-19 infection, obtained the following results:

| | | Age | | |
|---|---|---|---|---|
| | | $\leq 40$ | $(40, 65)$ | $\geq 65$ |
| Outcome | Recovered | 204 | 539 | 57 |
| | Died | 10 | 101 | 89 |

(20)   (a) Based on an adequate test, at the 1% level, say if you agree with the following statement: "for those who died of COVID-19 infection, the distribution of age is discrete uniform (each age group has probability $1/3$)."

$H_0 : P_1 = P_2 = P_3 = \frac{1}{3} \qquad O_1 = 10,\; O_2 = 101,\; O_3 = 89$

$E_1 = \dfrac{200}{3} = E_2 = E_3$

$Q_{obs} = \dfrac{(10 - 66.67)^2}{66.67} + \dfrac{(101 - 66.67)^2}{66.67} + \dfrac{(89 - 66.67)^2}{66.67} = 73.32 > 9.210$

$q_{0.99}^{\chi_2^2} = 9.210 \;\Rightarrow\; \underline{reject}\; H_0.$

(20) (b) Is the outcome of COVID-19 infection independent of age? Use $\alpha = 5\%$.

$$O_{11} = 204 \quad O_{12} = 539 \quad O_{13} = 57 \quad 800$$
$$O_{21} = 10 \quad O_{22} = 101 \quad O_{23} = 89 \quad 200$$
$$214 \quad 640 \quad 146 \quad 1000$$

$$E_{11} = \frac{800 * 214}{1000} = 171.2 \quad E_{12} = 512 \quad E_{13} = 116.8$$
$$E_{21} = 42.8 \quad E_{22} = 128 \quad E_{23} = 29.2$$

$$Q \underset{H_0}{\sim} \chi^2_2 \quad \text{as} \quad (r-1)(c-1) = 2 : \quad \chi^2_{2,0.95} = 5.991$$

$$Q_{obs} = \frac{(204 - 171.2)^2}{171.2} + \ldots + \frac{(89 - 29.2)^2}{29.2} = 191.6 > 5.991$$

we **reject** $H_0$

7. Consider the following models to explain the average monthly wage (in thousand euros) in terms of the number of years of education (estimated using OLS, standard errors in parentheses):

$$\ln \widehat{(wage_i)} = \underset{(0.4439)}{-1.9970} + \underset{(0.1756)}{1.2562} \ln (educ_i), \quad n = 428, \quad R^2 = 0.1073$$

$$\ln \widehat{(wage_i)} = \underset{(0.1852)}{-0.1852} + \underset{(0.0144)}{0.1086} educ_i, \quad n = 428, \quad R^2 = 0.1179$$

(15) (a) Interpret the estimated coefficients of the regressors $\ln (educ)$ and $educ$.

$\ln(educ)$: if I increase $educ$ with $\underline{1\%}$, wage (on avg) increases with $\underline{1.2562\%}$

$(educ)$: if I increase $educ$ with $\underline{1\ year}$, wage (on avg) increases with $\underline{10.86\%}$

(15) (b) Test the significance of $\ln (educ)$ (i.e. test $H_0 : \beta_{\ln(educ)} = 0$).

$$\begin{cases} H_0 : \beta_{\ln(educ)} = 0 \\ H_1 : \beta_{\ln(educ)} \neq 0 \end{cases}$$

we did not assume that $\varepsilon$ has a Normal distribution

$$T = \frac{\hat{\beta}_{\ln(educ)} - 0}{\sqrt{\widehat{Var}(\hat{\beta}_{\ln(educ)})}} \overset{approx}{\sim} N(0,1) : q^{N(0,1)}_{0.975} = 1.960$$

$$T_{obs} = \frac{1.2562 - 0}{0.1756} = 7.15 > 1.96$$

we  <u>reject</u>  $H_0$

(10)     (c) Based on your result in question b, can you say something about the confidence interval for $\beta_{\ln(educ)}$? Explain.

The 95% confidence interval for $\beta_{\ln(educ)}$ does <u>NOT</u> contain the value 0, because the test in Q7b rejected $H_0 : \beta_{\ln(educ)} = 0$.

**You can use the space below and the other side of this page to add to your answers. Clearly state which question you are referring to.**