

BIOSTATISTICS (MATH-449)

ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

Mock exam

Date: 27th of May, 2025

Time: 10:15–12:00

Name: _____

SCIPER: _____

INSTRUCTIONS TO CANDIDATES

- This is a mock exam. The final one will contribute 80% to your final grade. To obtain the maximum number of points you should be clear about your reasoning and present your arguments explicitly. ~~You have **3 hours** to complete the exam.~~ You have **2 hours** to complete the mock exam. (We are not expecting to finish all questions within the 2 hours, but we wanted to give more opportunities to practice material.)
- All that can be used for this exam is a pen. No notes, books, summaries, formula collections or calculators are allowed. All questions should be answered.
- The finest enumerated item in each question will be marked on a scale of 0 – 2 points, indicating an incorrect, partially correct and completely correct answer respectively (half-points are not given). **The mock exam has 4 questions with a total of 32 points.**
- **Write the answer to every question in *the other* booklet (Final exam - answers).** (This is a mock exam, you will not be provided with an answer booklet.) The questions are provided in this booklet. Scrap paper during the final exam will be provided for rough work, but only answers written in the other booklet will be marked.
- At the end of the exam, you will have to return everything: the booklet with the questions, the booklet with your answers, and the scrap paper.

Mark question 1 (TOT: 8 points):

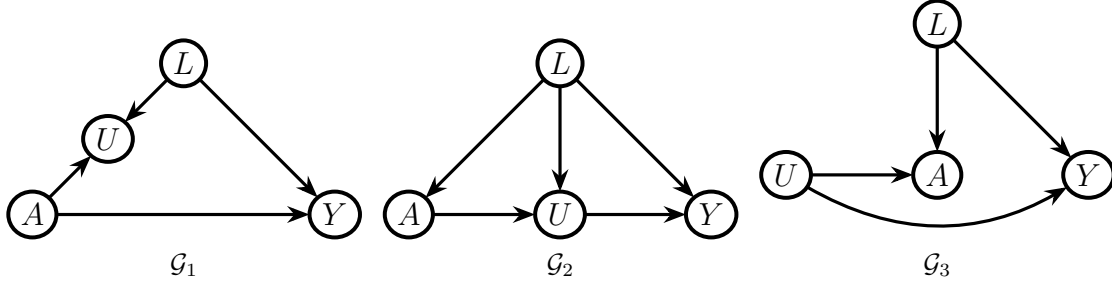
Mark question 2 (TOT: 6 points):

Mark question 3 (TOT: 6 points):

Mark question 4 (TOT: 12 points):

Question 1.

Let L, U, A and Y denote a measured baseline covariate, an unmeasured variable, treatment and outcome, respectively. Investigators 1, 2 and 3 propose causal models represented by the causal graphs $\mathcal{G}_1, \mathcal{G}_2$ and \mathcal{G}_3 respectively:



- (a) Does any of the graphs \mathcal{G}_1 , \mathcal{G}_2 and \mathcal{G}_3 imply independencies in the observed law $p(l, a, y)$? Can the observed law $p(l, a, y)$ be used to falsify \mathcal{G}_1 , \mathcal{G}_2 and \mathcal{G}_3 ? If yes, state the independencies that can be falsified.
- (b) Consider the causal effect $E[Y^{a=1} - Y^{a=0}]$, assuming that consistency and positivity hold, and that
 - (i) \mathcal{G}_1 is correct.
 - (ii) \mathcal{G}_2 is correct.
 - (iii) \mathcal{G}_3 is correct.

For parts (b) (i)–(iii), express $E[Y^{a=1} - Y^{a=0}]$ as a function of the observed law $p(l, a, y)$ if this is possible, otherwise explain why it is not possible. Your expression should be "minimal", in the sense that none of the terms should contain unnecessary variables.

Question 2. (a) Let $\alpha(t)$ be a hazard function and $S(t)$ the corresponding survival function for a (continuously distributed) survival time T . Show the identity

$$S(t) = \exp \left(- \int_0^t \alpha(s) ds \right).$$

- (b) Let $Z_j(t)$ denote the number of individuals at risk and $N_j(t)$ the counting processes of events in group j . State the formula for the Kaplan-Meier estimators of the survival functions for the groups $j = 1, 2$.
- (c) It is sometimes considered insufficient to just present Kaplan-Meier (survival) curves without a formal test of the difference between the groups. The most commonly used test for this purpose is the log-rank test, which is based on a test statistic defined on a time interval $[0, \tau]$. This statistic can be expressed as

$$Q(\tau) = N_2(\tau) - \int_0^\tau \frac{Z_2(t)}{Z_\bullet(t)} dN_\bullet(t),$$

where $Z_\bullet(t) = Z_1(t) + Z_2(t)$ and $N_\bullet(t) = N_1(t) + N_2(t)$.

Show that, under the null hypothesis of equal survival distributions in the two groups, it holds that $E(Q(\tau)) = 0$.

Hint: You can use that, under the null hypothesis, the counting process $N_j(t)$ had cumulative intensity process $\Lambda_j(t) = \int_0^t Z_j(s)\alpha(s) ds$ where $\alpha(s)$ is the common hazard function under the null hypothesis.

Question 3. We will consider the so-called accelerated failure time (AFT) model. In this model, the distribution of a survival time T conditional on a covariate x is given through the relation

$$\log(T) = \mu - \beta'x + \sigma W$$

where β' is a regression parameter, W is a random variable typically with mean zero, and μ and σ are so-called location and dispersion parameters generating a family of distributions $\mu + \sigma W$.

- (a) Show that the survival function for T given x can be expressed as

$$S(t | x) = S_0(\exp(\beta'x)t)$$

where $S_0(t) = P(\exp(\mu + \sigma W) > t)$.

Give an explanation for the name "AFT-model" based on this representation (1-3 sentences).

- (b) Show that the hazard function of T given x can be written as $\alpha(t | x) = \exp(\beta'x)\alpha_0(\exp(\beta'x)t)$, where $\alpha_0(t) = \alpha(t | 0) = -\frac{d}{dt} \log(S_0(t))$ is the hazard function when $x = 0$.
- (c) Assume now that $S_0(t) = \exp(-bt^k)$ for parameters $b > 0$ and $k > 0$, which means that $\exp(\mu + \sigma W)$ has a Weibull distribution (you do not need to know more about this distribution). Demonstrate that, in this case, the accelerated failure time model is also a proportional hazards model on the form $\alpha(t | x) = h_0(t) \exp(\gamma'x)$ with a constant hazard ratio.

Determine the relation between the coefficient β' in the accelerated failure time model and the coefficient γ' in the proportional hazards model.

Question 4. Consider a chain-binomial SIR model with discrete time dynamics, for a fixed population of size n . Infections are assumed to happen in (discrete) generations. In the following, we assume that once an individual becomes infective, they remain infective for one generation, and then they are immunized, hence, they are removed. Moreover, each susceptible from generation t remains susceptible in generation $t + 1$, if they avoided infection from all infectives of generation t . Infectious events at a given generation t occur independently from each other. Denote the number of susceptibles and infectives in generation t with S_t and I_t , respectively.

- (a) In the *Reed-Frost model*, it is assumed that each infective has the same per contact infection probability p in every generation. Moreover, each potential infection is not affected by any other contacts, neither from previous generations nor in the current generation. That is, if a susceptible individual escaped infections in the previous generations, or was not infected by other contacts in the same generation, they did not develop immunity, nor did they become more susceptible. Under these assumptions, write down the transition probability of the (Reed-Frost) model, that is, derive

$$P(S_{t+1} = s_{t+1}, I_{t+1} = i_{t+1} | S_0 = s_0, \dots, S_t = s_t, I_0 = i_0, \dots, I_t = i_t).$$

Does $\{S_t, I_t\}_{t=0,1,\dots}$ satisfy the Markov property? Justify!

- (b) Suppose you observe a particular chain of the number of infectives $\{i_0, i_1, \dots, i_T = 0\}$. Write down the probability of observing this chain, in other words, what is

$$P(I_0 = i_0, \dots, I_T = 0)?$$

- (c) Denote the final number of infected among the initially susceptible with Z , and the probability of observing $Z = j$ given $S_0 = k, I_0 = i$ with

$$P(Z = j | S_0 = k, I_0 = i) =: m_{ijk}.$$

Argue that the recursive expression

$$m_{ijk} = \binom{k}{j} m_{ijj} q^{(i+j)(k-j)}, \quad j < k, \quad \text{where } m_{ikk} = 1 - \sum_{j=0}^{k-1} m_{ijk}, \quad q = 1 - p,$$

is correct. A derivation of the formula is sufficient but not necessary for obtaining full points.

- (d) Suppose that the number of infectives in the new generation, does not depend on the exact number of the infectives in the previous generation, just on the presence of at least one infective. That is, if there is any infective present in a generation, the probability of infecting a given individual is p if there is any infective present, otherwise it is 0. Under this new assumption, corresponding to the *Greenwood model*, write down the transition probability

$$P(S_{t+1} = s_{t+1}, I_{t+1} = i_{t+1} | S_0 = s_0, \dots, S_t = s_t, I_0 = i_0, \dots, I_t = i_t).$$

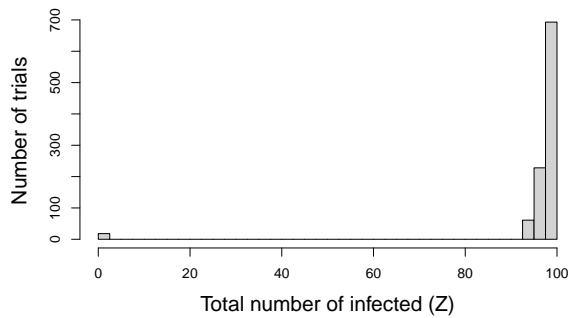
- (e) Give a formula for the reproduction number R_0 both in the Reed-Frost and the Greenwood model. Is it true that in these models, if $R_0 > 1$ the epidemic will always *take off*, that is, a large outbreak will occur? Justify!
- (f) We have simulated $N = 1000$ realizations from both the Reed-Frost and the Greenwood models, with fixed parameters and initial conditions, $I_0 = 1, n = 100$, under

different infection probabilities p . The distribution of the final size of the epidemic, can be estimated with the histograms, depicted below.

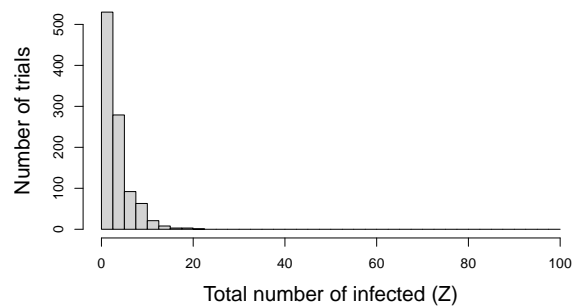
Match the different model and infection probability pairs to the histograms.

- (i) $p = 1 \cdot 10^{-2}$, *Reed-Frost*
- (ii) $p = 2 \cdot 10^{-2}$, *Reed-Frost*
- (iii) $p = 4 \cdot 10^{-2}$, *Reed-Frost*
- (iv) $p = 1 \cdot 10^{-2}$, *Greenwood*
- (v) $p = 4 \cdot 10^{-2}$, *Greenwood*

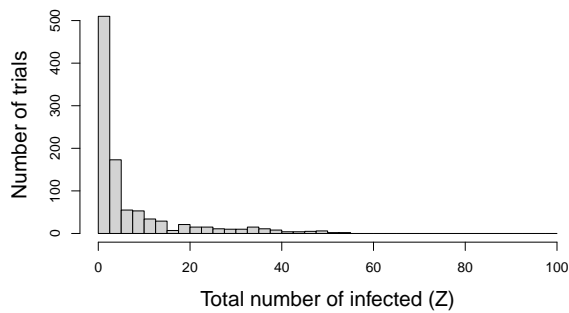
Give a brief justification (3-4 sentences in total) for your choice.



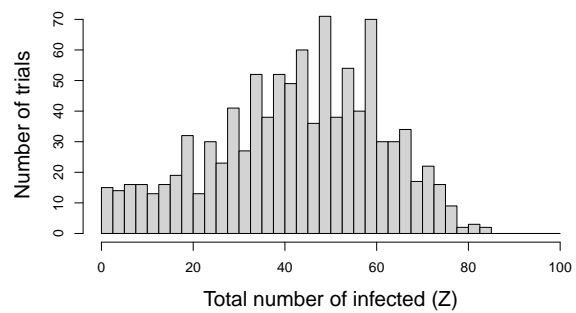
(A)



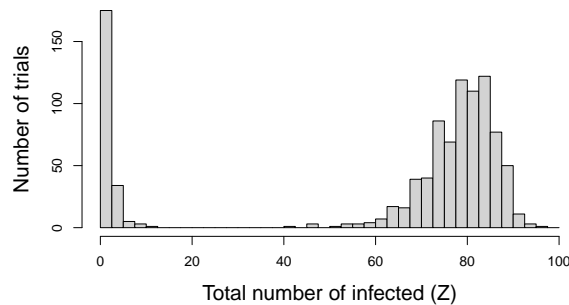
(B)



(C)



(D)



(E)