

## **EXAMINATION PAPER**

Examination Session: May/June

2025

Year:

Exam Code:

MATH2711-WE01

## Title:

Statistical Inference II

Time:	3 hours			
Additional Material provided:	Formula Sheet; Tables: Normal distribution, t-distribution, chi- squared distribution, signed-rank test statistic, rank-sum test statistic.			
Materials Permitted:				
Calculators Permitted:	Yes	Models Permitted: Casio FX83 series or FX85 series.		

Instructions to Candidates:	Answer all questions.				
	Section A is worth 40% and Section B is worth 60%. Within each section, all questions carry equal marks.				
	Write your answer in the white-covered answer booklet with barcodes.				
	Begin your answer to each question on a new page.				

**Revision:** 

## SECTION A

**Q1** Consider the random variables  $X_1, X_2, X_3$ , which are independent and identically exponentially distributed, each with marginal probability density function (p.d.f.)

$$f(x_i) = \lambda e^{-\lambda x_i},$$

for  $x_i > 0$ ,  $\lambda > 0$  and zero otherwise, for i = 1, 2, 3.

(a) Consider the transformations:

$$Y_1 = \frac{X_1}{X_1 + X_2 + X_3}, \qquad Y_2 = \frac{X_2}{X_1 + X_2 + X_3}, \qquad Y_3 = X_1 + X_2 + X_3.$$

Find the joint p.d.f. of  $Y_1, Y_2, Y_3$ .

- (b) Hence find the marginal p.d.f. of  $Y_3$ , and the marginal joint p.d.f. of  $[Y_1, Y_2]^T$ . Where appropriate, identify any standard distributions, clearly stating the corresponding parameters.
- Q2 In 2013, Fox et al. studied the levels of caffeine in coffee samples obtained from around the world. Assume that, given the values of  $\mu$  and  $\tau$ , these measurements are independent observations from a normal distribution with mean  $\mu$  and variance  $\sigma^2 = \tau^{-1}$ . From twenty-six samples, the sample mean caffeine concentration was found to be 12.04mg/g with a sample variance of 1.88.
  - (a) Assume that the value of σ<sup>2</sup> is known to be 2.5 and our prior distribution for μ is normal with mean 10 and standard deviation 1.
    Find the posterior distribution of μ given the data, and give a 95% highest posterior density (HPD) credible interval for mean caffeine concentration. You should state your results to 3 decimal places.
  - (b) Assume now that the value of the precision,  $\tau$ , is unknown but has a Ga(2, 6) prior distribution, and our conditional prior distribution for  $\mu$  given  $\tau$  is normal with mean 10 and variance  $(5\tau)^{-1}$ .

Find the posterior distributions of  $\mu \mid \tau$  and  $\tau$  given the data, and find a 95% HPD posterior credible interval for  $\mu$ . Your interval should not involve  $\tau$ , and you should state your results to 3 decimal places.

Q3 The chi-squared distribution is a continuous probability distribution with p.d.f. given by

$$f(x \mid k) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2},$$

where  $x \in (0, +\infty)$ ,  $k \in \{1, 2, ...\}$  is the degrees-of-freedom parameter and  $\Gamma(\cdot)$  is the gamma function.

- (a) Show that  $f(x \mid k)$  belongs to the 1-parameter exponential family and clearly identify all of the exponential family components.
- (b) Suppose that we observe a sample of independent and identically distributed (i.i.d.) observations  $\mathbf{x} = (x_1, \ldots, x_n)^T$  from the above chi-squared distribution. Show that  $f(\mathbf{x} \mid k)$  also belongs to the 1-parameter exponential family, identifying again all the relevant components.



**Q4** A study compared the performances of engine bearings made of two different types of compounds. Five bearings of each type were tested. The following table gives the times until failure (in units of millions of cycles):

Type A	11.75	13.37	11.33	16.19	13.66
Type B	9.36	11.97	12.48	12.15	10.39

The above data yield the following statistics:  $\bar{x} = 13.26$ ,  $\sum_{i=1}^{5} x_i^2 = 893.9$ , under Type A compound, and  $\bar{y} = 11.27$ ,  $\sum_{i=1}^{5} y_i^2 = 642.2155$ , under Type B compound.

- (a) Assuming that normal distributions with equal variances are appropriate models for the lifetimes of the two types of bearing, test the hypothesis that there is no difference in mean lifetime between the two types, at a significance level of 5%.
- (b) Without assuming normality, we can use the non-parametric Wilcoxon rank sum method to test the hypothesis that the two types of bearing have the same distribution. Perform this test at a significance level of 5%, using the exact tables of the rank-sum statistic.



## SECTION B

**Q5** It has been suggested that the sexes of successive births within families may not be independent, and various models that allow for dependence have ben proposed. One simple model states that there is a probability  $\theta$  that successive births are of the same sex.

In a study involving n three-child families, the frequencies a, b, c, d, e, f, g and h of the eight possible birth sequences and associated probabilities (based on a model in which the probability of a male on the first birth is 0.5 but subsequently depends on the sex of the previous child) are as follows:

Sequence	MMM	MMF	MFM	MFF
Probability	$\frac{1}{2}\theta^2$	$\frac{1}{2}\theta(1-\theta)$	$\frac{1}{2}(1-\theta)^2$	$\frac{1}{2}(1-\theta)\theta$
Observation	a	- b	- c	- d
Sequence	FMM	FMF	FFM	FFF
Probability	$\frac{1}{2}(1-\theta)\theta$	$\frac{1}{2}(1-\theta)^2$	$\frac{1}{2}\theta(1-\theta)$	$\frac{1}{2}\theta^2$
Observation	- e	$\tilde{f}$	$\tilde{g}$	h

(a) Show that the likelihood function for  $\theta$  is proportional to

$$\theta^s (1-\theta)^{2n-s},$$

find an expression for s in terms of the data frequencies, and hence find the maximum likelihood estimate of  $\theta$  in terms of s and n.

- (b) In a particular study involving 6906 three-child families, the observed frequencies were a = 953, b = 914, c = 846, d = 845, e = 825, f = 748, g = 852 and h = 923. Use a large-sample technique to calculate an approximate 95% confidence interval for  $\theta$ .
- (c) What do you infer from your confidence interval about the hypothesis that successive births of the same gender are more probable than successive births of opposite genders, in the above model?



- **Q6** A measurement device reports its measurements with errors, Z, that are normally distributed with known mean zero and an unknown variance v. In an experiment to estimate v, n independent evaluations of the measurement error,  $Z_i$ , are obtained, i = 1, 2, ..., n.
  - (a) A random variable Y has an Inverse Gamma distribution with parameters a and  $b, Y \sim IG(a, b)$ , if it has p.d.f.

$$f(y) = \frac{b^a}{\Gamma(a)} \frac{1}{y^{a+1}} \exp\left(-\frac{b}{y}\right), \qquad y > 0.$$

Show that the Inverse Gamma is the conjugate prior distribution for the unknown measurement error variance, v, of the device, and hence derive expressions for the posterior parameters  $a^*$  and  $b^*$ .

(b) Show that the posterior predictive distribution for the error,  $Z_{n+1}$ , of a further observation made by this device has a p.d.f proportional to

$$\left(1 + \frac{z_{n+1}^2}{2b^*}\right)^{-\left(a^* + \frac{1}{2}\right)}$$

where  $a^*$  and  $b^*$  are the posterior parameter values from part (a).

**Q7** An engineer is interested in ensuring the successful operation of the production lines in a manufacturing plant. Measuring the time between repairs of a given production line as X years, a model of an exponential distribution is agreed, with p.d.f. given by

$$f(x \mid \lambda) = \lambda e^{-\lambda x},$$

Exam code

MATH2711-WE01

where x > 0 and  $\lambda > 0$  is the rate parameter, so that  $\mathbb{E}[X] = 1/\lambda$ . Generally, the expected time between repairs is approximately three months.

- (a) The engineer suspects problems with this production line and records an i.i.d. sample of n inter-repair times  $\mathbf{x} = (x_1, \ldots, x_n)^T$ , intending to test the null hypothesis that the expected time between repairs is indeed three months (1/4 of a year) versus an alternative hypothesis under which the expected time between repairs is two months (1/6 of a year). Express this hypothesis test mathematically in terms of the parameter  $\lambda$ . Derive the most powerful test the engineer can use and show that this is equivalent to a rejection rule of the form  $\bar{x} < k$ , for an appropriate constant k.
- (b) Across the entire manufacturing plant, there are two production lines the original line X and a further line Y. Times between repairs were recorded for each line giving i.i.d. observations  $\mathbf{x} = (x_1, \ldots, x_n)^T$  from Line X and i.i.d. observations  $\mathbf{y} = (y_1, \ldots, y_m)^T$  from Line Y, where  $\mathbf{x}$  and  $\mathbf{y}$  are independent. Inter-repair times follow exponential distributions with respective rate parameters  $\lambda_X$  and  $\lambda_Y$ ; that is,  $x_i \sim \text{Exp}(\lambda_X)$  and  $y_j \sim \text{Exp}(\lambda_Y)$  for  $i = 1, \ldots, n$  and  $j = 1, \ldots, m$ .

Now, the engineer wishes to assess whether the two production lines behave differently; thus, testing  $\mathcal{H}_0: \lambda_X = \lambda_Y$  vs.  $\mathcal{H}_1: \lambda_X \neq \lambda_Y$ . Find the likelihood functions and the MLEs under the null and alternative hypotheses.

- (c) Derive the generalised likelihood ratio test for the hypothesis test in part (b).
- **Q8** Consider the setting where we have the following p.d.f. for  $x \in [0, \infty)$  conditional on parameter  $\lambda > 0$ :

$$f(x \mid \lambda) = \frac{3}{\lambda} x^2 e^{-x^3/\lambda}.$$

We are given a sample  $\mathbf{x} = (x_1, \dots, x_n)^T$  of *n* observations which we assume to be conditionally i.i.d. given parameter  $\lambda$ .

- (a) Consider the hypothesis test  $\mathcal{H}_0 : \lambda = 1$  vs.  $\mathcal{H}_1 : \lambda = 2$  and derive the Bayes factor in favour of  $\mathcal{H}_0$  against  $\mathcal{H}_1$ .
- (b) Assume now that we want to test the null hypothesis  $\mathcal{H}_0 : \lambda = 1$  vs. a more general alternative of the form  $\mathcal{H}_1 : \lambda > 0$ . Under the alternative we specify an inverse-Gamma distribution for parameter  $\lambda$ , with parameters a = b = 1. Derive the Bayes factor in favour of  $\mathcal{H}_0$  against  $\mathcal{H}_1$ .
- (c) Consider the framework where we want to compare a model  $\mathcal{M}_0$  with an inverse-Gamma prior on  $\lambda$ , as defined in part (b), versus a model  $\mathcal{M}_1$  with an inverse-Gamma prior with parameters a and b such that the prior mode and mean are b/(a+1) = 1 and b/(a-1) = 3, respectively. Find the corresponding prior distribution under  $\mathcal{M}_1$  and derive the Bayes factor in favour of  $\mathcal{M}_0$  against  $\mathcal{M}_1$ .