

Student Number

Semester 2 Final Exam, 2018

School of Mathematics and Statistics

MAST20005 Statistics

Writing time: 3 hours

Reading time: 15 minutes

This is NOT an open book exam

Common content with: MAST90058 Elements of Statistics

This paper consists of 6 pages (including this page)

Authorised Materials

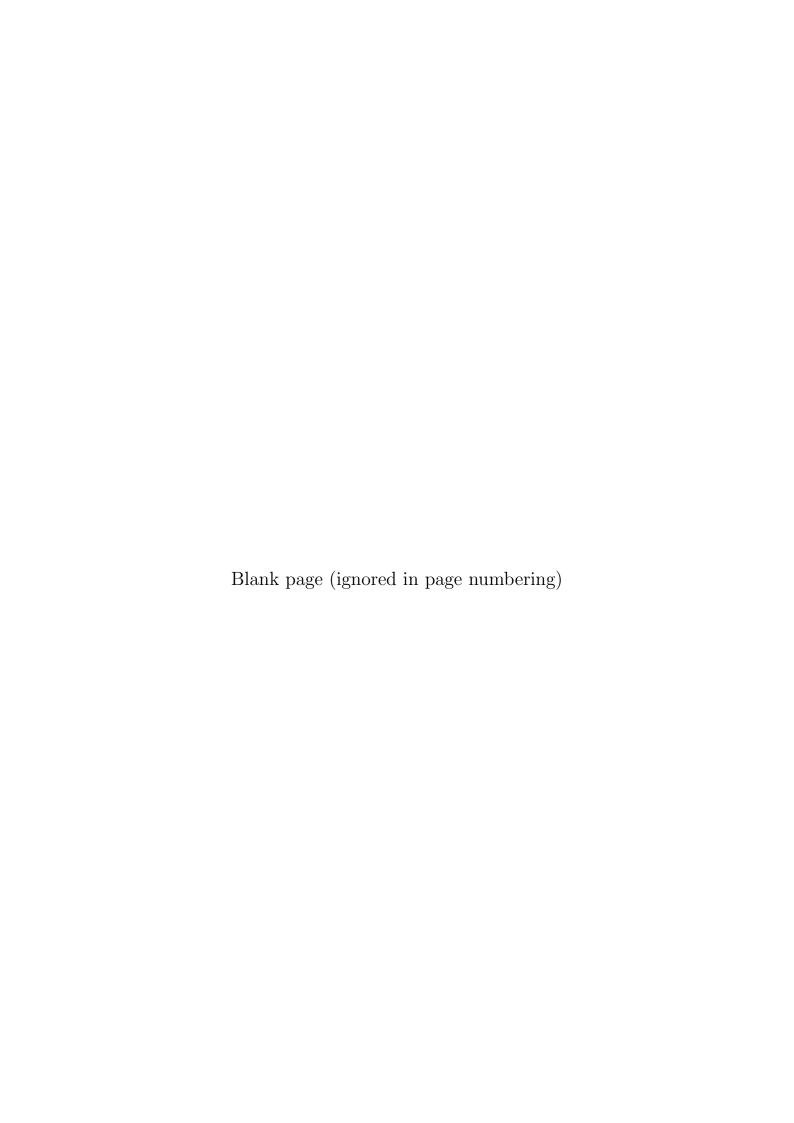
- Mobile phones, smart watches and internet or communication devices are forbidden.
- Only Casio FX82 (any suffix) calculators may be used.
- Students may bring one double-sided A4 sheet of handwritten notes.

Instructions to Students

- You must NOT remove this question paper at the conclusion of the examination.
- Some useful R output is given in an appendix on the last page.
- You should attempt all questions. Show full working for each of your answers.
- There are 8 questions with marks as shown. The total number of marks available is 90.

Instructions to Invigilators

- Students must NOT remove this question paper at the conclusion of the examination.
- All graphics or CAS calculators should be confiscated.
- Students may use one double-sided A4 sheet of handwritten notes.



Question 1 (10 marks) You have random samples from two groups, $X \sim N(\mu_1, \sigma_1^2)$ and $Y \sim N(\mu_2, \sigma_2^2)$. Some R output from analysing these data are below. The variable **x** contains the observations of X and the variable **y** contains the observations of Y.

```
> summary(x)
  Min. 1st Qu.
                 Median
                            Mean 3rd Qu.
                                            Max.
  2.620
          3.320
                  4.040
                           4.193
                                            6.300
                                   4.515
> summary(y)
  Min. 1st Qu. Median
                           Mean 3rd Qu.
                                            Max.
  4.150
          4.450
                  5.355
                           5.212
                                   5.595
                                            6.720
> sd(x)
[1] 1.200121
> sd(y)
[1] 0.8368034
> sort(y)
[1] 4.15 4.25 4.34 4.78 5.30 5.41 5.46 5.64 6.07 6.72
```

- (a) For each of the following quantities, state or calculate its value if possible, or otherwise explain why it is not possible.
 - (i) $x_{(1)}$
 - (ii) $y_{(4.5)}$
 - (iii) \bar{x}
 - (iv) σ_2
- (b) For each of the following statements, state whether they are true, false or if it is not possible to know from the given information. In each case state the values of the quantities, if possible.
 - (i) $x_{(1)} > y_{(1)}$
 - (ii) $\bar{x} > \bar{y}$
 - (iii) $\sigma_1 > \sigma_2$
- (c) For each of the following pairs of hypotheses, carry out the test if it is possible, using a 5% significance level, or otherwise explain what further information you need in order to do it.
 - (i) $H_0: \mu_1 = \mu_2 \text{ versus } H_1: \mu_1 \neq \mu_2$
 - (ii) H_0 : $\sigma_2 = 2$ versus H_1 : $\sigma_2 \neq 2$

Question 2 (9 marks) A random sample on X produced the following observations:

5.5 5.8 6.0 6.6 6.8 6.9 7.1 7.3 7.5 8.7

For these data, we have $\bar{x} = 6.82$ and s = 0.932.

- (a) Let $\mu = \mathbb{E}(X)$. Calculate a 95% confidence interval for μ , assuming that X is normally distributed.
- (b) Let $p = \Pr(X > 6.5)$. Calculate a 95% confidence interval for p.
- (c) Let m be the median of X. Calculate a distribution-free confidence interval for m, with an approximate confidence level of 90%.

Question 3 (11 marks) Consider a random sample of size n on X which has a geometric distribution with parameter θ . Its pmf is:

$$p_X(x) = \theta(1-\theta)^x, \quad x \in \{0, 1, 2, \dots\}$$

and it has mean $(1-\theta)/\theta$.

- (a) Determine a sufficient statistic for θ .
- (b) Find the method of moments estimator of θ .
- (c) Find the maximum likelihood estimator (MLE) of θ .
- (d) Find the Cramér–Rao lower bound for unbiased estimators of θ .
- (e) Derive an expression for the standard error of the MLE.
- (f) A random sample of size n = 20 produced the following observations:

```
3\ 2\ 0\ 1\ 1\ 3\ 0\ 0\ 0\ 0\ 0\ 3\ 0\ 3\ 1\ 1\ 0\ 2\ 2\ 0
```

Estimate θ and calculate an approximate 90% confidence interval.

Question 4 (12 marks) Laleh has bought a new toaster. It has a dial that allows her to set the 'strength' of toasting. She is not sure if it works very well and decides to run some experiments. She sets the dial to various values, x, and measures how long the toaster cooks the bread, Y, before it pops the bread out. She does a simple linear regression analysis of these data using the model $\mathbb{E}(Y \mid x) = \alpha + \beta x$. Some partial R output from her analysis is shown below.

- (a) How many experiments did Laleh carry out?
- (b) Carry out the following hypothesis tests, using a 5% significance level.
 - (i) $H_0: \alpha = 0$ versus $H_1: \alpha \neq 0$
 - (ii) H_0 : $\beta = 0$ versus H_1 : $\beta \neq 0$
- (c) Is there evidence that the dial is having an effect on the toasting time?
- (d) Write out the ANOVA table for this regression model fit. (Hint: the F statistic is shown in the above R output.)

Question 5 (13 marks) On his way to work in the mornings, Damjan records how long he has to wait for his train at the station. Over a series of days, he observes the following times (in minutes):

$$4.8 \quad 1.2 \quad 3.7 \quad 0.9 \quad 0.7 \quad 0.3 \quad 0.9 \quad 3.2 \quad 1.4$$

For these data we have $\bar{x} = 1.9$ and s = 1.58. Damjan decides to use an exponential distribution with mean θ as a model for these data. He would like to estimate his median waiting time, m.

- (a) Express m in terms of θ .
- (b) Damjan decides to use the sample median, \hat{M} , as his estimator.
 - (i) What is the asymptotic sampling distribution of \hat{M} ?
 - (ii) What is Damjan's estimate of m for this dataset?
 - (iii) Calculate a standard error for Damjan's estimate.
- (c) Damjan now considers using \bar{X} as his estimator of m.
 - (i) Show that this estimator is biased.
 - (ii) Let $T = c\bar{X}$ be an adjusted estimator. Find c so that T is unbiased.
 - (iii) Determine var(T).
 - (iv) Which of \hat{M} and T is the better estimator?
 - (v) What is the estimate, t, based on the data above?
 - (vi) Calculate a standard error for this estimate.

Question 6 (10 marks) On his way to the office early every morning, Allan walks past South Lawn and counts how many students he sees there. He decides to model these counts using a Poisson distribution with pmf,

$$p(x \mid \theta) = \frac{e^{-\theta}\theta^x}{r!}, \quad (x = 0, 1, \dots).$$

Across the first 45 days of semester, he observes on average 3.8 students per day. Robert says that he did a similar survey last year and got on average about 3.0 students per day, although he cannot remember over how many days he observed them. Allan would like to estimate θ by combining this information appropriately.

(a) Show that the gamma distribution is a conjugate prior for θ . Note that the pdf of $\theta \sim \text{Gamma}(\alpha, \beta)$ is,

$$f(\theta|\alpha,\beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}, \quad (\theta > 0)$$

and it has mean $\mathbb{E}(\theta) = \alpha/\beta$.

- (b) Allan decides to treat Robert's information as if it came from 10 days of sampling (i.e. as pseudodata). Determine the parameters of the prior that encode this information appropriately.
- (c) Using this prior, what is the posterior distribution of θ ?
- (d) Calculate the posterior mean.

Question 7 (13 marks) Ben runs a sports program for high school students. Each student that enters his program chooses one of the following three sports: volleyball, basketball and netball. The heights of the first 15 students, and the sports they chose, were:

Volleyball	183, 176, 170, 179
Basketball	165, 169, 171, 154, 165, 159
Netball	167, 160, 177, 173, 170

Ben runs an analysis of variance with these data. A partially complete ANOVA table from his analysis is given below.

Source	df	SS	MS	F
Treatment (sport) Error			38	
Total		872.4		

- (a) Is there evidence of a relationship between the students' heights and their choice of sport?
- (b) What is the sampling distribution of $\hat{\sigma}^2$?
- (c) Let μ_1 be the population mean of students who choose volleyball. What is the sampling distribution of $\hat{\mu}_1$?
- (d) Let X^* be the height of the next student who chooses volleyball. Show that,

$$X^{\star} - \hat{\mu}_1 \sim N\left(0, \frac{5\sigma^2}{4}\right).$$

(e) Calculate a 95% prediction interval for X^* .

Question 8 (12 marks) Each person has one of the genotypes A, B or C. According to the Hardy-Weinberg law in genetics, these three genotypes should occur in the population in the proportions θ^2 , $2\theta(1-\theta)$ and $(1-\theta)^2$, respectively, for some $\theta \in [0,1]$. In a sample of 600 individuals from the population, you observe 27 individuals with genotype A, 186 individuals with genotype B and 387 individuals with genotype C.

- (a) Find the maximum likelihood estimate of θ .
- (b) Calculate a standard error for this estimate.
- (c) Carry out a hypothesis test of the Hardy-Weinberg law, using a significance level of 5%.

End of exam questions—Total Available Marks = 90
Turn the page for appended material

Appendix (R output)

```
> p1 <- c(0.01, 0.025, 0.05, 0.1, 0.9, 0.95, 0.975, 0.99)
> qnorm(p1)
[1] -2.326 -1.960 -1.645 -1.282 1.282 1.645 1.960 2.326
> qt(p1, df = 8)
[1] -2.896 -2.306 -1.860 -1.397 1.397 1.860 2.306 2.896
> qt(p1, df = 9)
[1] -2.821 -2.262 -1.833 -1.383 1.383 1.833 2.262 2.821
> qt(p1, df = 10)
[1] -2.764 -2.228 -1.812 -1.372 1.372 1.812 2.228 2.764
> qt(p1, df = 11)
[1] -2.718 -2.201 -1.796 -1.363 1.363 1.796 2.201 2.718
> qt(p1, df = 12)
[1] -2.681 -2.179 -1.782 -1.356 1.356 1.782 2.179 2.681
> qt(p1, df = 13)
[1] -2.650 -2.160 -1.771 -1.350 1.350 1.771 2.160 2.650
> qt(p1, df = 14)
[1] -2.624 -2.145 -1.761 -1.345 1.345 1.761 2.145 2.624
> qt(p1, df = 28)
[1] -2.467 -2.048 -1.701 -1.313 1.313 1.701 2.048 2.467
> qt(p1, df = 29)
[1] -2.462 -2.045 -1.699 -1.311 1.311 1.699 2.045 2.462
> qchisq(p1, df = 1)
[1] 0.0001571 0.0009821 0.0039321 0.0157908 2.7055435 3.8414588 5.0238862 6.6348966
> qchisq(p1, df = 2)
 \hbox{\tt [1]} \ \ 0.02010 \ \ 0.05064 \ \ 0.10259 \ \ 0.21072 \ \ 4.60517 \ \ 5.99146 \ \ 7.37776 \ \ 9.21034 
> qchisq(p1, df = 3)
 \begin{bmatrix} 1 \end{bmatrix} \quad 0.1148 \quad 0.2158 \quad 0.3518 \quad 0.5844 \quad 6.2514 \quad 7.8147 \quad 9.3484 \quad 11.3449 
> qchisq(p1, df = 8)
[1] 1.646 2.180 2.733 3.490 13.362 15.507 17.535 20.090
> qchisq(p1, df = 9)
[1] 2.088 2.700 3.325 4.168 14.684 16.919 19.023 21.666
> qchisq(p1, df = 10)
[1] 2.558 3.247 3.940 4.865 15.987 18.307 20.483 23.209
> qf(p1, 1, 13)
 \hbox{\tt [1]} \ \ 0.0001632 \ \ 0.0010206 \ \ 0.0040868 \ \ 0.0164196 \ \ 3.1362051 \ \ 4.6671927 \ \ 6.4142543 \ \ 9.0738057 
> qf(p1, 1, 14)
 \hbox{\tt [1]} \ \ 0.0001628 \ \ 0.0010178 \ \ 0.0040756 \ \ 0.0163739 \ \ 3.1022134 \ \ 4.6001099 \ \ 6.2979386 \ \ 8.8615927 
> qf(p1, 2, 12)
[1] 0.01006 0.02537 0.05151 0.10629 2.80680 3.88529 5.09587 6.92661
> qf(p1, 2, 13)
[1] 0.01006 0.02537 0.05150 0.10622 2.76317 3.80557 4.96527 6.70096
> qf(p1, 3, 11)
[1] 0.03686 0.06957 0.11411 0.19148 2.66023 3.58743 4.63002 6.21673
> qf(p1, 3, 12)
[1] 0.03697 0.06975 0.11436 0.19173 2.60552 3.49029 4.47418 5.95254
> pbinom(0:5, 10, 0.5)
[1] 0.0009766 0.0107422 0.0546875 0.1718750 0.3769531 0.6230469
> pbinom(0:5, 9, 0.5)
[1] 0.001953 0.019531 0.089844 0.253906 0.500000 0.746094
```



Library Course Work Collections

Author/s:

Mathematics and Statistics

Title:

Statistics, 2018, Semester 2, MAST20005

Date:

2018

Persistent Link:

http://hdl.handle.net/11343/220963