



## EXAMINATION PAPER

<b>Examination Session:</b> May/June	<b>Year:</b> 2020	<b>Exam Code:</b> MATH2041-WE01
---	----------------------	------------------------------------

<b>Title:</b> Statistical Concepts II
--

Time (for guidance only):	3 hours	
Additional Material provided:	Tables: Normal distribution, t-distribution, chi-squared distribution, signed-rank test statistic, rank-sum test statistic.	
Materials Permitted:		
Calculators Permitted:	Yes	Models Permitted: There is no restriction on the model of calculator which may be used.

Instructions to Candidates:	<p>Credit will be given for your answers to all questions. All questions carry the same marks.</p> <p>Please start each question on a new page. Please write your CIS username at the top of each page.</p> <p>Show your working and explain your reasoning.</p>	
	<b>Revision:</b>	

Page number

2 of 7

Exam code

MATH2041-WE01

- Q1 1.1** Suppose we wish to test a null hypothesis  $H_0$  against an alternative hypothesis  $H_1$ , based on an observation  $Y$  which takes values 1, 2, 3, 4, 5. The probability distribution of  $Y$  under each hypothesis is as follows.

$y$	1	2	3	4	5
$\mathbb{P}[Y = y \mid H_0]$	0.05	0.1	0.15	0.3	0.4
$\mathbb{P}[Y = y \mid H_1]$	0.15	0.25	0.3	0.2	0.1

- (i) Suppose we seek a test which minimises the quantity

$$3\alpha + 2\beta,$$

where  $\alpha$  is the probability of making a Type I error and  $\beta$  is the probability of making a Type II error. Find the optimal test in this situation, and calculate the corresponding values of  $\alpha$  and  $\beta$ . *You should state, without proof, any general results that you require in order to demonstrate that the test is optimal.*

- (ii) An alternative version of this test is more concerned about the occurrence of a Type I error, and so places five times more weight on  $\alpha$  than  $\beta$ . Explain briefly why the resulting test may not be useful for this problem.
- (iii) For a fixed significance level of 0.15, find the most powerful test of these hypotheses and calculate the power of your test.
- 1.2** The average consumption of cigarettes for a random sample of 9 smokers during each month before and after their enrolment in a quit-smoking mentorship program are recorded in the table below.

Before	22	15	10	21	14	3	11	8	15
After	14	8	0	22	9	3	8	6	11

The researcher is interested in finding out whether the enrolment in a quit-smoking mentorship program causes a decline in cigarette smoking.

Use a Wilcoxon signed rank test to test this hypothesis at the 0.05 level of significance. Perform your test using:

- (i) tables of exact critical values, and
- (ii) using the normal approximation to the test statistic.

**Q2 2.1** Show that for a general independent sample of size  $n$ , denoted  $X_1, \dots, X_n$ , from a distribution  $f(x | \theta)$  with unknown parameter  $\theta$ , the expected Fisher information for  $\theta$  from the sample of  $n$  observations,  $\mathcal{I}_n(\theta)$ , can be expressed as

$$\mathcal{I}_n(\theta) = n\mathcal{I}_1(\theta),$$

where  $\mathcal{I}_1(\theta)$  is the expected Fisher information from a sample of size 1.

**2.2** The pdf for a random variable  $Y \sim N(\mu, v)$  is

$$f(y | \mu, v) = \frac{1}{\sqrt{2\pi v}} e^{-\frac{1}{2v}(y-\mu)^2}.$$

Suppose that  $X_1, \dots, X_n$  form an independent sample from a normal distribution with mean 0 and variance  $v$ .

- (i) Find a sufficient statistic for  $v$  given data  $(x_1, \dots, x_n)$ . Justify your answer.
- (ii) Derive the maximum likelihood estimator  $\hat{v}$  for  $v$ . Is  $\hat{v}$  an unbiased estimator for  $v$ ?
- (iii) Using the sampling distribution for  $\hat{v}$ , derive an exact  $100(1 - \alpha)\%$  confidence interval for  $v$  in terms of  $\hat{v}$ .
- (iv) Find the expected Fisher information for  $v$ .
- (v) Hence, for large  $n$ , calculate an approximate 95% confidence interval for  $v$  in terms of  $\hat{v}$ . Clearly state each approximation that is made in order to construct the interval.

**Q3** The number of large meteorites that strike the Earth is monitored by a particular observatory. The observatory records the annual counts of meteorite strikes,  $X$ , which can be described by a Poisson distribution with an unknown parameter  $\lambda$ .

**3.1** Given a random sample  $X_1, \dots, X_{50}$  of  $n = 50$  years of data, consider testing the null hypothesis  $H_0 : \lambda = 2$  against the alternative hypothesis  $H_1 : \lambda = 3$ .

- (i) Show that the optimal test of these hypotheses at significance level  $\alpha$  can be expressed in the form “Reject  $H_0$  if  $T > k^*$ ”, where  $T = \sum_{i=1}^n X_i$  and  $k^* > 0$  is a constant determined by  $\alpha$  and the sampling distribution of  $T$ .
- (ii) Find, approximately, the critical value  $k^*$  for the test in part (i) using a significance level of 0.002, and find the approximate power of the test with this critical value.

**3.2** A further independent sample  $Y_1, \dots, Y_{50}$  of size 50 of meteorite impacts is recorded from another observatory using a more sensitive detector.  $Y$  is also Poisson, but with an unknown parameter  $\mu$ , where it is believed that  $2\mu = 3\lambda$  to reflect the increased sensitivity.

Using both samples and assuming  $2\mu = 3\lambda$ :

- (i) Find the maximum likelihood estimator for  $\lambda$ .
- (ii) Show that the generalised likelihood ratio test statistic for the hypotheses  $H_0 : \lambda = 2$  against the alternative  $H_1 : \lambda \neq 2$  is of the form:

$$e^{S-250} \left( \frac{250}{S} \right)^S,$$

where  $S = \sum_{i=1}^{50} X_i + \sum_{j=1}^{50} Y_j$ .

- (iii) Hence find a simplified expression for the approximate large-sample generalised likelihood ratio test statistic for these hypotheses, and the corresponding critical value at significance level 0.05.

**Q4** The beta distribution, with parameters  $a > 0, b > 0$  has probability density function

$$f(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}, \quad 0 < x < 1.$$

- 4.1** Show that if the random quantity  $X$  has a beta distribution, with parameters  $a, b$ , then

$$E(X) = \frac{a}{a+b}, \quad Var(X) = \frac{ab}{(a+b)^2(a+b+1)}.$$

- 4.2** Explain what it means to say that a family of probability distributions is a conjugate prior family for sampling from a particular family of distributions. Show that the family of beta distributions is a conjugate prior family for sampling from the binomial distribution.
- 4.3** Suppose that a particular medical treatment has a probability  $q$  of success for a randomly selected patient with a particular disease. Suppose that the prior distribution for  $q$  is a beta distribution with parameters  $a = b = 4$ . Suppose that, in a random sample of 60 patients, the treatment is found to be successful for 36 of the patients. Find the posterior probability distribution for  $q$ . Therefore, find the probability that the next two patients that are treated will both have successful outcomes.
- 4.4** Explain what is meant by an  $\alpha$  level credible interval for  $q$ . Evaluate an approximate posterior credible interval for  $q$ , in the problem of question **4.3**, at level 0.95.
- 4.5** State the form of the limiting posterior distribution for sampling from a general likelihood as the sample size increases. Find the limiting posterior form for the problem of question **4.3**. For the given data, compare the mean and standard deviation of this approximation to the exact values for these quantities for the posterior distribution calculated in question **4.3**.

**Q5** The 4 C's (cut, colour, clarity, and carat weight) are the main criteria used to judge the quality of a diamond and determine its value (price). In order to understand how carat weight affects the price of a diamond, a random sample of 20 diamonds were obtained, and their carat weight and price (in dollars) are given below.

carat	1.02	1.54	0.66	0.37	1.08	0.72	0.53	0.31	0.71	0.41
price	5666	9562	2648	791	4685	2843	1948	680	2946	964
carat	0.91	0.9	0.5	1	1.37	0.33	0.53	0.64	0.48	0.54
price	3639	4229	963	6272	8577	594	1678	4281	990	1754

**5.1** Consider the simple linear regression model for the price on carat weight,  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ . Let the least-squares regression line be  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ .

- State the assumptions of the usual simple linear regression model with normal errors.
- Give expressions for  $\hat{\beta}_0$  and  $\hat{\beta}_1$ . Prove that they are unbiased estimators of  $\beta_0$  and  $\beta_1$ .

**5.2** A simple linear regression model was fitted to the above data using R, as shown on the next page. Use the data and the R output to answer the following.

- What are the estimates of the intercept and slope for the above data? Interpret the t-values for the coefficients. Give a 95% confidence interval for the slope.
- Explain how the ‘Multiple R-squared’ value is calculated, and interpret the corresponding value for this regression from the R output.
- For a diamond with carat weight  $x^* = 0.75$ , estimate the expected value of price,  $y$ , and give a 95% confidence interval for the expected value.
- Explain how the residuals are calculated. Use the residual plot to assess informally the adequacy of the linear regression model as a summary description of these data.

Page number

7 of 7

Exam code

MATH2041-WE01

```
> lmPrice = lm(price ~ carat, data=diamonds1)
> summary(lmPrice)
```

Call:

```
lm(formula = price ~ carat, data = diamonds1)
```

Residuals:

Min	1Q	Median	3Q	Max
-1212.31	-347.70	-53.19	254.39	1643.82

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-2104.8	354.1	-5.944	1.26e-05 ***
carat	7409.4	442.6	16.742	2.02e-12 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 659.4 on 18 degrees of freedom

Multiple R-squared: 0.9397, Adjusted R-squared: 0.9363

F-statistic: 280.3 on 1 and 18 DF, p-value: 2.023e-12

```
> mean(carat)
[1] 0.7275
> var(carat)
[1] 0.1168303
> plot(resid(lmPrice) ~ carat)
> abline(0,0)
```

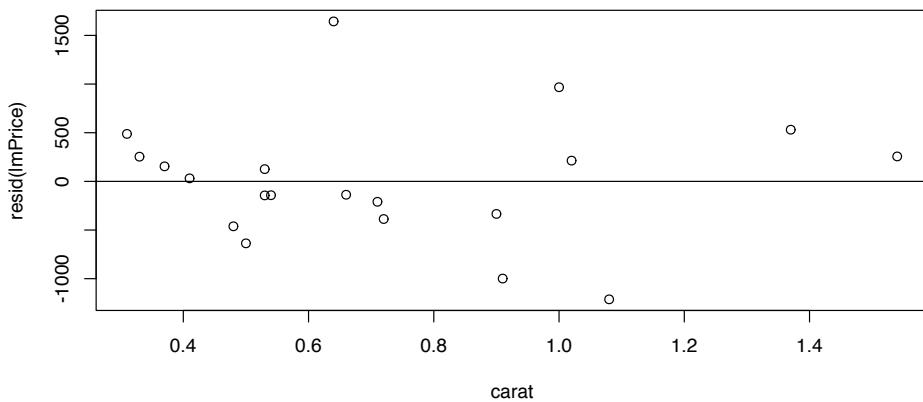


Table A: Probabilities for the standard normal distribution

Table entry for  $z$  is the probability lying to the left of  $z$  for a  $N(0, 1)$  distribution.

$$\Phi(z) = P[Z < z]$$

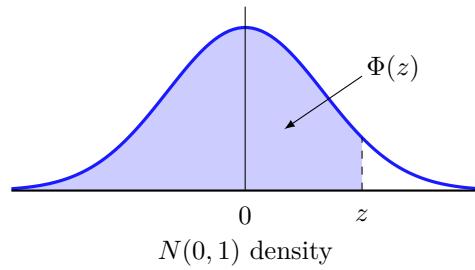


Table B: Probabilities for the  $t$ -distribution

Table entry for  $p$  and  $C$  is the critical value  $t^*$  with probability  $p$  lying to its right and probability  $C$  lying between  $-t^*$  and  $t^*$  for a  $t$ -distribution with  $k$  degrees of freedom.

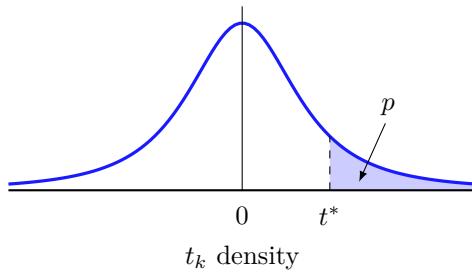


Table C: Probabilities for the  $\chi^2$ -distribution

Table entry for  $p$  is the point  $\chi^*$  with probability  $p$  lying above it for a  $\chi^2$  distribution with  $k$  degrees of freedom.

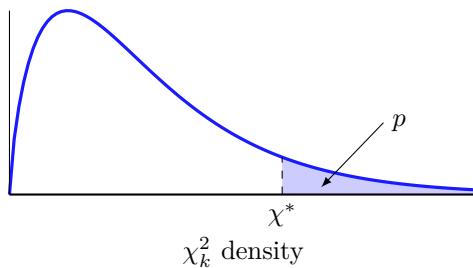


Table D: Critical values for the Rank-sum Test

Reject the hypothesis of identical populations in a two-sided test at level  $\alpha$  if the test statistic  $W$  from a group of size  $n$  is *less than* the value  $T_L$  shown in the following table, or greater than the value  $T_U$  where

$$T_U = n(n + m + 1) - T_L.$$

		$m$									
$\alpha = 0.01$		2	3	4	5	6	7	8	9	10	
2	-	-	-	-	-	-	-	-	-	-	-
3	-	-	-	-	-	-	-	-	7	7	
4	-	-	-	-	11	11	12	12	13		
5	-	-	-	16	17	17	18	19	20		
$n$	6	-	-	22	23	24	25	26	27	28	
	7	-	-	29	30	32	33	35	36	38	
	8	-	-	38	39	41	43	44	46	48	
	9	-	46	47	49	51	53	55	57	59	
	10	-	56	58	60	62	65	67	69	72	

		$m$									
$\alpha = 0.05$		2	3	4	5	6	7	8	9	10	
2	-	-	-	-	-	-	-	4	4	4	
3	-	-	-	7	8	8	9	9	10		
4	-	-	11	12	13	14	15	15	16		
5	-	16	17	18	19	21	22	23	24		
$n$	6	-	23	24	25	27	28	30	32	33	
	7	-	30	32	34	35	37	39	41	43	
	8	37	39	41	43	45	47	50	52	54	
	9	46	48	50	53	56	58	61	63	66	
	10	56	59	61	64	67	70	73	76	79	

		$m$									
$\alpha = 0.10$		2	3	4	5	6	7	8	9	10	
2	-	-	-	4	4	4	5	5	5	5	
3	-	7	7	8	9	9	10	11	11		
4	-	11	12	13	14	15	16	17	18		
5	16	17	18	20	21	22	24	25	27		
$n$	6	22	24	25	27	29	30	32	34	36	
	7	29	31	33	35	37	40	42	44	46	
	8	38	40	42	45	47	50	52	55	57	
	9	47	50	52	55	58	61	64	67	70	
	10	57	60	63	67	70	73	76	80	83	

Table E: Critical values for the Signed-rank Test

Reject the null hypothesis of identical populations if the test statistic  $V$  is *less than* the value  $T$  shown in the following table.

Sample size $n$	Level of significance for a two-sided test					
	20%	10%	5%	2%	1%	0.5%
	10%	5%	2.5%	1%	0.5%	0.25%
Level of significance for a one-sided test						
5	3	1	-	-	-	-
6	4	3	1	-	-	-
7	6	4	3	1	-	-
8	9	6	4	2	1	-
9	11	9	6	4	2	1
10	15	11	9	6	4	2
11	18	14	11	8	6	4
12	22	18	14	10	8	6
13	27	22	18	13	10	8
14	32	26	22	16	13	10
15	37	31	26	20	16	13
16	43	36	30	24	20	16
17	49	42	35	28	24	20
18	56	48	41	33	28	24
19	63	54	47	38	33	28
20	70	61	53	44	38	33
21	78	68	59	50	43	38
22	87	76	66	56	49	43
23	95	84	74	63	55	49
24	105	92	82	70	62	55
25	114	101	90	77	69	61
26	125	111	99	85	76	68
27	135	120	108	93	84	75
28	146	131	117	102	92	83
29	158	141	127	111	101	91
30	170	152	138	121	110	99