

THE UNIVERSITY OF SYDNEY  
SCHOOL OF MATHEMATICS AND STATISTICS

**Quiz 2a**

---

MATH1905: Statistics Advanced

Semester 1, 2017

---

Full Name..... SID.....  
Day..... Time ..... Room.....  
Tutor..... Signature.....

**Time allowed: 40 minutes**

1. **This quiz is closed book. You may not use the computer.**
2. Full marks will only be given if you obtain the correct answer **and** your working is sufficient to justify your answer.
3. Partial marks may be awarded for working.
4. Please write carefully and legibly.
5. All of your answers should be written using ink and **not** pencil, with your final answer placed in the answer box.
6. All working must be done on the quiz paper in the indicated space.
7. Each question is worth **2 marks**.
8. Only University of Sydney approved calculators may be used (must have a sticker).
9. All pages (including working) of the quiz paper must be handed in at the end of the quiz.

This quiz paper has 12 pages (this cover sheet + 10 pages of questions + 1 page of statistical formulae) and 10 questions.

1. Let  $X_1, X_2, \dots, X_{39}$  be a random sample of geometric random variables with expectation 4 and variance 12. A normal approximation to  $P\left(113 \leq \sum_{i=1}^{39} X_i < 188\right)$ , with continuity correction, is of the form  $P(a \leq Z \leq b)$  (where  $Z \sim N(0, 1)$ ) for some constants  $a$  and  $b$ . Write these down (to 2 decimal places) where indicated in the boxes below.

|   |
|---|
| $a =$<br><br><br><div style="text-align: center;">-2.01</div> |
|---|

|  |
|--|
| $b =$<br><br><br><div style="text-align: center;">1.46</div> |
|--|

*Working for Question 1.*

Writing  $S = \sum_{i=1}^{39} X_i$  we have  $E(S) = 39 \times 4 = 156$  and  $Var(S) = 39 \times 12 = 468$ . Since  $S$  is integer-valued we can rewrite the initial probability as

$$P(113 \leq S < 188) = P(112.5 \leq S \leq 187.5)$$

and it is this last version we approximate by replacing  $S$  with  $Y \sim N(156, 468)$ . The approximation is thus

$$\begin{aligned} P(112.5 \leq Y \leq 187.5) &= P\left(\frac{112.5 - 156}{\sqrt{468}} \leq \frac{Y - 156}{\sqrt{468}} \leq \frac{187.5 - 156}{\sqrt{468}}\right) \\ &\approx P(-2.01 \leq Z \leq 1.46) \quad \text{where } Z = (Y - 156)/\sqrt{468} \sim N(0, 1) \end{aligned}$$

**The next two questions relate to the following scenario:** A random sample of size 20 is taken from a population with variance known to be 2.6 and unknown mean  $\mu$ . The following R output may be useful

```
> qnorm(c(0.8, 0.9, 0.95, 0.975, 0.98, 0.99, 0.995))
```

```
[1] 0.8416212 1.2815516 1.6448536 1.9599640 2.0537489 2.3263479 2.5758293
```

2. The p-value for a two-sided  $Z$ -test of  $H_0: \mu = 46$  is given by 0.646. Given that  $P(Z \leq 0.46) = 0.677$  (where  $Z \sim N(0, 1)$ ), determine to 3 decimal places the value taken by the sample mean given that it is greater than 46.

Answer is

46.166

---

*Working for Question 2.*

Whatever the value of the sample mean  $\bar{x}$ , the observed value of the test statistic is

$$z = \frac{\bar{x} - 46}{\sqrt{2.6/20}}$$

and the p-value is  $2P(Z \geq |z|)$  where  $Z \sim N(0, 1)$ . Thus

$$2P(Z \geq |z|) = 0.646$$

$$P(Z \geq |z|) = 0.323$$

$$P(Z \leq |z|) = 1 - 0.323 = 0.677$$

Thus using the information above,  $|z| = 0.46$ .

Note that  $|z| = z$  since we are told the difference  $\bar{x} - 46$  is positive. Thus

$$\bar{x} = 46 + 0.46\sqrt{2.6/20} \approx 46.1659$$

3. A 90% confidence interval for  $\mu$  is of the form  $\bar{x} \pm a$  where  $\bar{x}$  is your answer to the previous question. Write down the value of  $a$  (to 3 decimal places) in the box below.

Ans

0.593

---

*Working for Question 3.*

The value of  $a$  is  $c \times$  s.e. where

- s.e. (the standard error) is  $\sigma/\sqrt{n} = \sqrt{2.6/20}$ ;
- the table value  $c$  satisfies  $P(Z > c) = \alpha/2$  where  $100(1 - \alpha)\% = 90\%$  is the confidence level, i.e.  $\alpha = 0.1$ .

The table value thus satisfies  $P(Z \geq c) = 0.05$  and the R output above shows us that  $c = 1.645$ . Thus

$$a = 1.645 \times \sqrt{2.6/20} \approx 0.5931$$

**The following two questions relate to the following scenario:** A random sample of size 9 is taken from a population, assumed to be  $N(\mu, \sigma^2)$  for  $\mu$  and  $\sigma^2$  both unknown. The sample mean takes the value 13.37 while the sample standard deviation takes the value 5.686. The following R output may be useful:

```
> n
[1] 9

> qt(c(0.9, 0.95, 0.975, 0.99, 0.995), df=n-1)
[1] 1.396815 1.859548 2.306004 2.896459 3.355387

> qt(c(0.9, 0.95, 0.975, 0.99, 0.995), df=n)
[1] 1.383029 1.833113 2.262157 2.821438 3.249836
```

4. Suppose it is desired to test the null hypothesis  $H_0: \mu = 15$ . Compute the value taken by the appropriate  $t$ -statistic to 3 decimal places.

*Ans*

−0.86

---

*Working for Question 4.*

The value of the statistic is

$$t = \frac{13.37 - 15}{5.686/\sqrt{9}} \approx -0.860007034822371$$

5. Compute a 95% upper confidence limit for  $\mu$  (to 2 decimal places).

*Ans*

16.89

---

*Working for Question 5.*

The upper confidence limit is  $13.37 + (c \times 5.686)/\sqrt{9}$  where  $P(t_8 > c) = 0.05$ . From the R output above (using 8 degrees of freedom) we see that  $c = 1.8595$ . So we get

$$13.37 + (1.8595 \times 5.686)/\sqrt{9} \approx 16.89437.$$

**The next three questions relate to the following scenario:** Two random samples are drawn independently from two normal populations with unknown means and variances. The first sample (sample “A”) is of size 5, has mean 7.811 and standard deviation 1.95. The second sample (sample “B”) is of size 11, has mean 19.591 and standard deviation 7.81.

6. If it is assumed that the two populations have the same variance  $\sigma^2$ , compute and write down the pooled estimate of  $\sigma^2$  (to 3 decimal places).

Ans

44.655

---

*Working for Question 6.*

The pooled estimate of the variance is given by

$$s_p^2 = \frac{(5-1)1.95^2 + (11-1)7.81^2}{5+11-2} \approx 44.6550714285714$$

7. Still assuming equal population variances, the appropriate  $t$ -statistic for testing the hypothesis that the population means are equal is of the form  $b/s_p$  where  $s_p$  is the square root of your answer to the previous question. Write down the value of  $b$  to 3 decimal places in the box below.

*Ans*

−21.841

---

*Working for Question 7.*

The statistic is

$$t = \frac{7.811 - 19.591}{s_p \sqrt{\frac{1}{5} + \frac{1}{11}}} \quad (\text{or minus this})$$

so the value of  $b$  is given by

$$\frac{7.811 - 19.591}{\sqrt{\frac{1}{5} + \frac{1}{11}}} \approx -21.8407045444967 \quad (\text{or minus this}).$$



8. Suppose now that it is *not* assumed that the population variances are equal (so a Welch test is to be performed). Write down to 3 decimal places the standard error of the mean difference.

|                                       |
|---------------------------------------|
| <p><i>Ans</i></p> <p><i>2.511</i></p> |
|---------------------------------------|

---

*Working for Question 8.*

The standard error is

$$\sqrt{\frac{1.95^2}{5} + \frac{7.81^2}{11}} \approx 2.51109537851512$$

**The final two questions relate to the following scenario.** Bivariate data  $(x_1, y_1), \dots, (x_{14}, y_{14})$  is collected and each  $y_i$  is modelled as the value taken by a normal random variable  $Y_i \sim N(\alpha + \beta x_i, \sigma^2)$  for unknown parameters  $\alpha, \beta$  and  $\sigma > 0$ . All random variables are assumed independent.

After computing the least-squares regression line the residual sum of squares is 1.869 and the sample standard deviation of the  $x_i$ 's is 0.434. The scatterplot appears roughly linear and there is no obvious pattern in the residual plot.

The following R output may be useful:

```
> n
[1] 14

> qt(c(0.9, 0.95, 0.98, 0.975, 0.99, 0.995), df=n)
[1] 1.345030 1.761310 2.263781 2.144787 2.624494 2.976843

> qt(c(0.9, 0.95, 0.98, 0.975, 0.99, 0.995), df=n-1)
[1] 1.350171 1.770933 2.281604 2.160369 2.650309 3.012276

> qt(c(0.9, 0.95, 0.98, 0.975, 0.99, 0.995), df=n-2)
[1] 1.356217 1.782288 2.302722 2.178813 2.680998 3.054540
```

9. Write down to 3 decimal places the value of  $\hat{\sigma}^2$ , the estimate of the error variance.

*Ans*

*0.156*

---

*Working for Question 9.*

The estimate of  $\sigma^2$  is given by

$$\frac{1.869}{14 - 2} \approx 0.15575.$$

10. A 99% lower confidence limit for  $\beta$  is of the form  $\hat{\beta} - c\hat{\sigma}$  where  $\hat{\sigma}$  is the square root of your answer to the previous question. Write the value of  $c$  to 3 decimal places in the box below.

*Ans*

*1.713*

---

*Working for Question 10.*

The lower confidence limit is given by  $\hat{\beta} - t \frac{\hat{\sigma}}{\sqrt{S_{xx}}}$  where

$$P(t_{14-2} > t) = 0.01$$

and  $S_{xx} = (14 - 1)0.434^2$ . From the R output above we see that the value of  $t$  is 2.681. Thus the value of  $c$  is given by

$$c = \frac{t}{\sqrt{S_{xx}}} \approx \frac{2.681}{0.434\sqrt{13}} \approx 1.71330786414728.$$

# FORMULA SHEET FOR MATH1905 STATISTICS

## • Calculation formulae:

– For a sample  $x_1, x_2, \dots, x_n$

|                       |  |
|-----------------------|--|
| Sample mean $\bar{x}$ | $\frac{1}{n} \sum_{i=1}^n x_i$   |
| Sample variance $s^2$ | $\frac{1}{n-1} \left[ \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \right] = \frac{1}{n-1} S_{xx}$ |

– For paired observations  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

|          |  |
|----------|--|
| $S_{xy}$ | $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left( \sum_{i=1}^n x_i \right) \left( \sum_{i=1}^n y_i \right)$ |
| $S_{xx}$ | $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2$  |
| $S_{yy}$ | $\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - \frac{1}{n} \left( \sum_{i=1}^n y_i \right)^2$  |
| $r$      | $\frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$  |

|   |                         |
|---|-------------------------|
| For the regression line<br>$y = a + bx$ : |                         |
| $b$                                       | $\frac{S_{xy}}{S_{xx}}$ |
| $a$                                       | $\bar{y} - b\bar{x}$    |

## • Some probability results:

|                                       |   |
|---------------------------------------|---|
| For any two events $A$ and $B$        | $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ and<br>$P(A \cap B) = P(A)P(B A)$ |
| If $A$ and $B$ are mutually exclusive | $P(A \cap B) = 0$ and $P(A \cup B) = P(A) + P(B)$                           |
| If $A$ and $B$ are independent        | $P(A \cap B) = P(A)P(B)$  |

- If  $Y \sim \text{Pois}(\lambda)$ , then  $P(Y = k) = \frac{e^{-\lambda} \lambda^k}{k!}$  for  $k = 0, 1, 2, \dots$ . Furthermore,  $\mathbb{E}(Y) = \lambda$  and  $\text{var}(Y) = \lambda$ .
- If  $X \sim B(n, p)$ , then,  $P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}$ , for  $x = 0, 1, \dots, n$ . Furthermore,  $E(X) = np$  and  $\text{var}(X) = np(1-p)$ .

## • Some test statistics and sampling distributions under appropriate assumptions and hypotheses:

|  |
|--|
| $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ |
| $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$ |
| $\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$      |

|   |
|---|
| $\frac{\bar{X} - \bar{Y}}{S_p \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}} \sim t_{n_x + n_y - 2}$ , where<br>$S_p^2 = \frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2}$  |
| $\hat{\alpha} \sim N\left(\alpha, \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right] \right); \hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{S_{xx}}\right); \hat{\sigma}^2 = \frac{\sum_i \hat{\epsilon}_i^2}{n-2} \sim \frac{\sigma^2 \chi_{n-2}^2}{n-2}$ |
| $\sum_i \frac{(O_i - E_i)^2}{E_i} = \sum_i \frac{O_i^2}{E_i} - n \sim \chi_\nu^2$ , for appropriate $\nu$   |