

1.

Data Analysis

Mr Faruk

Teacher of Mathematics
BSc/MSc/PGCE Mathematics

✉ cieigcsesolutions@gmail.com



SECTION A: STATISTICS

Answer ALL questions. Write your answers in the spaces provided.

1. Sara is investigating the variation in daily maximum gust, t kn, for Camborne in June and July 1987.

She used the large data set to select a sample of size 20 from the June and July data for 1987. Sara selected the first value using a random number from 1 to 4 and then selected every third value after that.

(a) State the sampling technique Sara used. (1)

(b) From your knowledge of the large data set, explain why this process may not generate a sample of size 20. (1)

The data Sara collected are summarised as follows

$$n = 20 \quad \sum t = 374 \quad \sum t^2 = 7600$$

(c) Calculate the standard deviation. (2)

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

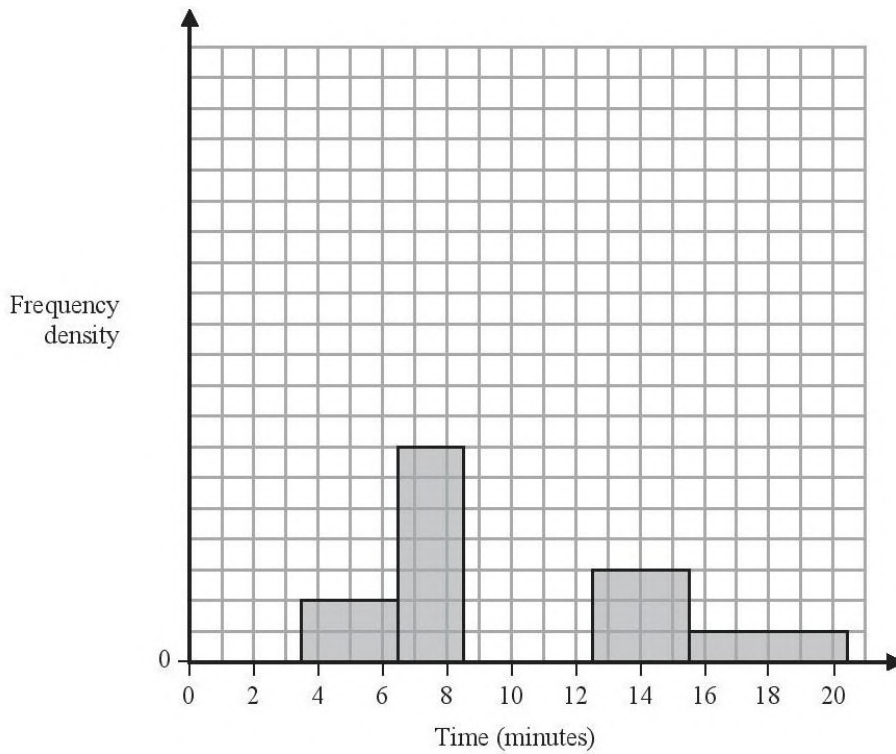
.....

.....

.....

.....

2. The partially completed histogram and the partially completed table show the time, to the nearest minute, that a random sample of motorists were delayed by roadworks on a stretch of motorway.



| Delay (minutes) | Number of motorists |
|-----------------|---------------------|
| 4 – 6 | 6 |
| 7 – 8 | |
| 9 | 17 |
| 10 – 12 | 45 |
| 13 – 15 | 9 |
| 16 – 20 | |

Estimate the percentage of these motorists who were delayed by the roadworks for between 8.5 and 13.5 minutes.

(5)

.....

.....

.....

SECTION A: STATISTICS

Answer ALL questions. Write your answers in the spaces provided.

1. The number of hours of sunshine each day, y , for the month of July at Heathrow are summarised in the table below.

| | | | | | |
|------------------|----------------|----------------|-----------------|------------------|------------------|
| Hours | $0 \leq y < 5$ | $5 \leq y < 8$ | $8 \leq y < 11$ | $11 \leq y < 12$ | $12 \leq y < 14$ |
| Frequency | 12 | 6 | 8 | 3 | 2 |

A histogram was drawn to represent these data. The $8 \leq y < 11$ group was represented by a bar of width 1.5 cm and height 8 cm.

- (a) Find the width and the height of the $0 \leq y < 5$ group. (3)

- (b) Use your calculator to estimate the mean and the standard deviation of the number of hours of sunshine each day, for the month of July at Heathrow.
Give your answers to 3 significant figures. (3)

(1)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

Question 1 continued

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

(Total for Question 1 is 13 marks)

SECTION A: STATISTICS

Answer ALL questions.

1. *Kaff* coffee is sold in packets. A seller measures the masses of the contents of a random sample of 90 packets of *Kaff* coffee from her stock. The results are shown in the table below.

| Mass w (g) | Midpoint y (g) | Frequency f |
|--------------------|------------------|---------------|
| $240 \leq w < 245$ | 242.5 | 8 |
| $245 \leq w < 248$ | 246.5 | 15 |
| $248 \leq w < 252$ | 250.0 | 35 |
| $252 \leq w < 255$ | 253.5 | 23 |
| $255 \leq w < 260$ | 257.5 | 9 |

(You may use $\sum fy^2 = 5\,644\,171.75$)

A histogram is drawn and the class $245 \leq w < 248$ is represented by a rectangle of width 1.2 cm and height 10 cm.

- (a) Calculate the width and the height of the rectangle representing the class $255 \leq w < 260$.

(3)

4. Charlie is studying the time it takes members of his company to travel to the office. He stands by the door to the office from 0840 to 0850 one morning and asks workers, as they arrive, how long their journey was.

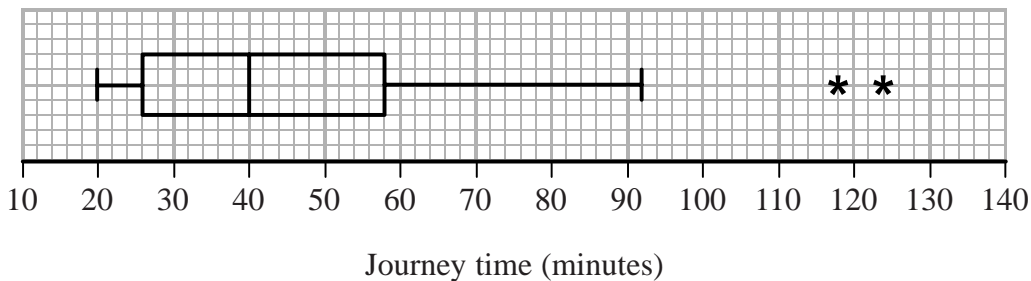
(a) State the sampling method Charlie used. (1)

(b) State and briefly describe an alternative method of non-random sampling Charlie could have used to obtain a sample of 40 workers. (2)

Taruni decided to ask every member of the company the time, x minutes, it takes them to travel to the office.

(c) State the data selection process Taruni used. (1)

Taruni's results are summarised by the box plot and summary statistics below.



$$n = 95 \quad \sum x = 4133 \quad \sum x^2 = 202294$$

(d) Write down the interquartile range for these data. (1)

(e) Calculate the mean and the standard deviation for these data. (3)

(f) State, giving a reason, whether you would recommend using the mean and standard deviation or the median and interquartile range to describe these data. (2)

Rana and David both work for the company and have both moved house since Taruni collected her data.

Rana's journey to work has changed from 75 minutes to 35 minutes and David's journey to work has changed from 60 minutes to 33 minutes.

Taruni drew her box plot again and only had to change two values.

(g) Explain which two values Taruni must have changed and whether each of these values has increased or decreased. (3)

4. Joshua is investigating the daily total rainfall in Hurn for May to October 2015

Using the information from the large data set, Joshua wishes to calculate the mean of the daily total rainfall in Hurn for May to October 2015

- (a) Using your knowledge of the large data set, explain why Joshua needs to clean the data before calculating the mean. (1)

Using the information from the large data set, he produces the grouped frequency table below.

| Daily total rainfall (r mm) | Frequency | Midpoint (x mm) |
|--------------------------------|-----------|--------------------|
| $0 \leq r < 0.5$ | 121 | 0.25 |
| $0.5 \leq r < 1.0$ | 10 | 0.75 |
| $1.0 \leq r < 5.0$ | 24 | 3.0 |
| $5.0 \leq r < 10.0$ | 12 | 7.5 |
| $10.0 \leq r < 30.0$ | 17 | 20.0 |

You may use $\sum fx = 539.75$ and $\sum fx^2 = 7704.1875$

- (b) Use linear interpolation to calculate an estimate for the upper quartile of the daily total rainfall. (2)
- (c) Calculate an estimate for the standard deviation of the daily total rainfall in Hurn for May to October 2015 (2)
- (d) (i) State the assumption involved with using class midpoints to calculate an estimate of a mean from a grouped frequency table.
- (ii) Using your knowledge of the large data set, explain why this assumption does not hold in this case.
- (iii) State, giving a reason, whether you would expect the actual mean daily total rainfall in Hurn for May to October 2015 to be larger than, smaller than or the same as an estimate based on the grouped frequency table. (3)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

2.

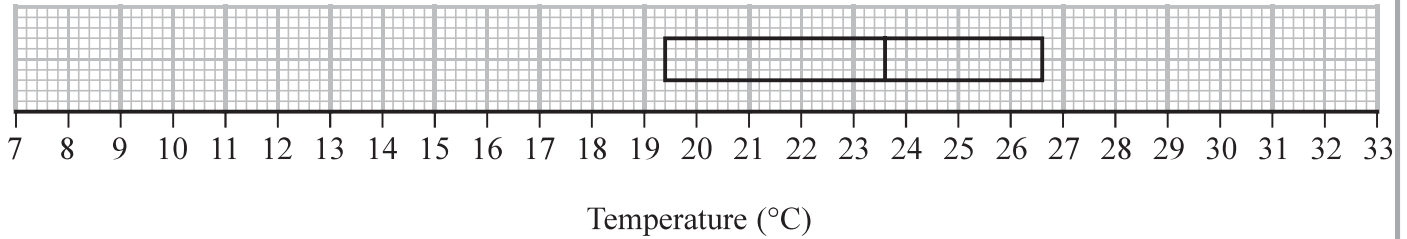


Figure 1

The partially completed box plot in Figure 1 shows the distribution of daily mean air temperatures using the data from the large data set for Beijing in 2015

An outlier is defined as a value
 more than $1.5 \times \text{IQR}$ below Q_1 or
 more than $1.5 \times \text{IQR}$ above Q_3

The three lowest air temperatures in the data set are 7.6°C , 8.1°C and 9.1°C
 The highest air temperature in the data set is 32.5°C

(a) Complete the box plot in Figure 1 showing clearly any outliers. (4)

(b) Using your knowledge of the large data set, suggest from which month the two outliers are likely to have come. (1)

Using the data from the large data set, Simon produced the following summary statistics for the daily mean air temperature, $x^\circ\text{C}$, for Beijing in 2015

$$n = 184 \quad \sum x = 4153.6 \quad S_{xx} = 4952.906$$

(c) Show that, to 3 significant figures, the standard deviation is 5.19°C (1)

Simon decides to model the air temperatures with the random variable

$$T \sim N(22.6, 5.19^2)$$

(d) Using Simon's model, calculate the 10th to 90th interpercentile range. (3)

Simon wants to model another variable from the large data set for Beijing using a normal distribution.

(e) State two variables from the large data set for Beijing that are **not** suitable to be modelled by a normal distribution. Give a reason for each answer. (2)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

4. A lake contains three different types of carp.

There are an estimated 450 mirror carp, 300 leather carp and 850 common carp.

Tim wishes to investigate the health of the fish in the lake.

He decides to take a sample of 160 fish.

- (a) Give a reason why stratified random sampling cannot be used. (1)
- (b) Explain how a sample of size 160 could be taken to ensure that the estimated populations of each type of carp are fairly represented.

You should state the name of the sampling method used. (2)

As part of the health check, Tim weighed the fish.

His results are given in the table below.

| Weight (w kg) | Frequency (f) | Midpoint (m kg) |
|------------------|-------------------|--------------------|
| $2 \leq w < 3.5$ | 8 | 2.75 |
| $3.5 \leq w < 4$ | 32 | 3.75 |
| $4 \leq w < 4.5$ | 64 | 4.25 |
| $4.5 \leq w < 5$ | 40 | 4.75 |
| $5 \leq w < 6$ | 16 | 5.5 |

(You may use $\sum fm = 692$ and $\sum fm^2 = 3053$)

- (c) Calculate an estimate for the standard deviation of the weight of the carp. (2)

Tim realised that he had transposed the figures for 2 of the weights of the fish.

He had recorded in the table 2.3 instead of 3.2 and 4.6 instead of 6.4

- (d) Without calculating a new estimate for the standard deviation, state what effect
- (i) using the correct figure of 3.2 instead of 2.3
- (ii) using the correct figure of 6.4 instead of 4.6

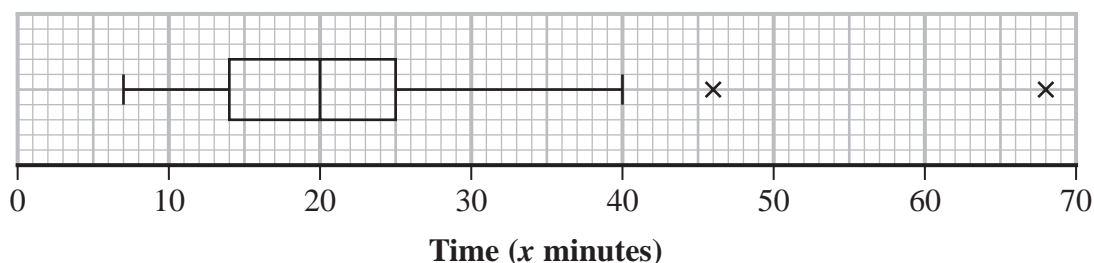
would have on your estimated standard deviation.

Give a reason for each of your answers. (2)

3. Each member of a group of 27 people was timed when completing a puzzle.

The time taken, x minutes, for each member of the group was recorded.

These times are summarised in the following box and whisker plot.



- (a) Find the range of the times. (1)

- (b) Find the interquartile range of the times. (1)

For these 27 people $\sum x = 607.5$ and $\sum x^2 = 17\,623.25$

- (c) calculate the mean time taken to complete the puzzle, (1)

- (d) calculate the standard deviation of the times taken to complete the puzzle. (2)

Taruni defines an outlier as a value more than 3 standard deviations above the mean.

- (e) State how many outliers Taruni would say there are in these data, giving a reason for your answer. (1)

Adam and Beth also completed the puzzle in a minutes and b minutes respectively, where $a > b$.

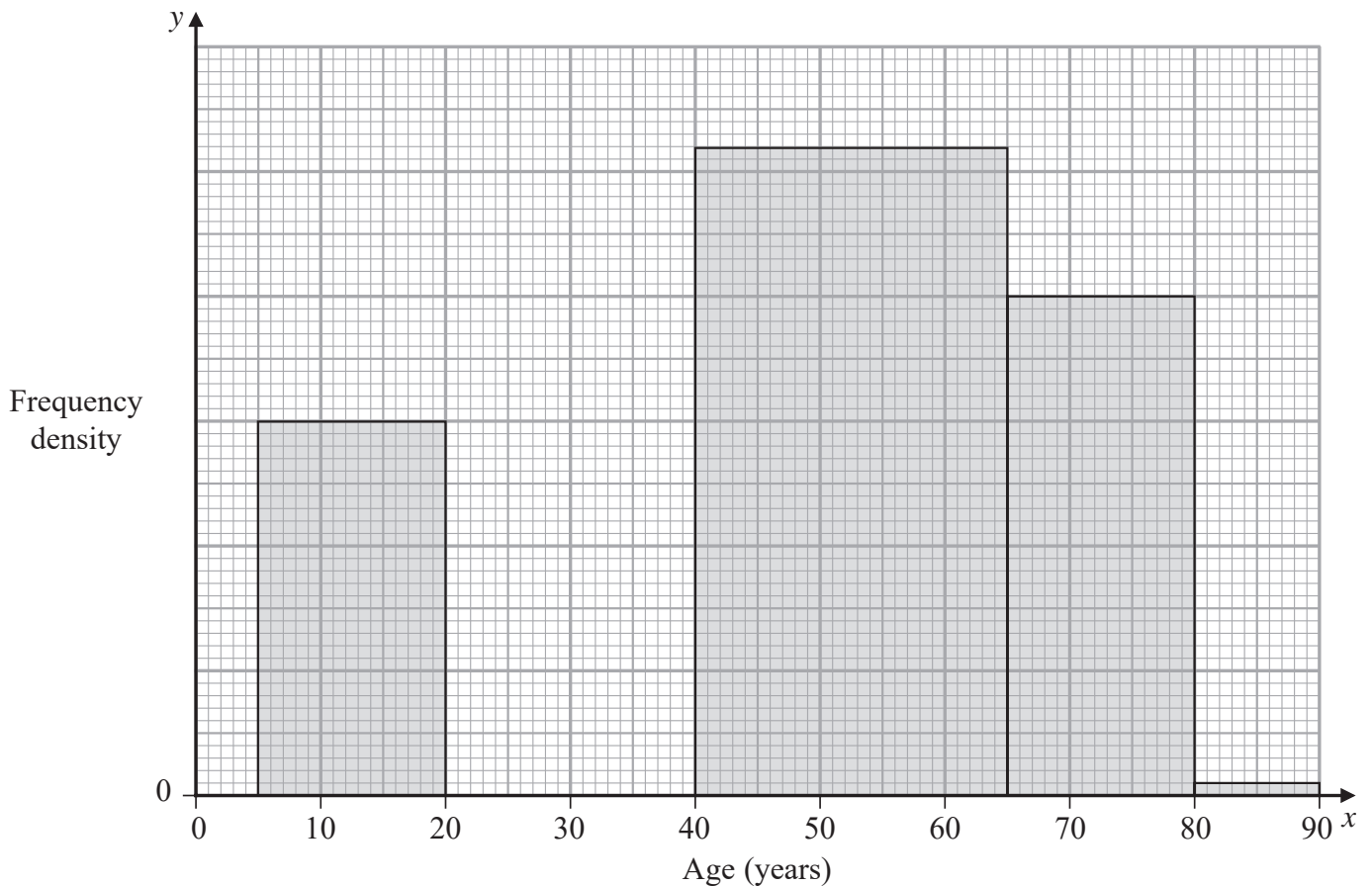
When their times are included with the data of the other 27 people

- the median time increases
 - the mean time does not change
- (f) Suggest a possible value for a and a possible value for b , explaining how your values satisfy the above conditions. (3)
- (g) Without carrying out any further calculations, explain why the standard deviation of all 29 times will be lower than your answer to part (d). (1)

2. The partially completed table and partially completed histogram give information about the ages of passengers on an airline.

There were no passengers aged 90 or over.

| Age (x years) | $0 \leq x < 5$ | $5 \leq x < 20$ | $20 \leq x < 40$ | $40 \leq x < 65$ | $65 \leq x < 80$ | $80 \leq x < 90$ |
|------------------|----------------|-----------------|------------------|------------------|------------------|------------------|
| Frequency | 5 | 45 | 90 | | | 1 |



- (a) Complete the histogram. (3)
- (b) Use linear interpolation to estimate the median age. (4)

An outlier is defined as a value greater than $Q_3 + 1.5 \times$ interquartile range.

Given that $Q_1 = 27.3$ and $Q_3 = 58.9$

- (c) determine, giving a reason, whether or not the oldest passenger could be considered as an outlier. (2)

3. Helen is studying one of the qualitative variables from the large data set for Heathrow from 2015.

She started with the data from 3rd May and then took every 10th reading.

There were only 3 different outcomes with the following frequencies

| | | | |
|------------------|----------|----------|----------|
| Outcome | <i>A</i> | <i>B</i> | <i>C</i> |
| Frequency | 16 | 2 | 1 |

- (a) State the sampling technique Helen used. (1)

- (b) From your knowledge of the large data set
- (i) suggest which variable was being studied,
- (ii) state the name of outcome *A*. (2)

George is also studying the same variable from the large data set for Heathrow from 2015. He started with the data from 5th May and then took every 10th reading and obtained the following

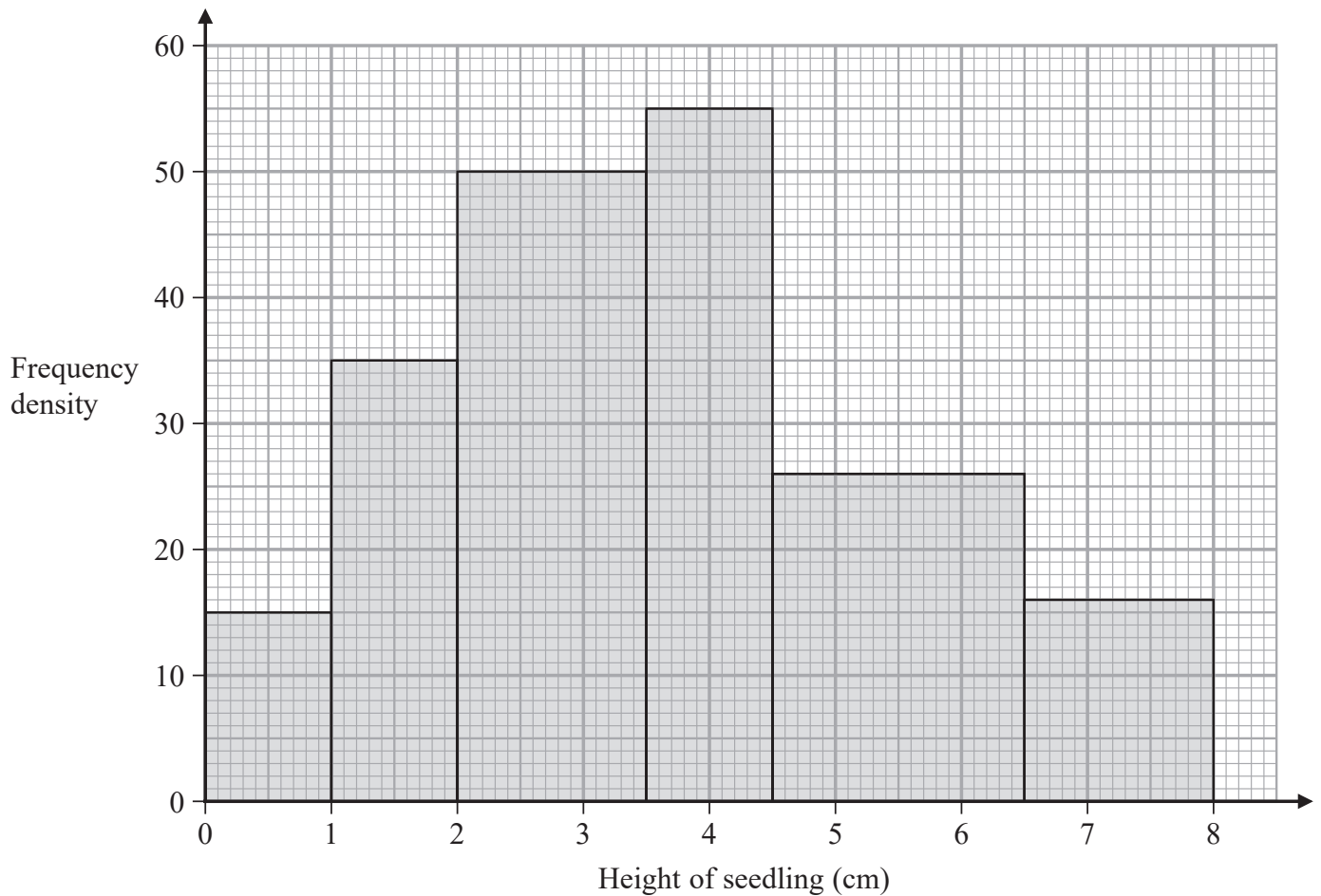
| | | | |
|------------------|----------|----------|----------|
| Outcome | <i>A</i> | <i>B</i> | <i>C</i> |
| Frequency | 16 | 1 | 1 |

Helen and George decided they should examine all of the data for this variable for Heathrow from 2015 and obtained the following

| | | | |
|------------------|----------|----------|----------|
| Outcome | <i>A</i> | <i>B</i> | <i>C</i> |
| Frequency | 155 | 26 | 3 |

- (c) State what inference Helen and George could reliably make from their original samples about the outcomes of this variable at Heathrow, for the period covered by the large data set in 2015. (1)

3. The histogram summarises the heights of 256 seedlings two weeks after they were planted.



- (a) Use linear interpolation to estimate the median height of the seedlings.

(4)

Chris decides to model the **frequency density** for these 256 seedlings by a curve with equation

$$y = kx(8 - x) \quad 0 \leq x \leq 8$$

where k is a constant.

- (b) Find the value of k

(3)

Using this model,

- (c) write down the median height of the seedlings.

(1)

Question 1 continued

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

(Total for Question 1 is 5 marks)

Question 2 continued

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

(Total for Question 2 is 3 marks)

3. Customers in a shop have to queue to pay.

The partially completed table below and partially completed histogram opposite, give information about the time, x minutes, spent in the queue by each of 112 customers one day.

| Time in queue (x minutes) | Frequency |
|------------------------------|-----------|
| 1–2 | 64 |
| 2–3 | |
| 3–4 | 13 |
| 4–6 | |
| 6–8 | 3 |

No customer spent less than 1 minute or longer than 8 minutes in the queue.

(a) Complete the table.

(2)

(b) Complete the histogram.

(2)

Ting decides to model the **frequency density** for these 112 customers by a curve with equation

$$y = \frac{k}{x^2} \quad 1 \leq x \leq 8$$

where k is a constant.

(c) Find the value of k

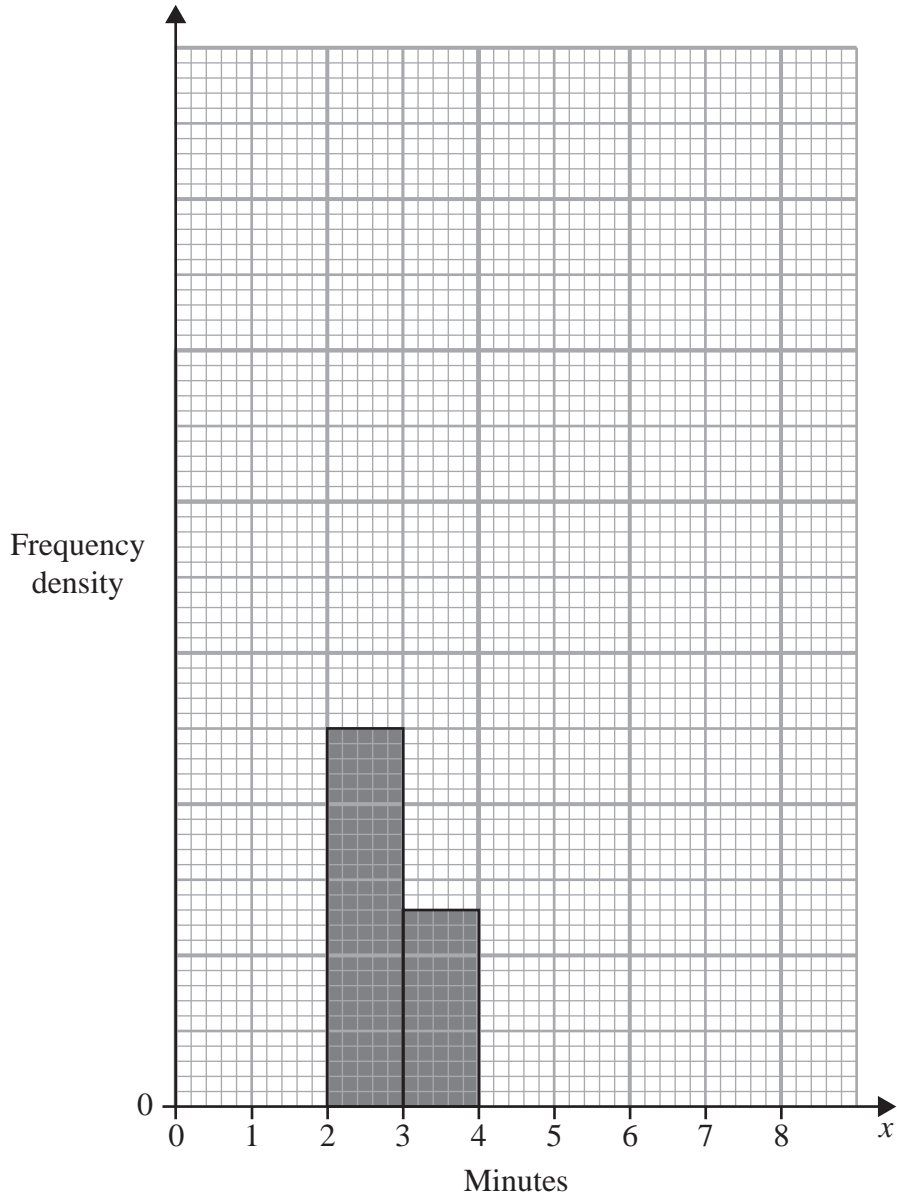
(3)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

Question 3 continued



Turn over for a spare grid if you need to redraw your histogram.

