

Edexcel Maths S1

Topic Questions from Papers

Representation & Summary of Data

2. The following table summarises the distances, to the nearest km, that 134 examiners travelled to attend a meeting in London.

Distance (km)	Number of examiners
41–45	4
46–50	19
51–60	53
61–70	37
71–90	15
91–150	6

- (a) Give a reason to justify the use of a histogram to represent these data. (1)
- (b) Calculate the frequency densities needed to draw a histogram for these data.  
**(DO NOT DRAW THE HISTOGRAM)** (2)
- (c) Use interpolation to estimate the median  $Q_2$ , the lower quartile  $Q_1$ , and the upper quartile  $Q_3$  of these data. (4)

The mid-point of each class is represented by  $x$  and the corresponding frequency by  $f$ . Calculations then give the following values

$$\Sigma fx = 8379.5 \quad \text{and} \quad \Sigma fx^2 = 557489.75$$

- (d) Calculate an estimate of the mean and an estimate of the standard deviation for these data. (4)

One coefficient of skewness is given by

$$\frac{Q_3 - 2Q_2 + Q_1}{Q_3 - Q_1}.$$

- (e) Evaluate this coefficient and comment on the skewness of these data. (4)
- (f) Give another justification of your comment in part (e). (1)

---

---

---

---

---



Leave  
blank

### Question 2 continued



(a) On the graph paper opposite draw a box plot to represent these data.

**(7)**

(b) Comment on the distribution of delays. Justify your answer.

(2)

(c) Suggest how the distribution might be interpreted by a passenger who frequently flies from City  $A$  to City  $B$ .

(1)



Leave  
blank

**Question 4 continued**

[illegible]

**(Total 10 marks)**

**Q4**



1. (a) Describe the main features and uses of a box plot.

---

---

---

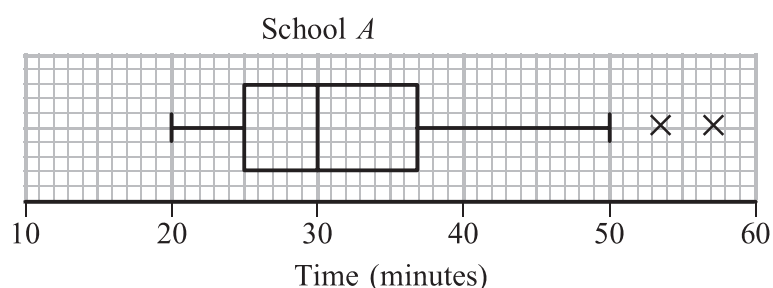
---

---

(3)

Children from schools *A* and *B* took part in a fun run for charity. The times, to the nearest minute, taken by the children from school *A* are summarised in Figure 1.

**Figure 1**



- (b) (i) Write down the time by which 75% of the children in school *A* had completed the run.

---

---

- (ii) State the name given to this value.

---

(2)

- (c) Explain what you understand by the two crosses (X) on Figure 1.

---

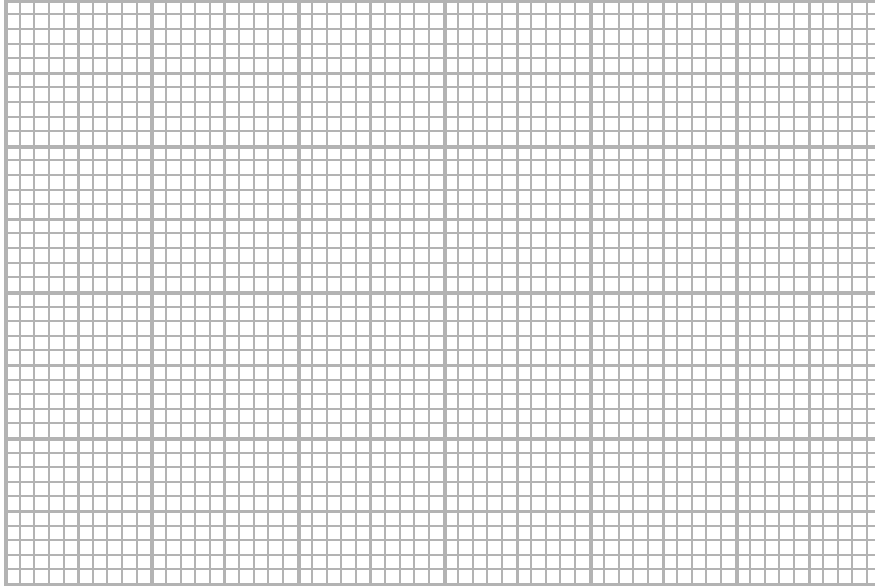
---

(2)

**Question 1 continued**

For school *B* the least time taken by any of the children was 25 minutes and the longest time was 55 minutes. The three quartiles were 30, 37 and 50 respectively.

- (d) Draw a box plot to represent the data from school *B*.



(4)

- (e) Compare and contrast these two box plots.

---

---

---

---

---

(4)

(Total 15 marks)

Q1

--	--



2. Sunita and Shelley talk to one another once a week on the telephone. Over many weeks they recorded, to the nearest minute, the number of minutes spent in conversation on each occasion. The following table summarises their results.

Time (to the nearest minute)	Number of Conversations
5–9	2
10–14	9
15–19	20
20–24	13
25–29	8
30–34	3

Two of the conversations were chosen at random.

- (a) Find the probability that both of them were longer than 24.5 minutes. (2)

The mid-point of each class was represented by  $x$  and its corresponding frequency by  $f$ , giving  $\Sigma fx = 1060$ .

- (b) Calculate an estimate of the mean time spent on their conversations. (2)

During the following 25 weeks they monitored their weekly conversations and found that at the end of the 80 weeks their overall mean length of conversation was 21 minutes.

- (c) Find the mean time spent in conversation during these 25 weeks. (4)

- (d) Comment on these two mean values. (2)

---

---

---

---

---

---

---

---

---

---





Leave  
blank

### Question 2 continued



4. Summarised below are the distances, to the nearest mile, travelled to work by a random sample of 120 commuters.

Distance (to the nearest mile)	Number of commuters
0–9	10
10–19	19
20–29	43
30–39	25
40–49	8
50–59	6
60–69	5
70–79	3
80–89	1

For this distribution,

- (a) describe its shape, (1)

- (b) use linear interpolation to estimate its median. (2)

The mid-point of each class was represented by  $x$  and its corresponding frequency by  $f$  giving

$$\Sigma fx = 3550 \quad \text{and} \quad \Sigma fx^2 = 138020$$

- (c) Estimate the mean and the standard deviation of this distribution. (3)

One coefficient of skewness is given by

$$\frac{3(\text{mean} - \text{median})}{\text{standard deviation}}.$$

- (d) Evaluate this coefficient for this distribution. (3)

- (e) State whether or not the value of your coefficient is consistent with your description in part (a). Justify your answer. (2)



Leave  
blank

- (f) State, with a reason, whether you should use the mean or the median to represent the data in this distribution.

(2)

- (g) State the circumstance under which it would not matter whether you used the mean or the median to represent a set of data.

(1)



Leave  
blank

Question 4 continued

Lined area for writing the answer to Question 4.



5. A teacher recorded, to the nearest hour, the time spent watching television during a particular week by each child in a random sample. The times were summarised in a grouped frequency table and represented by a histogram.

One of the classes in the grouped frequency distribution was 20–29 and its associated frequency was 9. On the histogram the height of the rectangle representing that class was 3.6 cm and the width was 2 cm.

- (a) Give a reason to support the use of a histogram to represent these data. (1)

- (b) Write down the underlying feature associated with each of the bars in a histogram. (1)

- (c) Show that on this histogram each child was represented by  $0.8 \text{ cm}^2$ . (3)

The total area under the histogram was  $24 \text{ cm}^2$ .

- (d) Find the total number of children in the group. (2)



Leave  
blank

Question 5 continued

Lined area for writing the answer to Question 5.

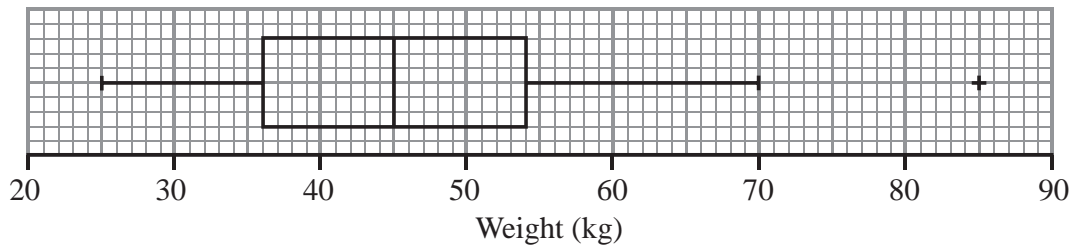
(Total 7 marks)

Q5

Mark box for Q5.



2. The box plot in Figure 1 shows a summary of the weights of the luggage, in kg, for each musician in an orchestra on an overseas tour.



**Figure 1**

The airline's recommended weight limit for each musician's luggage was 45 kg. Given that none of the musicians' luggage weighed exactly 45 kg,

- (a) state the proportion of the musicians whose luggage was below the recommended weight limit.

(1)

A quarter of the musicians had to pay a charge for taking heavy luggage.

- (b) State the smallest weight for which the charge was made.

(1)

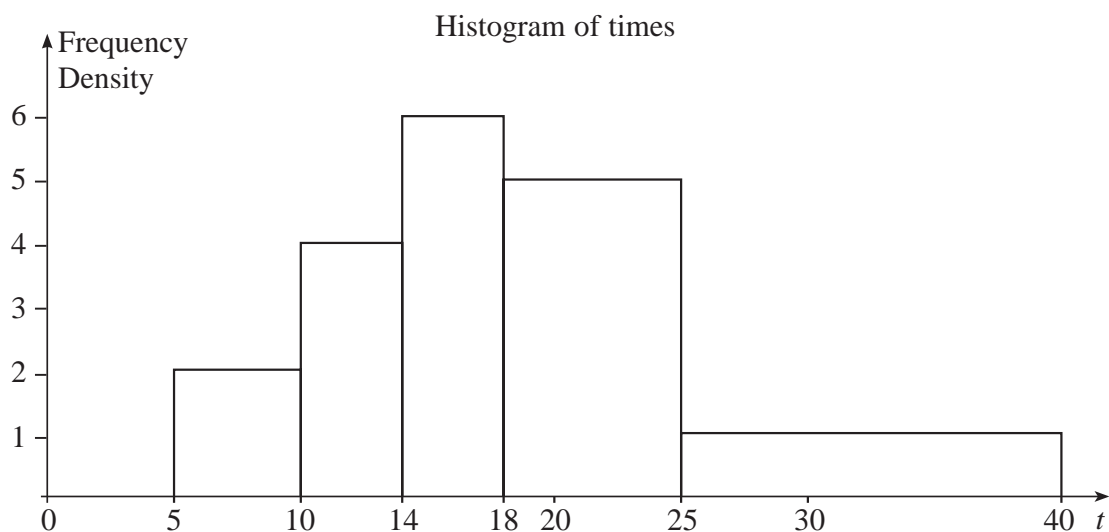
- (c) Explain what you understand by the + on the box plot in Figure 1, and suggest an instrument that the owner of this luggage might play.

(2)





5.



**Figure 2**

Figure 2 shows a histogram for the variable  $t$  which represents the time taken, in minutes, by a group of people to swim 500m.

(a) Complete the frequency table for  $t$ .

$t$	5–10	10–14	14–18	18–25	25–40
Frequency	10	16	24		

(2)

(b) Estimate the number of people who took longer than 20 minutes to swim 500m.

(2)

(c) Find an estimate of the mean time taken.

(4)

(d) Find an estimate for the standard deviation of  $t$ .

(3)

(e) Find the median and quartiles for  $t$ .

(4)

One measure of skewness is found using  $\frac{3(\text{mean} - \text{median})}{\text{standard deviation}}$ .

(f) Evaluate this measure and describe the skewness of these data.

(2)

---



---



---



2. Cotinine is a chemical that is made by the body from nicotine which is found in cigarette smoke. A doctor tested the blood of 12 patients, who claimed to smoke a packet of cigarettes a day, for cotinine. The results, in appropriate units, are shown below.

Patient	$A$	$B$	$C$	$D$	$E$	$F$	$G$	$H$	$I$	$J$	$K$	$L$
Cotinine level, $x$	160	390	169	175	125	420	171	250	210	258	186	243

[You may use  $\sum x^2 = 724\,961$ ]

- (a) Find the mean and standard deviation of the level of cotinine in a patient's blood. (4)

- (b) Find the median, upper and lower quartiles of these data. (3)

A doctor suspects that some of his patients have been smoking more than a packet of cigarettes per day. He decides to use  $Q_3 + 1.5(Q_3 - Q_1)$  to determine if any of the cotinine results are far enough away from the upper quartile to be outliers.

- (c) Identify which patient(s) may have been smoking more than a packet of cigarettes a day. Show your working clearly. (4)

Research suggests that cotinine levels in the blood form a skewed distribution.

One measure of skewness is found using  $\frac{(Q_1 - 2Q_2 + Q_3)}{(Q_3 - Q_1)}$ .

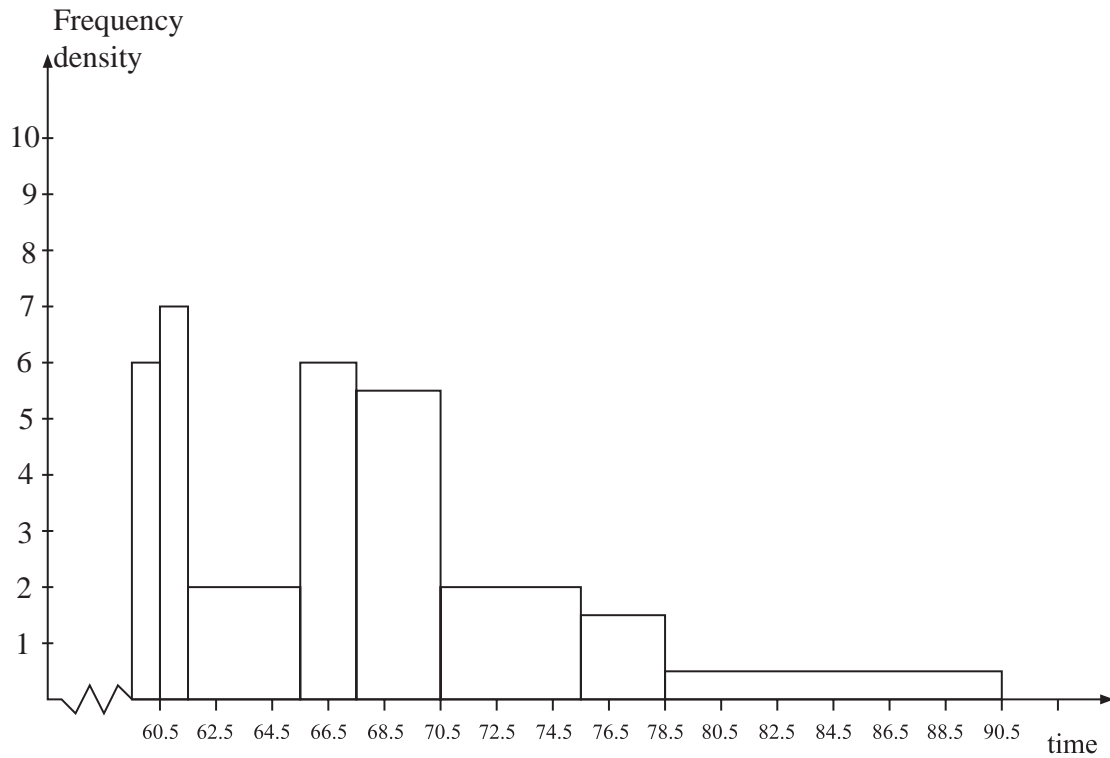
- (d) Evaluate this measure and describe the skewness of these data. (3)



11/11/2019



3. The histogram in Figure 1 shows the time taken, to the nearest minute, for 140 runners to complete a fun run.



**Figure 1**

Use the histogram to calculate the number of runners who took between 78.5 and 90.5 minutes to complete the fun run.

(5)

2. The age in years of the residents of two hotels are shown in the back to back stem and leaf diagram below.

Abbey Hotel    8|5|0 means 58 years in Abbey hotel and 50 years in Balmoral hotel    Balmoral Hotel

(1)	2	0		
(4)	9751	1		
(4)	9831	2	6	(1)
(11)	99997665332	3	447	(3)
(6)	987750	4	005569	(6)
(1)	8	5	000013667	(9)
		6	233457	(6)
		7	015	(3)

For the Balmoral Hotel,

- (a) write down the mode of the age of the residents, (1)
- (b) find the values of the lower quartile, the median and the upper quartile. (3)
- (c) (i) Find the mean,  $\bar{x}$ , of the age of the residents.
- (ii) Given that  $\sum x^2 = 81\,213$  find the standard deviation of the age of the residents. (4)

One measure of skewness is found using

$$\frac{\text{mean} - \text{mode}}{\text{standard deviation}}$$

- (d) Evaluate this measure for the Balmoral Hotel. (2)

For the Abbey Hotel, the mode is 39, the mean is 33.2, the standard deviation is 12.7 and the measure of skewness is  $-0.454$

- (e) Compare the two age distributions of the residents of each hotel. (3)

---

---

---

---

---

---

---

---



Leave  
blank

### Question 2 continued



5. A person's blood group is determined by whether or not it contains any of 3 substances  $A$ ,  $B$  and  $C$ .

A doctor surveyed 300 patients' blood and produced the table below.

Blood contains	No. of Patients
only $C$	100
$A$ and $C$ but not $B$	100
only $A$	30
$B$ and $C$ but not $A$	25
only $B$	12
$A$ , $B$ and $C$	10
$A$ and $B$ but not $C$	3

- (a) Draw a Venn diagram to represent this information.

(4)







4. In a study of how students use their mobile telephones, the phone usage of a random sample of 11 students was examined for a particular week.

The total length of calls,  $y$  minutes, for the 11 students were

17, 23, 35, 36, 51, 53, 54, 55, 60, 77, 110

- (a) Find the median and quartiles for these data.

(3)

A value that is greater than  $Q_3 + 1.5 \times (Q_3 - Q_1)$  or smaller than  $Q_1 - 1.5 \times (Q_3 - Q_1)$  is defined as an outlier.

- (b) Show that 110 is the only outlier.

(2)

- (c) Using the graph paper on page 15 draw a box plot for these data indicating clearly the position of the outlier.

(3)

The value of 110 is omitted.

- (d) Show that  $S_{yy}$  for the remaining 10 students is 2966.9

(3)

These 10 students were each asked how many text messages,  $x$ , they sent in the same week.

The values of  $S_{xx}$  and  $S_{xy}$  for these 10 students are  $S_{xx} = 3463.6$  and  $S_{xy} = -18.3$ .

- (e) Calculate the product moment correlation coefficient between the number of text messages sent and the total length of calls for these 10 students.

(2)

A parent believes that a student who sends a large number of text messages will spend fewer minutes on calls.

- (f) Comment on this belief in the light of your calculation in part (e).

(1)

---

---

---

---

---

---

---

---

---

---





Leave  
blank

5. In a shopping survey a random sample of 104 teenagers were asked how many hours, to the nearest hour, they spent shopping in the last month. The results are summarised in the table below.

Number of hours	Mid-point	Frequency
0 – 5	2.75	20
6 – 7	6.5	16
8 – 10	9	18
11 – 15	13	25
16 – 25	20.5	15
26 – 50	38	10

A histogram was drawn and the group (8 – 10) hours was represented by a rectangle that was 1.5 cm wide and 3 cm high.

- (a) Calculate the width and height of the rectangle representing the group (16 – 25) hours. (3)
- (b) Use linear interpolation to estimate the median and interquartile range. (5)
- (c) Estimate the mean and standard deviation of the number of hours spent shopping. (4)
- (d) State, giving a reason, the skewness of these data. (2)
- (e) State, giving a reason, which average and measure of dispersion you would recommend to use to summarise these data. (2)



11/11/2016



114

---



Foot length, $l$ , (cm)	Number of children
$10 \leq l < 12$	5
$12 \leq l < 17$	53
$17 \leq l < 19$	29
$19 \leq l < 21$	15
$21 \leq l < 23$	11
$23 \leq l < 25$	7

- (a) Use interpolation to estimate the median of this distribution. (2)

- (b) Calculate estimates for the mean and the standard deviation of these data. (6)

One measure of skewness is given by

$$\text{Coefficient of skewness} = \frac{3(\text{mean} - \text{median})}{\text{standard deviation}}$$

- (c) Evaluate this coefficient and comment on the skewness of these data. (3)

Greg suggests that a normal distribution is a suitable model for the foot lengths of ten-year-old children.

- (d) Using the value found in part (c), comment on Greg's suggestion, giving a reason for your answer. (2)



Leave  
blank

**Question 4 continued**





Leave  
blank

2. The 19 employees of a company take an aptitude test. The scores out of 40 are illustrated in the stem and leaf diagram below.

	2 6 means a score of 26	
0	7	(1)
1	88	(2)
2	4468	(4)
3	2333459	(7)
4	00000	(5)

Find

- the median score, (1)
- the interquartile range. (3)

The company director decides that any employees whose scores are so low that they are outliers will undergo retraining.

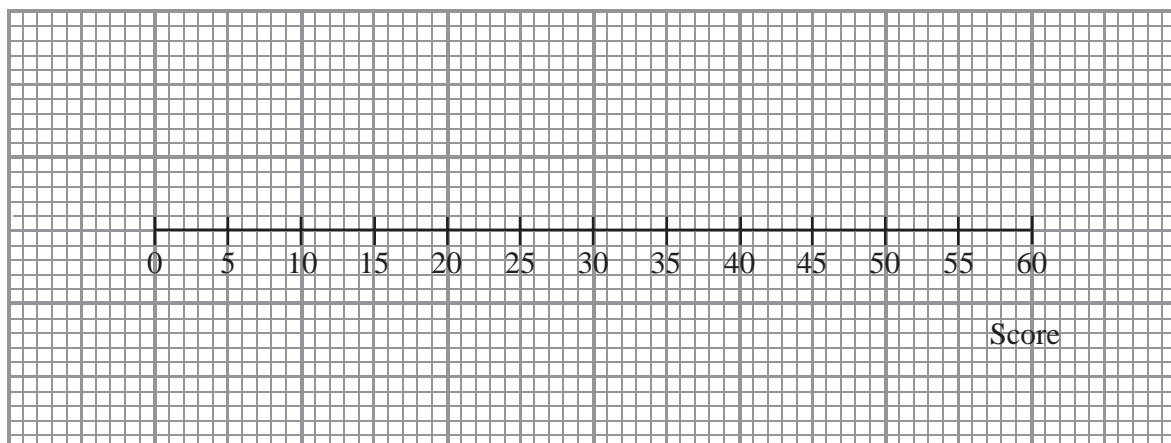
An outlier is an observation whose value is less than the lower quartile minus 1.0 times the interquartile range.

- (c) Explain why there is only one employee who will undergo retraining. (2)
- (d) On the graph paper on page 5, draw a box plot to illustrate the employees' scores. (3)

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

Leave  
blank

### Question 2 continued



Q2

**(Total 9 marks)**



3. The birth weights, in kg, of 1500 babies are summarised in the table below.

Weight (kg)	Midpoint, $x$ kg	Frequency, $f$
0.0 – 1.0	0.50	1
1.0 – 2.0	1.50	6
2.0 – 2.5	2.25	60
2.5 – 3.0		280
3.0 – 3.5	3.25	820
3.5 – 4.0	3.75	320
4.0 – 5.0	4.50	10
5.0 – 6.0		3

[You may use  $\sum fx = 4841$  and  $\sum fx^2 = 15\,889.5$ ]

- (a) Write down the missing midpoints in the table above. (2)
- (b) Calculate an estimate of the mean birth weight. (2)
- (c) Calculate an estimate of the standard deviation of the birth weight. (3)
- (d) Use interpolation to estimate the median birth weight. (2)
- (e) Describe the skewness of the distribution. Give a reason for your answer. (2)

---

---

---

---

---

---

---

---

---

---

---

---



Leave  
blank

### Question 3 continued



Hours	1–10	11–20	21–25	26–30	31–40	41–59
Frequency	6	15	11	13	8	3
Mid-point	5.5	15.5		28		50

- A histogram was drawn to represent these data. The 11–20 group was represented by a bar of width 4 cm and height 6 cm.

- (c) Estimate the mean and standard deviation of the time spent watching television by these students. (5)

- (d) Use linear interpolation to estimate the median length of time spent watching television by these students. (2)

The teacher estimated the lower quartile and the upper quartile of the time spent watching television to be 15.8 and 29.3 respectively.

- (e) State, giving a reason, the skewness of these data. (2)

This image shows a single page of white paper with horizontal grey lines, resembling notebook paper. The lines are evenly spaced and run across the width of the page. There is no handwriting or other markings on the paper.

Leave  
blank

**Question 5 continued**

[illegible]

2. Keith records the amount of rainfall, in mm, at his school, each day for a week. The results are given below.

2.8      5.6      2.3      9.4      0.0      0.5      1.8

Jenny then records the amount of rainfall,  $x$  mm, at the school each day for the following 21 days. The results for the 21 days are summarised below.

$$\Sigma x = 84.6$$

- (a) Calculate the mean amount of rainfall during the whole 28 days. (2)

Keith realises that he has transposed two of his figures. The number 9.4 should have been 4.9 and the number 0.5 should have been 5.0  
Keith corrects these figures.

- (b) State, giving your reason, the effect this will have on the mean. (2)



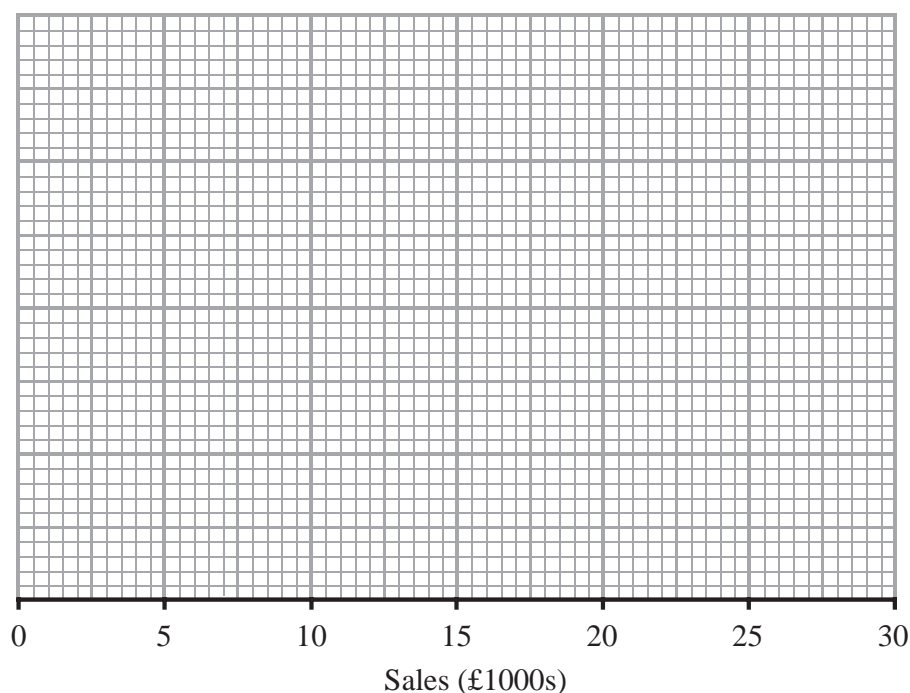
3. Over a long period of time a small company recorded the amount it received in sales per month. The results are summarised below.

	Amount received in sales (£1000s)
Two lowest values	3, 4
Lower quartile	7
Median	12
Upper quartile	14
Two highest values	20, 25

An outlier is an observation that falls either  $1.5 \times$  interquartile range above the upper quartile or  $1.5 \times$  interquartile range below the lower quartile.

- (a) On the graph paper below, draw a box plot to represent these data, indicating clearly any outliers.

(5)



- (b) State the skewness of the distribution of the amount of sales received. Justify your answer.

(2)

- (c) The company claims that for 75 % of the months, the amount received per month is greater than £10000. Comment on this claim, giving a reason for your answer.

(2)





Leave  
blank

Question 3 continued

Lined area for writing the answer to Question 3.

Q3

(Total 9 marks)



5. On a randomly chosen day, each of the 32 students in a class recorded the time,  $t$  minutes to the nearest minute, they spent on their homework. The data for the class is summarised in the following table.

Time, $t$	Number of students
10 – 19	2
20 – 29	4
30 – 39	8
40 – 49	11
50 – 69	5
70 – 79	2

- (a) Use interpolation to estimate the value of the median.

(2)

Given that

$$\sum t = 1414 \quad \text{and} \quad \sum t^2 = 69378$$

- (b) find the mean and the standard deviation of the times spent by the students on their homework.

(3)

- (c) Comment on the skewness of the distribution of the times spent by the students on their homework. Give a reason for your answer.

(2)

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

Leave  
blank

Question 5 continued

Lined area for writing the answer to Question 5.

(Total 7 marks)

Q5

Mark box for Q5.



5. A class of students had a sudoku competition. The time taken for each student to complete the sudoku was recorded to the nearest minute and the results are summarised in the table below.

Time	Mid-point, $x$	Frequency, $f$
2 - 8	5	2
9 - 12		7
13 - 15	14	5
16 - 18	17	8
19 - 22	20.5	4
23 - 30	26.5	4

(You may use  $\sum fx^2 = 8603.75$ )

- (a) Write down the mid-point for the 9 - 12 interval. (1)
- (b) Use linear interpolation to estimate the median time taken by the students. (2)
- (c) Estimate the mean and standard deviation of the times taken by the students. (5)

The teacher suggested that a normal distribution could be used to model the times taken by the students to complete the sudoku.

- (d) Give a reason to support the use of a normal distribution in this case. (1)

On another occasion the teacher calculated the quartiles for the times taken by the students to complete a different sudoku and found

$$Q_1 = 8.5 \quad Q_2 = 13.0 \quad Q_3 = 21.0$$

- (e) Describe, giving a reason, the skewness of the times on this occasion. (2)



Leave  
blank

**Question 5 continued**

Lined area for writing the answer to Question 5.



Leave  
blank

4. The marks,  $x$ , of 45 students randomly selected from those students who sat a mathematics examination are shown in the stem and leaf diagram below.

Mark	Totals	Key	(3 6 means 36)
3	6 9 9	(3)	
4	0 1 2 2 3 4	(6)	
4	5 6 6 6 8	(5)	
5	0 2 3 3 4 4	(6)	
5	5 5 6 7 7 9	(6)	
6	0 0 0 0 1 3 4 4 4	(9)	
6	5 5 6 7 8 9	(6)	
7	1 2 3 3	(4)	

- (a) Write down the modal mark of these students. **(1)**
- (b) Find the values of the lower quartile, the median and the upper quartile. **(3)**

For these students  $\sum x = 2497$  and  $\sum x^2 = 143\,369$

- (c) Find the mean and the standard deviation of the marks of these students. (3)
- (d) Describe the skewness of the marks of these students, giving a reason for your answer. (2)

The mean and standard deviation of the marks of all the students who sat the examination were 55 and 10 respectively. The examiners decided that the total mark of each student should be scaled by subtracting 5 marks and then reducing the mark by a further 10 %.

- (e) Find the mean and standard deviation of the scaled marks of all the students. (4)

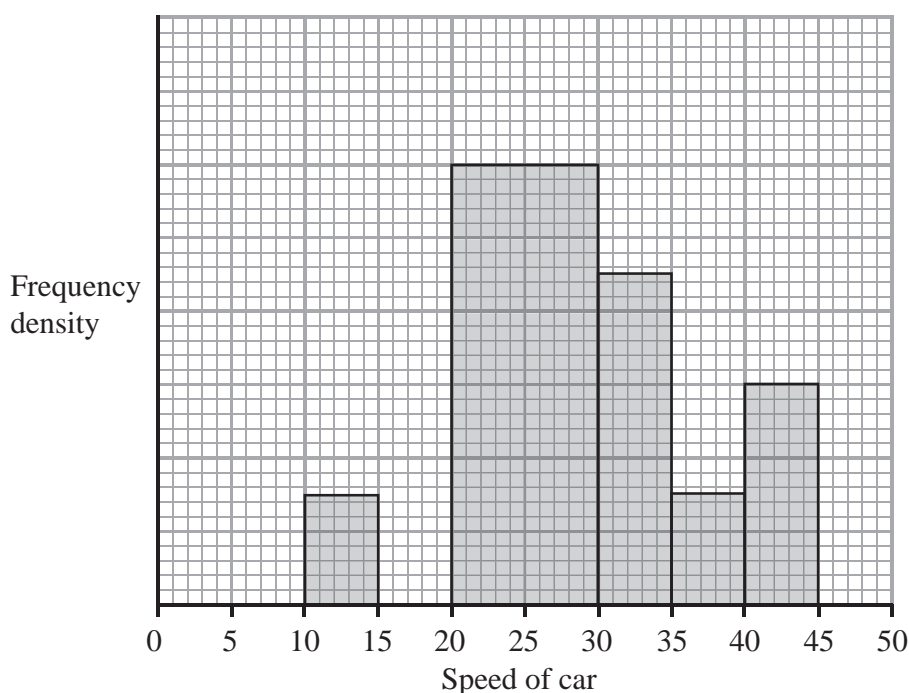


Leave  
blank

**Question 4 continued**



5.



**Figure 2**

A policeman records the speed of the traffic on a busy road with a 30 mph speed limit. He records the speeds of a sample of 450 cars. The histogram in Figure 2 represents the results.

- Calculate the number of cars that were exceeding the speed limit by at least 5 mph in the sample. (4)
- Estimate the value of the mean speed of the cars in the sample. (3)
- Estimate, to 1 decimal place, the value of the median speed of the cars in the sample. (2)
- Comment on the shape of the distribution. Give a reason for your answer. (2)
- State, with a reason, whether the estimate of the mean or the median is a better representation of the average speed of the traffic on the road. (2)

---

---

---

---

---

---





**Question 5 continued**



5. A survey of 100 households gave the following results for weekly income £y.

Income y (£)	Mid-point	Frequency $f$
$0 \leq y < 200$	100	12
$200 \leq y < 240$	220	28
$240 \leq y < 320$	280	22
$320 \leq y < 400$	360	18
$400 \leq y < 600$	500	12
$600 \leq y < 800$	700	8

(You may use  $\sum fy^2 = 12\,452\,800$ )

A histogram was drawn and the class  $200 \leq y < 240$  was represented by a rectangle of width 2 cm and height 7 cm.

- (a) Calculate the width and the height of the rectangle representing the class  $320 \leq y < 400$  (3)
- (b) Use linear interpolation to estimate the median weekly income to the nearest pound. (2)
- (c) Estimate the mean and the standard deviation of the weekly income for these data. (4)

One measure of skewness is  $\frac{3(\text{mean} - \text{median})}{\text{standard deviation}}$ .

- (d) Use this measure to calculate the skewness for these data and describe its value. (2)

Katie suggests using the random variable  $X$  which has a normal distribution with mean 320 and standard deviation 150 to model the weekly income for these data.

- (e) Find  $P(240 < X < 400)$ . (2)
- (f) With reference to your calculations in parts (d) and (e) and the data in the table, comment on Katie's suggestion. (2)

---

---

---

---

---

---

---

---

---

---

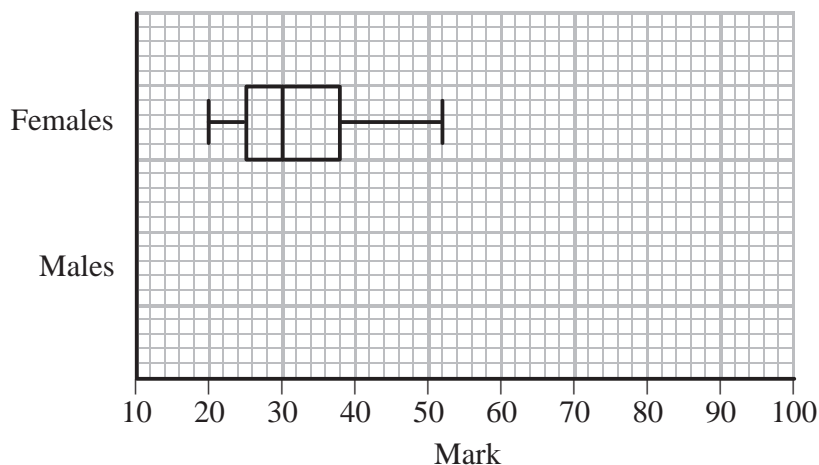


Leave  
blank

**Question 5 continued**



2. The marks of a group of female students in a statistics test are summarised in Figure 1



**Figure 1**

- (a) Write down the mark which is exceeded by 75% of the female students.

(1)

The marks of a group of male students in the same statistics test are summarised by the stem and leaf diagram below.

Mark	(2 6 means 26)	Totals
1	4	(1)
2	6	(1)
3	4 4 7	(3)
4	0 6 6 7 7 8	(6)
5	0 0 1 1 1 3 6 7 7	(9)
6	2 2 3 3 3 8	(6)
7	0 0 8	(3)
8	5	(1)
9	0	(1)

- (b) Find the median and interquartile range of the marks of the male students.

(3)

An outlier is a mark that is

either more than  $1.5 \times$  interquartile range above the upper quartile

or more than  $1.5 \times$  interquartile range below the lower quartile.



**(2)**



4. The following table summarises the times,  $t$  minutes to the nearest minute, recorded for a group of students to complete an exam.

Time (minutes) $t$	11 – 20	21 – 25	26 – 30	31 – 35	36 – 45	46 – 60
Number of students $f$	62	88	16	13	11	10

[You may use  $\sum ft^2 = 134281.25$ ]

- Estimate the mean and standard deviation of these data. (5)
- Use linear interpolation to estimate the value of the median. (2)
- Show that the estimated value of the lower quartile is 18.6 to 3 significant figures. (1)
- Estimate the interquartile range of this distribution. (2)
- Give a reason why the mean and standard deviation are not the most appropriate summary statistics to use with these data. (1)

The person timing the exam made an error and each student actually took 5 minutes less than the times recorded above. The table below summarises the actual times.

Time (minutes) $t$	6 – 15	16 – 20	21 – 25	26 – 30	31 – 40	41 – 55
Number of students $f$	62	88	16	13	11	10

- (f) Without further calculations, explain the effect this would have on each of the estimates found in parts (a), (b), (c) and (d). (3)



### Question 4 continued



## Statistics S1

### Probability

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A \cap B) = P(A)P(B | A)$$

$$P(A | B) = \frac{P(B | A)P(A)}{P(B | A)P(A) + P(B | A')P(A')}$$

### Discrete distributions

For a discrete random variable  $X$  taking values  $x_i$  with probabilities  $P(X = x_i)$

Expectation (mean):  $E(X) = \mu = \sum x_i P(X = x_i)$

Variance:  $\text{Var}(X) = \sigma^2 = \sum (x_i - \mu)^2 P(X = x_i) = \sum x_i^2 P(X = x_i) - \mu^2$

For a function  $g(X)$ :  $E(g(X)) = \sum g(x_i) P(X = x_i)$

### Continuous distributions

Standard continuous distribution:

Distribution of $X$	P.D.F.	Mean	Variance
Normal $N(\mu, \sigma^2)$	$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$	$\mu$	$\sigma^2$



### ***Correlation and regression***

For a set of  $n$  pairs of values  $(x_i, y_i)$

$$S_{xx} = \Sigma(x_i - \bar{x})^2 = \Sigma x_i^2 - \frac{(\Sigma x_i)^2}{n}$$

$$S_{yy} = \Sigma(y_i - \bar{y})^2 = \Sigma y_i^2 - \frac{(\Sigma y_i)^2}{n}$$

$$S_{xy} = \Sigma(x_i - \bar{x})(y_i - \bar{y}) = \Sigma x_i y_i - \frac{(\Sigma x_i)(\Sigma y_i)}{n}$$

The product moment correlation coefficient is

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

The regression coefficient of  $y$  on  $x$  is  $b = \frac{S_{xy}}{S_{xx}} = \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{\Sigma(x_i - \bar{x})^2}$

Least squares regression line of  $y$  on  $x$  is  $y = a + bx$  where  $a = \bar{y} - b\bar{x}$

# THE NORMAL DISTRIBUTION FUNCTION

The function tabulated below is  $\Phi(z)$ , defined as  $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{1}{2}t^2} dt$ .

$z$	$\Phi(z)$	$z$	$\Phi(z)$	$z$	$\Phi(z)$	$z$	$\Phi(z)$	$z$	$\Phi(z)$
0.00	0.5000	0.50	0.6915	1.00	0.8413	1.50	0.9332	2.00	0.9772
0.01	0.5040	0.51	0.6950	1.01	0.8438	1.51	0.9345	2.02	0.9783
0.02	0.5080	0.52	0.6985	1.02	0.8461	1.52	0.9357	2.04	0.9793
0.03	0.5120	0.53	0.7019	1.03	0.8485	1.53	0.9370	2.06	0.9803
0.04	0.5160	0.54	0.7054	1.04	0.8508	1.54	0.9382	2.08	0.9812
0.05	0.5199	0.55	0.7088	1.05	0.8531	1.55	0.9394	2.10	0.9821
0.06	0.5239	0.56	0.7123	1.06	0.8554	1.56	0.9406	2.12	0.9830
0.07	0.5279	0.57	0.7157	1.07	0.8577	1.57	0.9418	2.14	0.9838
0.08	0.5319	0.58	0.7190	1.08	0.8599	1.58	0.9429	2.16	0.9846
0.09	0.5359	0.59	0.7224	1.09	0.8621	1.59	0.9441	2.18	0.9854
0.10	0.5398	0.60	0.7257	1.10	0.8643	1.60	0.9452	2.20	0.9861
0.11	0.5438	0.61	0.7291	1.11	0.8665	1.61	0.9463	2.22	0.9868
0.12	0.5478	0.62	0.7324	1.12	0.8686	1.62	0.9474	2.24	0.9875
0.13	0.5517	0.63	0.7357	1.13	0.8708	1.63	0.9484	2.26	0.9881
0.14	0.5557	0.64	0.7389	1.14	0.8729	1.64	0.9495	2.28	0.9887
0.15	0.5596	0.65	0.7422	1.15	0.8749	1.65	0.9505	2.30	0.9893
0.16	0.5636	0.66	0.7454	1.16	0.8770	1.66	0.9515	2.32	0.9898
0.17	0.5675	0.67	0.7486	1.17	0.8790	1.67	0.9525	2.34	0.9904
0.18	0.5714	0.68	0.7517	1.18	0.8810	1.68	0.9535	2.36	0.9909
0.19	0.5753	0.69	0.7549	1.19	0.8830	1.69	0.9545	2.38	0.9913
0.20	0.5793	0.70	0.7580	1.20	0.8849	1.70	0.9554	2.40	0.9918
0.21	0.5832	0.71	0.7611	1.21	0.8869	1.71	0.9564	2.42	0.9922
0.22	0.5871	0.72	0.7642	1.22	0.8888	1.72	0.9573	2.44	0.9927
0.23	0.5910	0.73	0.7673	1.23	0.8907	1.73	0.9582	2.46	0.9931
0.24	0.5948	0.74	0.7704	1.24	0.8925	1.74	0.9591	2.48	0.9934
0.25	0.5987	0.75	0.7734	1.25	0.8944	1.75	0.9599	2.50	0.9938
0.26	0.6026	0.76	0.7764	1.26	0.8962	1.76	0.9608	2.55	0.9946
0.27	0.6064	0.77	0.7794	1.27	0.8980	1.77	0.9616	2.60	0.9953
0.28	0.6103	0.78	0.7823	1.28	0.8997	1.78	0.9625	2.65	0.9960
0.29	0.6141	0.79	0.7852	1.29	0.9015	1.79	0.9633	2.70	0.9965
0.30	0.6179	0.80	0.7881	1.30	0.9032	1.80	0.9641	2.75	0.9970
0.31	0.6217	0.81	0.7910	1.31	0.9049	1.81	0.9649	2.80	0.9974
0.32	0.6255	0.82	0.7939	1.32	0.9066	1.82	0.9656	2.85	0.9978
0.33	0.6293	0.83	0.7967	1.33	0.9082	1.83	0.9664	2.90	0.9981
0.34	0.6331	0.84	0.7995	1.34	0.9099	1.84	0.9671	2.95	0.9984
0.35	0.6368	0.85	0.8023	1.35	0.9115	1.85	0.9678	3.00	0.9987
0.36	0.6406	0.86	0.8051	1.36	0.9131	1.86	0.9686	3.05	0.9989
0.37	0.6443	0.87	0.8078	1.37	0.9147	1.87	0.9693	3.10	0.9990
0.38	0.6480	0.88	0.8106	1.38	0.9162	1.88	0.9699	3.15	0.9992
0.39	0.6517	0.89	0.8133	1.39	0.9177	1.89	0.9706	3.20	0.9993
0.40	0.6554	0.90	0.8159	1.40	0.9192	1.90	0.9713	3.25	0.9994
0.41	0.6591	0.91	0.8186	1.41	0.9207	1.91	0.9719	3.30	0.9995
0.42	0.6628	0.92	0.8212	1.42	0.9222	1.92	0.9726	3.35	0.9996
0.43	0.6664	0.93	0.8238	1.43	0.9236	1.93	0.9732	3.40	0.9997
0.44	0.6700	0.94	0.8264	1.44	0.9251	1.94	0.9738	3.50	0.9998
0.45	0.6736	0.95	0.8289	1.45	0.9265	1.95	0.9744	3.60	0.9998
0.46	0.6772	0.96	0.8315	1.46	0.9279	1.96	0.9750	3.70	0.9999
0.47	0.6808	0.97	0.8340	1.47	0.9292	1.97	0.9756	3.80	0.9999
0.48	0.6844	0.98	0.8365	1.48	0.9306	1.98	0.9761	3.90	1.0000
0.49	0.6879	0.99	0.8389	1.49	0.9319	1.99	0.9767	4.00	1.0000
0.50	0.6915	1.00	0.8413	1.50	0.9332	2.00	0.9772		

## PERCENTAGE POINTS OF THE NORMAL DISTRIBUTION

The values  $z$  in the table are those which a random variable  $Z \sim N(0, 1)$  exceeds with probability  $p$ ; that is,  $P(Z > z) = 1 - \Phi(z) = p$ .

$p$	$z$	$p$	$z$
0.5000	0.0000	0.0500	1.6449
0.4000	0.2533	0.0250	1.9600
0.3000	0.5244	0.0100	2.3263
0.2000	0.8416	0.0050	2.5758
0.1500	1.0364	0.0010	3.0902
0.1000	1.2816	0.0005	3.2905