

Y1S1 XMQs and MS

(Total: 17 marks)

1. P3_2018 Q4 . 13 marks - Y1S1 Data collection
2. P31(AS)_2021 Q3 . 4 marks - Y1S1 Data collection

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 08 40 to 08 50 one morning and asks workers, as they arrive, how long their journey was.

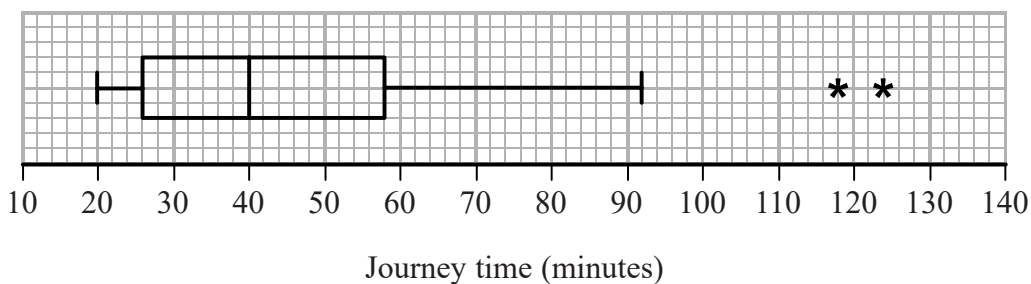
(a) State the sampling method Charlie used. (1)

(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)



DO NOT WRITE IN THIS AREA

Qu 4	Scheme	Marks	AO
(a)	Convenience <u>or</u> opportunity [sampling]	B1 (1)	1.2
(b)	Quota [sampling] e.g. Take 4 people every 10 minutes	B1 B1 (2)	1.1a 1.1b
(c)	Census	B1 (1)	1.2
(d)	[58 – 26 =] 32 (min)	B1 (1)	1.1b
(e)	$\mu = \frac{4133}{95} = 43.505263\dots$ $\sigma_x = \sqrt{\frac{202\,294}{95} - \mu^2} = \sqrt{236.7026\dots}$ $= 15.385\dots$ awrt 15.4 (min)	awrt 43.5 (min) B1 M1 A1 (3)	1.1b 1.1b 1.1b
(f)	There are outliers in the data (or data is skew) which will affect mean and sd Therefore use median and IQR	B1 dB1 (2)	2.4 2.4
(g)	Value of 20, LQ at 26 and outliers will not change <u>or</u> state that median and upper quartile are the values that <u>do</u> change <u>More values now below 40 than above so Q_2 or Q_3 will change and be lower</u> <u>Both Q_2 and Q_3 will be lower</u>	B1 M1 A1 (3)	1.1b 2.1 2.4
		(13 marks)	
Notes			
(b)	1 st B1 for quota (sampling) mentioned (“Stratified” or “systematic” or “random” are B0B0) 2 nd B1 for a description of how such a system might work, requires suitable strata or categories e.g. time slots, departments, gender, age groups, distance travelled etc Suggestion of randomness is B0		
(e)	B1 for a correct mean (awrt 43.5) M1 for a correct expression for the sd (including $\sqrt{\quad}$)ft their mean A1 for awrt 15.4 (Allow $s = 15.4667\dots$ awrt 15.5)		
(f)	1 st B1 for acknowledging <u>outliers</u> or <u>skewness</u> are a problem for <u>mean and sd</u> “extreme values”/“anomalies” OK May be implied by saying median and IQR not affected by.. We need to see mention of “outliers”, “skewness” and the problem so “data is skewed so use median and IQR” is B0 unless mention that they are not affected by extreme values <u>or</u> mean and standard deviation can be “inflated” by the positive skew etc 2 nd dB1 dep on 1 st B1 for therefore choosing <u>median and IQR</u>		
(g)	B1 for identifying 2 of these 3 groups of unchanged values or stating only Q_2 and Q_3 change M1 for <u>explaining</u> that median or UQ should be lower. E.g. the 2 values have moved to below 40 (or 58) and therefore more than 50% below 40 or (more than 75% below 58) <u>or</u> an argument to show that the other 3 values are the same. (o.e.) Allow arrows on box plot provided statement in words about increased % below 40 or 58 etc A1 for stating median <u>and</u> UQ are both lower with clear evidence of M1 scored [If lots of values on 40 then median might not change but, since two values <u>do</u> change then UQ would change. If this meant that 92 became an outlier then we would have a new value for upper whisker and an extra outlier so effectively 3 values are altered. So median changes]		

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)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA



Qu	Scheme	Marks	AO
3. (a)	Systematic (sampling)	B1 (1)	1.2
(b)(i)	[Daily Mean] Wind Speed	B1	2.2a
(ii)	Light	B1 (2)	1.2
(c)	Variable A occurs most (around 80~90%) of the time	B1 (1)	2.2b
Notes			
(a)	B1 for identifying the correct sampling technique Allow slight misspelling e.g. “sysmatic”, “sytmatic” Do NOT allow “systemic”		
(b)(i)	B1 for identifying appropriate qualitative variable. {LDS mark} Allow “Wind speed” or “Wind strength” but NOT just “wind” or “wind direction”		
(ii)	B1 for realising that modal wind speed is “Light” {LDS mark} Allow just “light” or “most light”		
NB	These two B marks are independent so can score B0B1 for e.g. “rainfall” and “light”		
(c)	B1 for inferring that frequency of A can be estimated fairly reliably: {underestimates B and over estimates C} e.g. “A is the most frequent” [can then ignore comments about B and C]		