Research Methods 2 Week 7: Exercise Sheet 1 Solution sheet

Question 1

The data from the question have been typed into column C2 and the data from 1977-83 into column C5. Using **Display Descriptive Statistics** the following screen is obtained. (Note that the descriptive statistics for both columns have been found using one application of the command: this has been done by selecting first C2 and then selecting C5 before clicking on **O**K)

ΣΜΙΝΙΤ	AB - Untit	led								_ 8 ×	
<u>File Edit Manip Calc Stat Graph Editor Window H</u> elp											
🖺 Session 📃 🗌 🔀											
Desc	riptive \$	Statistics	s: C2, C	:5							
Varia	able		N	Mean	Median	TrMean		StDev	SE Mean		
C2			34	1.229	1.155	1.144		0.667	0.114		
C5			7	1.673	1.290	1.673		1.201	0.454		
Varia	able	Minim	imum Maximum		01	03					
C2	C2 0		470 4.000		0.815	1.405					
C5		0.5	00	4.000	0.820	2.500					
										_	
•										<u>•</u> //.	
sex r	atio.MTW	***							<i>0</i> 2	- 🗆 ×	
+	C1	C2	C3	C4	C5	C6	C7	C8	C9	C1L	
1	1950	0.80		1977	1.17						
2	1951	0.67		1978	0.50			1			
3	1952	0.82		1979	4.00						
4	1953	1.42		1980	0.82						
5	1954	1.11		1981	2.50						
								le v			
Current Worksheet: sex ratio.MTW [Editable] 11:49											
Start Start	📇 Micros	oft Scen] 🤇	🔟 Week/	🛄 M	icrosoft Wor		B - U	ULotus U	irqanize 💓 🤇	2 73 11:49	

For the full data, the 95% confidence interval for the mean sex ratio is computed as

Lower limit = $1.229 - 2 \times 0.114 = 1.001$ Upper Limit = $1.229 + 2 \times 0.114 = 1.457$

For the data from 1977-83, the 95% confidence interval is found as

Lower limit = $1.673 - 2 \times 0.454 = 0.765$ Upper Limit = $1.673 + 2 \times 0.454 = 2.581$

Question 2

Clicking the sequence <u>Stat</u> -> <u>Basic</u> Statistics -> <u>1</u>-sample t... and selecting C2 into the <u>Variables</u>: box and clicking <u>OK</u> gives the following screen.

MINI	TAB - Untit	led								- 8 ×
<u>File E</u> dit <u>M</u> anip <u>C</u> alc <u>S</u> tat <u>G</u> raph E <u>d</u> itor <u>W</u> indow <u>H</u> elp										
2					1 4 8	0 ? 🗗	9	<u>+C</u>	602	🗒 🎦 🖄
🖺 Session 📃 🗖 🕅										
MTB > OneT C2.										
One	-Sample	T: C2								
C2 MTB MTB	> >	34	1.22	9 0.66	57 O.I	114 (93.0 0.996,	1.462)		
MTB	>									- - -
MTB	>	***								
MTB	> ratio.MTW C1	*** C2	C3	C4	C5	C6	C7	C8	C9	
MTB	> ratio.MTW C1 1950	•••• C2 0.80	C3	C4 1977	C5	C6	C7	C8	C9	• • • • •
MTB Sex 1 1 2	> ratio.MTW C1 1950 1951	C2 0.80 0.67	C3	C4 1977 1978	C5 1.17 0.50	C6	C7	C8	C9	↓ ↓ C1L
MTB sex + 1 2 3	> ratio.MTW C1 1950 1951 1952	*** C2 0.80 0.67 0.82	C3	C4 1977 1978 1979	C5 1.17 0.50 4.00	C6	C7	C8	C9	• • • • •
MTB * * 1 2 3 4	> ratio.MTW C1 1950 1951 1952 1953	C2 0.80 0.67 0.82 1.42	C3	C4 1977 1978 1979 1980	C5 1.17 0.50 4.00 0.82	C6	C7	C8	C9	
MTB Sex 1 2 3 4 5	> ratio.MTW C1 1950 1951 1952 1953 1954	ZZZ J C2 0 0.80 0 0.67 0 0.82 1 1.42 1	C3	C4 1977 1978 1979 1980 1981	C5 1.17 0.50 4.00 0.82 2.50	C6	C7	C8	C9	
MTB Sex 1 2 3 4 5 (> ratio.MTW C1 1950 1951 1952 1953 1954	C2 I 0.80 0.67 0.82 1.42 1.11	C3	C4 1977 1977 1978 1979 1980 1981	C5 1.17 0.50 4.00 0.82 2.50	C6	C7	C8	C9	
MTB • Sex • 1 1 2 3 4 5 • Current V	> ratio.MTW C1 1950 1951 1952 1953 1954 Worksheet s	C2 0.80 0.67 0.82 1.42 1.11 ex ratio.MTW	C3	C4 1977 1978 1979 1980 1981	C5 1.17 0.50 4.00 0.82 2.50	C6	C7	C8 Editable	C9	

The output from the command repeats the values of the sample mean and SD and the SE and also produces the limits of the confidence interval.

Repeating the procedure with column C5 gives the interval (0.561, 2.784)

The intervals found in questions 1 and two are as follows

Method	Full data	Data 1977-83
Question 1	(1.001, 1.457)	(0.765, 2.581)
Question 2	(0.996, 1.462)	(0.561,2.784)

Two comments are made.

- 1. The intervals based on the reduced sample are much wider than those for the full sample. It is not surprising that we should have a better idea about the location of the mean of a population if we base the inference on 34 rather than just 7 points.
- 2. The methods in the two questions produced very similar answers when applied to the full sample. When applied to the reduced sample there is a greater discrepancy. This relates to the point made in the note 'Confidence intervals: a technicality' in the study document for this week. For small samples smaller than about 25-30 the method based on mean $\pm 2 \times SE$ gives intervals that are too narrow and the exact method obtained from Minitab directly should be used.

Question 3

Entering the command in the question gives the following screen, once the Tally command has been applied to the indicator column, C9.

≥ MIN⊓	FAB - Unti	tled								- 6 X
Ele E	dit <u>M</u> anip	Calc :	Stat Graph	Editor	Window Help					
6	-	2 8 -				6 A 6 4	2440	11 -0 0	602	
E Ses	sion				- R R.					
i txt MTE : SUBC: Tally C9 0 1	1000. simi endmas Tally Court for Dis Count 49 951	.00 ulation c9; nts. crete V	s complete	ed C9						
iv≕ ∙ I Worl	ibuu									2 2 1
+	C1	C2	C3	C4	C5	C6	C7	C8	C9	C1_
						mean	Lower_CI	Upper_CI	Indicator	5
1						1.03723	0.92981	1.14466	1	8
2						1.02850	D.95271	1.10430	1	
3						0.92390	0.81440	1.03339	0	8
.4						1.06376	N 96324	1 16427	1	्यां
Durrent M	vorksheet '	Warksheel	1						13	54
Start	BMicro	soft Scen	Whicrosoft	Wor.	MSc Oncology	My Co	imputer	MINITAB -	U. 🕅 🏹 🗄	13:54

Note: in the Data Window the population mean, 1.05, is between the upper and lower limits in rows 1, 2 and 4 but not in row 3. Hence column C9 contains 1s on rows 11, 2 and 4 and 0 in row 3. Note also that the column names have been added for clarity

The claim for a 95% confidence interval is that on 95% of occasions the interval will include the population mean. Of course, in practice, it will not be known whether or not a given interval will include the population mean, as this quantity is unknown. It is simply that it is likely to include the population mean - indeed, with probability 95%. However, if data are generated from a known population, the population mean *is* known and for each of many similarly generated intervals it can be decided whether or not the interval includes that population mean.

From the above, the number of 1s is 951, so out of 1000 intervals, only 49 fail to include the population mean. This is close the value of $95\% \times 1000 = 950$ intervals containing the mean.

Trying various different populations and sample sizes should give similar results. The percentage of intervals containing the population mean, using the same mean and SD above but with different sample sizes and different numbers of samples, are given below.

	Number of sample generated						
Sample size	100	1000	10000				
10	0.94	0.955	0.9470				
20	0.95	0.957	0.9505				
100	0.98	0.943	0.9517				

Question 4

The 95% confidence interval for the sex ratio, computed using the Minitab command from the whole sample is (0.996, 1.462). It is probably sensible to round this to 2 decimal places, giving (1.00, 1.46).

The sample mean (to 2 d.p.) is 1.23. If we had to estimate the mean sex ratio by a single value, then we would use this value. A more complete estimate, which acknowledges the imprecision which necessarily arises from basing our inference on a sample of 34 values, is to say that we believe the population mean sex ratio is between 1.00 and 1.46, and we have a 5% chance of being in error in this assertion.

Confidence intervals are valuable in two respects. First, we have a reasonable belief that the sex ratio is certainly not below 1 or above 1.5 (to 1 d.p.), i.e. the values outside the interval have reasonably been excluded. Second, our sample has not allowed us to pin the population mean down any better than between 1 and 1.5. In other words, a confidence interval allows us to quantify the imprecision in our knowledge in a way that an estimate based on a single value cannot do. Moreover, it is often the case that the intervals are wider than we might have imagined, i.e. our intuitions are often over-confident.

End of solution sheet