Research Methods 2

Week 6: Exercise Sheet 1

Solution sheet

[Remember: for questions 1 to 4, the exercises are based on randomly generated data, so while the numerical answers you obtain ought to be quite similar to those given here, they will not agree precisely. However, the important points of principle behind the numerical examples will be the same]

Question 1

Generating data as suggested in the question, and then using <u>**D**</u>isplay **Descriptive Statistics...** item gives the following screen.

× II	ITAB - Unli	lied														18
is į	dit Menip I	Qalo <u>S</u> lat J	Breph Eglé	or Wind	law Help											
9		中国	-7 🛄		E B	114	9	1				1	6 🖶 🖂	02		
I Se	ssion															
_	0	5/01/02 10	:48:23 -	2				- 83 -								4
Uelo ISTB ISTB ISTB ISTB ISTB	ome to Mi > Bare 12 > Random > Borma > Describ	nitab, pro 3. 90 cl; 1 13.5 l.1 c Cl.	023 Fl fd D.	or help												
Des	criptive S	itatistics:	C1													
Var: Cl	able	M 90	И: 13,5	ean) 505	Hedian 13.516	Trffean 13,549	StDev 1.029	55 B 0,	Dean 108							
Var: CL	uble	Miniawa 9.994	Nex13 15.3	8000 927	01 12.918	03 14.151										
ETE	5															- 1
1000	283															
4	_	100			-					-						1
1 ~	ark riverst 1 *															
+	CI	a	C3	C4	C5	06	7	C8	0	C10	C11	C12	C13	C14	C15	C16_
1	15.0250				-		_									-
2	11.5395															
3	15.2438															
4	13.4236															
5	14.6792															1
6	13.8976															
7	14.4946															
8	15.1812															
9	13.0919															
•																. PE
H	2000 Million	e III							_		_			_	_	
Avert	Wolksheet W	folksheet 1												Editable	1	0.53
-			Die	for Arrest	ant 301	Cateriantes		1 m	Minnesky	Vined - Evening		NITAR - II	titled	1	-	00 1050
		ALC: 1995	0		21.52 SZ										A	10000

The standard deviation of this sample is 1.028 g/dl. You can estimate the standard error of the mean of this sample by dividing 1.028 by $\sqrt{90}$, which gives

Standard Error = $1.028/\sqrt{90} = 1.028/9.487 = 0.108$.

Notice that this is the same value that Minitab prints under SE Mean. In other words Minitab automatically calculates the Standard error for you.

Note that your sample will have a slightly different SD and SE, but in all cases the result of dividing the SD by $\sqrt{90}$ will be exactly the SE.

le j	<u>E</u> dit <u>M</u> anip	<u>C</u> alc <u>S</u> ta	t <u>G</u> raph	E <u>d</u> itor <u>W</u> ir	ndow <u>H</u> elp					
1					1466	0 ? 🗗		* C	₿ ₿ 0 ≥	
Se	ssion									
Des	criptive S	tatistics	: C1							
Var:	lable		N	Mean	Median	TrMean	st	Dev	SE Mean	
Cl		90		.3.603	13.715	13.602	1.	057	0.111	
Variable		Minimum		ximum	Q1	Q3				
UL		11.31	L7 1	.6.208	12.031	14.330				
MTB	>									
мтв	>									
мтв	>									F
MTB J Wo	> rksheet 1 ***	2								<u> </u>
ATB J ₩o	> 	C2	C3	C4	C5	C6	C7	C8	C9) C
4TB ∫ ₩0 ↓ 8	> rksheet 1 *** C1 13.8183	C2	C3	C4	C5 13.599	C6	C7	C8	C9	
4TB ↓ ↓ 8 9	> rksheet 1 *** C1 13.8183 13.3521	C2	C3	C4	C5 13.599 13.449	C6	C7	Ca	C9	
4TB + 8 9 10	> rksheet 1 *** C1 13.8183 13.3521 12.7441	C2	C3	C4	C5 13.599 13.449 13.603	C6	C7	C8	C9	
111 111 111 111	> rksheet 1 *** C1 13.8183 13.3521 12.7441 13.7877	C2	C3	C4	C5 13.599 13.449 13.603	C6	C7	C8	C9	
™® ₩0 ↓ 8 9 10 11 12	> rksheet 1 *** C1 13.8183 13.3521 12.7441 13.7877 13.3208	C2	C3	C4	C5 13.599 13.449 13.603	C6	C7	C8	C9	
111 11 12 11 12 11 12 11 12 11 12 11 12 11 12 11 12 11 12 11 12 11 12 11 12 11 12 12	Ksheet 1 *** C1 13.8183 13.8521 12.7441 13.7877 13.3208	C2	C3	C4	C5 13.599 13.449 13.603	C6	C7	C8	C9	

The above screen shows the result of generating 10 samples, each of size 90 from the population used in question 1. Each sample is stored in C1, the mean of C1 is found using the well-trodden route through <u>Stat</u> \rightarrow <u>Basic Statistics</u> \rightarrow <u>Display</u> **Descriptive Statistics**. The mean is found and copied into the next available row of column C5. The screen shows the situation where C1 holds the tenth sample to be generated, the Session window gives the statistics for C1, including the mean of 13.603, which has been copied into the last row of C5.

At this stage C5 contains the means of ten samples, each of size 90, from the Normal population with mean 13.5 and SD 1. The next stage is to find the sample SD of these ten means. This is done by applying **Display Descriptive Statistics** to C5. This gives output

Descriptive Statistics: C5

Variable	N	Mean	Median	TrMean	StDev	SE Mean
C5	10	13.560	13.590	13.568	0.098	0.031
Variable C5	Minimum 13.364	Maximum 13.694	Q1 13.491	Q3 13.635		

Thus the SD of this collection of sample mean is 0.098 g/dl. This is an estimate of the SE of the Mean of a single sample and is similar to that obtained in question 1, namely 0.108 g/dl. {note that the SE Mean entry in the above (i.e. 0.031) has no meaning as C5 is not raw data but a collection of sample means}.

MINITAB - Ur	titled								- 8
⊡le <u>E</u> dit <u>M</u> an ≊I ⊡I Æsi X	ip <u>C</u> alc <u>S</u> ta I⊡al ∩al ⊳ I	t <u>G</u> raph	Editor <u>W</u> ii	ndow <u>H</u> elp ↑	<u> </u>				al⇒al⊽
Session									
MTB > %samp Executing : MTB > endma MTB > Desca	omns 13.5 : from file: acro ribe C5.	1 90 100 D:\MTBW	0 c5 IN\MACRO	S\sampmns.M	AC				1
Descriptive	e Statistics	: C5							
Variable C5	100	N)0 1	Mean 3.497	Median 13.498	TrMean 13.496	stD 0.1	ev S D5	E Mean 0.003	
Variable C5	Minim 13.1	um Ma 69 1	ximum 3.856	Q1 13.428	Q3 13.566				
<[- I
+ C1	C2	C3	C4	C5	C6	C7	C8	C9	C
8				13.5994					
9				13.4486					
10				13.6028					
urrent Workshee	t: Worksheet 1						Editable	[[1	6:14
Start Start	rosoft S 🔄 M	ISc Oncol	🔄 Week6	1 Micros	oft 🔀 MII	NITAB		V	16:1

By using the macro sampmns macro the above screen is obtained and the SD of the 1000 sample means is 0.105 g/dl. As this SD is based on 1000 means, rather than the 10 used in question 2, it should provide a more precise estimate of the SE of the mean. Indeed, this is the case as the estimate based on 1000 means, namely 0.105 g/dl is closer to that found from the formula used in question 1 (0.108) than the value based on just ten means in question 2.

The histogram of the 1000 sample means is below



The histogram shows that if you take the mean of a sample of 90 haemoglobin concentrations, then this will very likely be within 0.4 g/dl of the population mean. This is because the 1000 means you have generated all lie within a range of 0.8 g/dl, between 13.1 and 13.9 g/dl. Moreover, most of them lie in a narrower central band and the sample means themselves seem to follow a Normal distribution.

The following histogram is analogous to the one above but shows means of samples of size $\ensuremath{\mathsf{9}}$



The broad picture is the same, namely the sample means seem to have a Normal distribution. It is also centred on what we know is the population mean (13.5 g/dl). However, it is much more dispersed than the previous histogram. Therefore, if you take a mean of a sample of size 9, rather than 90, then it will very likely lie within 1 g/dl of the population mean, as the above shows that the 1000 means you have generated lie between 12.5 and 14.5 g/dl.

If you type the data given in the question into column C1 and then apply <u>Stat</u> \rightarrow <u>Basic Statistics</u> \rightarrow <u>Display Descriptive Statistics</u> to C1, you get the following screen

<mark>></mark> MIN	ITAB - Untit	led								_ 8 ×
<u>F</u> ile	<u>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					
B					H 🔏 📲 📲		Ø Ø. 0	18 1	3 6 6 6	
🖽 Se	ssion				~					
Des	scriptive	Statistic	s: polyp	s						
Var.	iable		N	Mean	Median	TrMean	. 3	tDev	SE Mean	
pol	yps		9	25.11	26.00	25.11	1	2.02	4.01	
Var. pol	iable yps	Minin 11	mum M .00	laximum 45.00	Q1 15.00	Q3 36.00				
										- I
للكار										
PO 🖽	LYPOSIS.M	TW ***								_ 🗆 ×
P0	LYPOSIS.M C1	C2	C3	C4	C5	C6	C7	C8	C9	C1
₩ P0 +	C1 C1 polyps	C2	C3	C4	C5	C6	C7	C8	C9	
+ 1	C1 polyps 26	C2	C3	C4	C5	C6	C7	C8	C9	
+ 1 2	C1 polyps 26 16	C2	C3	C4	C5	C6	C7	C8	C9	C1L
+ 1 2 3	C1 polyps 26 16 40	C2	C3	C4	C5	C6	C7	C8	C9	C1
+ 1 2 3 4	C1 polyps 26 16 40 14	C2	C3	C4	C5	C6	C7	C8	C9	C1
+ 1 2 3 4 5	LYPOSIS.M C1 polyps 26 16 40 14 16	C2	C3	C4	C5	C6	C7	C8	C9	
+ 1 2 3 4 5 6	C1 polyps 26 16 40 14 16 11	C2	C3	C4	C5	C6	C7	C8	C9	
+ 1 2 3 4 5 6 7	LYPOSIS.M C1 polyps 26 16 40 14 16 11 26	C2	C3	C4	C5	C6	C7	C8	C9	
+ 1 2 3 4 5 6 4 5 0 0	C1 polyps 26 16 40 14 16 11 26 Worksheet F	C2 C2 POLYPOSIS	C3	C4	C5	C6	C7	C8	C9	

The mean number of polyps per patient in this sample is 25.1. The SE of this mean is 4.01. Therefore the mean of the population from which the sample has been drawn is estimated to be 25.1. The standard error, which measures how good is this estimate of the mean, is 4.01.

Just quoting a SE in this way is not very meaningful. More helpful ways of using the SE to quantify the uncertainty in a sample mean will be discussed next week.

Running the macro three times as described in the question, and saving the three sets of means in columns C1, C2 and C3 (which have been names 'size 10', 'size 50' and 'size 250' respectively) gives the following screen.

≥м	NITAB - Unt	itled								_ 8 ×
<u>F</u> ile	<u>E</u> dit <u>M</u> anip	o <u>C</u> alc <u>S</u> t	tat <u>G</u> raph E	E <u>d</u> itor <u>W</u> inc	łow <u>H</u> elp					
B		te 💼 🔊		II 🗷 📕			0 3 0	0 ?	-0 - 0	2 3 10 10
E, S	ession									
MT	в >									
MT	B > %lgnm	ns 10 100)O c1							
Ex	ecuting f	rom file:	: D:\MTBWI	N\MACROS	\lgnmns.	MAC				
MT.	B > endma B > %larrm	cro 50 100	10 ~2							
Ex	ecuting f	rom file:	: D:\MTBWI	N\MACROS	\lanmns.	MAC				
MT	B > endma	cro								
MT	B > %lgnm	ns 250 10	000 c3			120012				
EX MT	ecuting f B > endma	rom file: cro	: D:/MTBWI	N\MACROS	\lgnmns.	MAC				
MT	B >	010								
11										-
•										• //
W	/orksheet 1 ^s	***								_ 🗆 ×
+	C1	C2	C3	C4	C5	C6	C7	CE	B C9	C1
	size 10	size 50	size 250							
1	2.55084	1.87330	1.68669							
2	1.17187	1.64760	1.74379							
3	1.59003	1.55393	1.70170							
4	1.75261	1.41372	1.59691							
5	1.39769	1.30627	1.45943							-
										▶ /h
Curre	nt Worksheet:	Worksheet 1								21:08

Histograms of these three sets of means are shown below.



In the study document and in the previous questions, the underlying population was always Normal. If this is the case then distribution of the sample means will also follow a Normal distribution. This is fine if the population is Normal, but what about other forms of population, such as that which is illustrated by the histogram in this question?

This question illustrates an important feature of the distribution of sample means which was not covered in the study document. This is the fact that the distribution of sample means is approximately Normal *whatever the shape of the distribution of the individual observations.* The approximation gets closer as the size of the sample gets bigger, as illustrated by the three histograms above.

A consequence is that methods for statistical inference based on assumptions of Normality have wider applicability than might be thought.

End of solution sheet