Unit 2: Some Basics

On the accuracy of statistical procedures in Microsoft Excel 97

B. D. McCullough^{𝔅, ⋈} and Berry Wilson

Federal Communications Commission, 445 12th St. SW Room 2C-134, Washington, DC 20554, USA

Received 1 June 1998; revised 1 December 1998. Available online 10 August 1999.

Abstract

The reliability of statistical procedures in Excel are assessed in three areas: estimation (both linear and nonlinear); random number generation; and statistical distributions (e.g., for calculating *p*-values). Excel's performance in all three areas is found to be inadequate. Persons desiring to conduct statistical analyses of data are advised not to use Excel.

Author Keywords: DIEHARD; ELV; Numerical accuracy; Software testing; StRD

Example: Hair color at NYS Fair

G		÷				Book1	- Microso	ft Excel						-	= x
	Home Insert	Page Layou	t Formu	las Dat	a Revie	w Vie	w Get	Started						0 -	σx
	J 🎹 🔜				1	-	1	. 0		Α				Ω	
Piv	otTable Table Pictur	e Clip Shape	es SmartArt	Column	Line Pie	Bar	Area Scat	tter Other	Hyperlink	Text H	eader Wor	dArt Signatu	re Object :	Symbol	
	Tables	Illustration	\$			Charts		G	Links	DOX CL	rooter	Text			
	A1 -	• (• f	blond(e	•)											×
4	A	В	С	D	E	F	G	н	1	J	К	L	М	1	N
1	blond(e)	15489													
2	black	17831													
3	brown	27044													
4	red	5281			30000	,									
5	green/blue	255			25000	,									
6	gray	25689			20000	-					-				
7	none	14921			15000	,									
8					10000	, <u> </u>	_	-		-		Series1			
9					5000										
					0		_								
						(e)	- the	un ed	Jue	134	ane				
						blonos	DI. Y	50.	reenlow	\$ 0	0				

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Numerical variables: Graphs of frequencies I



Match Variable to Dotplot

- A. jersey numbers of 2011 Colgate football players
- B. annual snowfall amounts for a sample of U.S. cities
- C. margins of victory in a sample of Major League Baseball games
- D. prices of properties in the Monopoly board game
- E. ages at which a sample of mothers had their first child
- F. scores on a statistics exam
- G. weights of 1999 cars

1.					.: .					
2.	•••		•••	••••			••••	•		•
3.		<u></u>	<u></u>	<u></u> :.	<u></u>	••••	••••	•••	••	•••
4.										
				÷					:	•
5.				::	•		•	•		
6.		:			:.	.:	::::		.:	.: .
7.				:::	•	•	:.:.			

▲ロト ▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ● 臣 ■ ● の Q (2)



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで

Numerical variables: Graphs of frequencies II



Excel's "histograms" aren't, unless bar widths are equal.

Density scale

Weekly salaries in a company: Vertical axis is in % per \$200.

Where is the high point?



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─ のへで

Reading a histogram

Hrs slept by CU students

(Data from questionnaire)

Hrs	Areas	
1,2,3	$3 \times 1 = 3$	3/20 = 15%
4,5	$2 \times 3 = 6$	6/20 = 30%
6	$1 \times 4 = 4$	4/20 = 20%
7	$1 \times 3 = 3$	3/20 = 15%
8,9	$2 \times 2 = 4$	4/20 = 20%
	Sum = 20	

So maybe 5% slept 1hr, 5% 2hr and 5% 3hr;

or maybe 0% slept 1hr, 7% 2hr, 8% 3hr



- Which given interval contains the most students?
- 2. Which 1-hr period contains the most students (i.e., is most "crowded")?
- 3. About what % slept 8 hr?
- About what % slept 3 or 4 hr?
- 5. If there were 240 surveyed, about how many slept 6–7 hr?

Drawing a histogram: Horses' weights in kg



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Frequencies by area?







ъ

Digital ratios spreadsheet

Weekly salaries in a small factory: 30 workers 17 get \$200, 5 \$400, 6 \$500, 1 \$2000, 1 \$4600 avg = \$500, median = \$200



Test for skew

median < avg: skewed to right median > avg: skewed to left (These can be used as the <u>def</u> of left or right skew.)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Changing the list (I)

"Linear" changes of variable (i.e., changing units of measure)

Basic list: 6, 9, 15

$$\mu = \frac{6+9+15}{3} = 10$$

$$\sigma = \sqrt{\frac{(6-10)^2 + (9-10)^2 + (15-10)^2}{3}} = \sqrt{14}$$

Add 5: 11, 14, 20

$$\mu = \frac{11 + 14 + 20}{3} = 15$$

$$\sigma = \sqrt{\frac{(11 - 15)^2 + (14 - 15)^2 + (20 - 15)^2}{3}} = \sqrt{14}$$

Multiply by 12: 72, 108, 180

$$\mu = \frac{72 + 108 + 180}{3} = 120$$

$$\sigma = \sqrt{\frac{(72 - 120)^2 + (108 - 120)^2 + (180 - 120)^2}{3}} = 12\sqrt{14}$$

Changing the list (II):

Adjoining copies of the average: 6, 9, 15, 10, 10, 10, 10

$$\mu = \frac{6+9+15+4(10)}{7} = 10$$

$$\sigma = \sqrt{\frac{(6-10)^2 + (9-10)^2 + (15-10)^2 + 4(10-10)^2}{7}}$$

$$= \sqrt{\frac{3}{7}} \cdot \sqrt{14} < \sqrt{14}$$

Combining two lists:

A: 25, 35 B: 40, 40, 50, 50

 $\mu_A = 30, \qquad \sigma_A = 5, \qquad \qquad \mu_B = 45, \qquad \sigma_B = 5$

A and B: 25, 35, 40, 40, 50, 50

$$\mu = 40, \qquad \sigma = 5\sqrt{3} > 5$$

Box-&-whisker plots

Include min, first quartile, median, third quartile, max — "normal" box-&-whisker plot just shows them all.

	U5	 A =16+1.5 	(H8-H6)	E	E	0	ч				м	N
	Eavpti	an Skull Measure	ements	C	r	0			K	L	M	N
	from A	. Thomson and I	R. Randall-Maci	iver, Ar	ncient Ra	ces of the	Thebaid	(1905)			
3	Year	Max Breadth										
4	-4000	131	MaxBreadth									
5	-4000	125	Year	Max	75th %ile	Median	25th &ile	Min	Hi bd	Hi fence	Lo fence	Lo bd
6	-4000	131	-4000	141	134.75	131	128	119	155	144.88	117.88	107.75
7	-4000	119	-3300	148	134.75	132	130.75	123	147	140.75	124.75	118.75
в	-4000	136	-1850	140	137	136	132.25	126	151	144.13	125.13	118
9	-4000	138	-200	144	138.75	135	132.25	129	158	148.5	122.5	112.75
0	-4000	139	150	147	139	137	132.25	126	159	149.13	122.13	112
1	-4000	125	overall	148	137	134	131	119	155	146	122	113
2	-4000	131										
3	-4000	134	overall	117	0			0		nalas		
4	-4000	129	frequencies	120	1	40		overa	eque			
5	-4000	134		123	1	40				Л		
6	-4000	126		126	11	30						
7	-4000	132		129	11	30						
8	-4000	141		132	35	20						
9	-4000	131		135	31	20						
10	-4000	135		138	38	10						
1	-4000	132		141	15	10						
2	-4000	139		144	4	0				1117		
13	-4000	132		147	2	0	4 0 0	A 0	0.6	~ ~ ~	4	
4	-4000	126		More	1		N & & .	Q Q ,	Sr 19 1	a the the t	N. Mar	
:5	-4000	135										
6	-4000	134							÷			
27	-4000	128			Tow	ard box & v	vhisker plot	8				
8	-4000	130		150 ~			_	_	[
9	-4000	138			0		n		- [
0	-4000	128		145 -		Π						
11	-4000	127		140 -					к [
2	-4000	131		135 -		h-h-	- h- t	7 5	th %ile			
13	-4000	124		1.20	h h	n r		D 0 10	idian			
14	-3300	124		1	n li			= H	n			
15	-3300	133		125 -								
6	-3300	138		120 -				-				
17	-3300	148		115								
8	-3300	126			-4000 -3300	-1850 -20	0 150 over	all				
19	-3300	135										
10	-3300	132										
11	-3300	133										
ő	-3300	131										



æ.

"Modified" plot puts limit on whisker length: 1.5 IQR.

- ▶ 3rd quartile +1.5 IQR = "(inner) fence"
- whisker ends at last value before or on fence
- datum beyond the fence is an "outlier" [which we reject only "for cause" (?)]
- beyond 3rd quartile + 3 IQR ("outer fence" or "bound") is "extreme outlier"



Overall frequencies

Handspans (cm)



◆□ ▶ ◆圖 ▶ ◆ 圖 ▶ ◆ 圖 ▶ ○ 圖 ○