Author: Leonore Findsen, Cheng Li

1. Numerical Summaries

Example 1. Time to Start a Business (Data: eg01-23time24.txt)

An entrepreneur faces many bureaucratic and legal hurdles when starting a new business. The World Bank collects information about starting businesses throughout the world. It has determined the time, in days, to complete all of the procedures required to start a business. Data for 195 countries are included in the data set. For this section we will examine data for a sample of 24 of these countries. Here are the data:

13	66	36	12	8	27	6	7	5	7	52	48
15	7	12	94	28	5	13	60	5	5	18	18

a) Find the mean and the standard deviation for the time to start for the new businesses.

- b) Find the five number summary for the time to start for the new businesses.
- c) Create a histogram for the length of all of the flowers.
- d) Create a boxplot (modified) for the new businesses.

Solution

```
* reading in the data;
data TimeStart;
 infile 'W:\PC-text\eg01-23time24.txt' delimiter = '09'x firstobs = 2;
 input country $ time;
run;
*a) and b);
proc univariate data = TimeStart;
 var time; /*this means that only the variable time will be
             analyzed */
run;
*c);
proc sqplot data = TimeStart;
 histogram time / binwidth = 10 ;
  /*The binwidth is the width of each box in the histogram
   For large data sets, there should be no more than 20 classes
   You will need to change the binwidth to get this number based on
       the range. See the information from proc univariate*/
run ;
*d);
data TimeStart1;
 set TimeStart;
 index = 1;
run:
/*Explanation of what was done above:
We created another data set called 'TimeStart1'.
We set the original data set 'TimeStart' as the input file.
This means that all variables in 'TimeStart' remain in the new data
   set.
```

Author: Leonore Findsen, Cheng Li

Another variable called 'index ' was created and added to 'TimeStart1'. According to the code, every observation of 'index' has value '1'. This was done because the proc boxplot requires two variables, one variable has all of the quantitative values and the other categorical variable that indicates which category each of the values are in. If there is only one category, we need to create the second variable and arbitrarily indicate the label of the one category. */

```
proc boxplot data = TimeStart1;
  plot time * index/boxstyle=schematic idsymbol=circle;
/* This creates a modified boxplot(s) of the score variable for each
  group in the group variable.
Note, if there is only one group it will produce a single boxplot (the
group variable is still required), if there are multiple groups it will
create side-by-side boxplots
  boxstyle=schematic: modified boxplots, that is outliers are points
  idsymbol = circle: the outliers are circles.
  The diamond in the plot is the location of the mean.*/;
run;
```

a) Find the mean and the standard deviation for the time to start for the new businesses.

b) Find the five number summary for the time to start for the new businesses.

Solution

			The	e SAS	System										
		The	UNI \	IVARIA1 /ariable	E Procedur time	e									
				Mome	ents										
N				24 S	um Weights	8		24							
Mea	m		23	8.625 S	um Observa	tions		567			_				
Std Deviation		23	828	7596 V	Variance		567.809783		Quantiles (Definition 5)						
Skowness		1,6	008	1155 K	Kurtosis		2.07559957		Quantile	Estimate					
Uncorrected SS		SS	2	6455 C	orrected SS		13059.625		100% Max	94					
Coeff Variation 100.852474 Std Error Mean			au .	4.8540	2518	99%	94								
	10170	Ba	Basic Statistical Measures						95%	66	;				
	Loc	cation Variability			ity			90%	60		Extre	me O	bservatio		
	Mean Z		3.00000 Vari		ance 567		57.80978			75% Q3		Lowest		Highe	
	Median 13 Mode 5	5.000	5.00000 Rai		ge 8		89.00000		50% Median	13		Value	Obs	Value	C
				Interqu	artile Rang	e 2	5.00000		25% Q1	7		5	22	48	Г
		Tests for Location: Mu0=0							10%	5		5	21	52	t
	Test		Statistic p Value				5%	5		5	18	60	t		
	Stude	nt's t	t	4.8570	088 Pr > t	<	0001		1%	5	-	5	Q	66	+
	Sign	Sign M			12 Pr >=	- M <0	0001			5			-	00	+
	Signed Ra		S	1	150 Pr >- 1	SI <	0001		0% Min	5		6	7	94	

Though I provided all of the tables in this tutorial, you will lose points if you provide more tables than are required to explain your answer.

2

STAT 350: Introduction to Statistics

Department of Statistics, Purdue University, West Lafayette, IN 47907

Author: Leonore Findsen, Cheng Li

Mean = 23.625, Standard deviation = 23.82876 The five number summary: Max = 94, Q_3 = 32, Median = 13, Q_1 = 7, Min = 5

2. Creating Histograms

c) Create a histogram for the length of all of the flowers.

Solution

Remember that the histogram should have the appropriate number of bins for small data sets.

number of binx $\approx \sqrt{number of data points}$

For large data sets, the best number is around 10 – 20 or what 'looks good'.

The following shows the resulting histogram:



I strongly recommend that you change the size of the graph so that it fits better on the page.

3. Boxplots

d) Create a boxplot (modified) for the new businesses.

Solution:

Remember that you need a grouping variable even if there is no categorical variable in the problem.

Author: Leonore Findsen, Cheng Li

