for

An Introduction to Statistical Methods and Data Analysis

6[™] EDITION

R. Lyman Ott

Michael Longnecker

Texas A&M University

Prepared by

Jackie Miller

The Ohio State University

John Draper

The Ohio State University





© 2010 Brooks/Cole, Cengage Learning

ALL RIGHTS RESERVED. No part of this work covered by the copyright herein may be reproduced, transmitted, stored, or used in any form or by any means graphic, electronic, or mechanical, including but not limited to photocopying, recording, scanning, digitizing, taping, Web distribution, information networks, or information storage and retrieval systems, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without the prior written permission of the publisher except as may be permitted by the license terms below.

For product information and technology assistance, contact us at Cengage Learning Customer & Sales Support, 1-800-354-9706

For permission to use material from this text or product, submit all requests online at www.cengage.com/permissions

Further permissions questions can be emailed to permissionrequest@cengage.com

ISBN-13: 978-0-495-10915-0 ISBN-10: 0-495-10915-0

Brooks/Cole

20 Channel Center Street Boston, MA 02210 USA

Cengage Learning is a leading provider of customized learning solutions with office locations around the globe, including Singapore, the United Kingdom, Australia, Mexico, Brazil, and Japan. Locate your local office at: international.cengage.com/region

Cengage Learning products are represented in Canada by Nelson Education, Ltd.

For your course and learning solutions, visit academic.cengage.com

Purchase any of our products at your local college store or at our preferred online store www.ichapters.com

NOTE: UNDER NO CIRCUMSTANCES MAY THIS MATERIAL OR ANY PORTION THEREOF BE SOLD, LICENSED, AUCTIONED, OR OTHERWISE REDISTRIBUTED EXCEPT AS MAY BE PERMITTED BY THE LICENSE TERMS HEREIN.

READ IMPORTANT LICENSE INFORMATION

Dear Professor or Other Supplement Recipient:

Cengage Learning has provided you with this product (the "Supplement") for your review and, to the extent that you adopt the associated textbook for use in connection with your course (the "Course"), you and your students who purchase the textbook may use the Supplement as described below. Cengage Learning has established these use limitations in response to concerns raised by authors, professors, and other users regarding the pedagogical problems stemming from unlimited distribution of Supplements.

Cengage Learning hereby grants you a nontransferable license to use the Supplement in connection with the Course, subject to the following conditions. The Supplement is for your personal, noncommercial use only and may not be reproduced, posted electronically or distributed, except that portions of the Supplement may be provided to your students IN PRINT FORM ONLY in connection with your instruction of the Course, so long as such students are advised that they

may not copy or distribute any portion of the Supplement to any third party. You may not sell, license, auction, or otherwise redistribute the Supplement in any form. We ask that you take reasonable steps to protect the Supplement from unauthorized use, reproduction, or distribution. Your use of the Supplement indicates your acceptance of the conditions set forth in this Agreement. If you do not accept these conditions, you must return the Supplement unused within 30 days of receipt.

All rights (including without limitation, copyrights, patents, and trade secrets) in the Supplement are and will remain the sole and exclusive property of Cengage Learning and/or its licensors. The Supplement is furnished by Cengage Learning on an "as is" basis without any warranties, express or implied. This Agreement will be governed by and construed pursuant to the laws of the State of New York, without regard to such State's conflict of law rules.

Thank you for your assistance in helping to safeguard the integrity of the content contained in this Supplement. We trust you find the Supplement a useful teaching tool.

Table of Contents

CHAPTER 1: Statistics and the Scientific Method	1
CHAPTER 2: Using Surveys and Experimental Studies to Gather Data	3
CHAPTER 3: Data Description	11
CHAPTER 4: Probability and Probability Distributions	51
CHAPTER 5: Inferences about Population Central Values	75
CHAPTER 6: Inferences Comparing Two Population Central Values	97
CHAPTER 7: Inferences about Population Variances	119
CHAPTER 8: Inferences about More Than Two Population Central Values	139
CHAPTER 9: Multiple Comparisons	169
CHAPTER 10:Categorical Data	179
CHAPTER 11:Linear Regression and Correlation	223
CHAPTER 12: Multiple Regression and the General Linear Model	265
CHAPTER 13:Further Regression Topics	303
CHAPTER 14: Analysis of Variance for Completely Randomized Designs	365
CHAPTER 15: Analysis of Variance for Blocked Designs	391
CHAPTER 16:The Analysis of Covariance	411
CHAPTER 17: Analysis of Variance for Some Fixed-, Random-, and Mixed-Effects Models	433
CHAPTER 18:Split-Plot, Repeated Measures, and Crossover Designs	453
CHAPTER 19: Analysis of Variance for Some Unbalanced Designs	481

Chapter 1

Statistics and the Scientific Method

1.1

- **a.** The population of interest is the weight of the shrimp maintained on the specific diet for a period of 6 months.
- **b.** The sample is the 100 shrimp selected from the pond and maintained on the specific diet for a period of 6 months.
- **c.** The weight gain of the shrimp over 6 months.
- **d.** Since the sample is only a small proportion of the whole population, it is necessary to evaluate what the mean weight may be for any other randomly selected 100 shrimps.

1.2

- **a.** The amount of radioactivity at all points in the suspect area.
- **b.** The 200 randomly selected points in the suspect area.
- **c.** The level of radioactivity in the suspect area.
- **d.** We want to relate the level of radioactivity of the 200 points in the sample to the level in the whole suspect area. Thus we need to know how accurate a portrayal of the population is provided by the 200 points in the sample.

1.3

- **a.** All households in the city that receive welfare support.
- **b.** The 400 households selected from the city welfare rolls.
- c. The number of children per household for those households in the city which receive welfare.
- **d.** In order to evaluate how closely the sample of 400 households matches the number of children in all households in the city receiving welfare.

1.4

- **a.** All football helmets produced by the five companies over a given period of time.
- **b.** The 540 helmets selected from the output of the five companies.
- **c.** The amount of shock transmitted to the neck when the helmet's face mask is twisted.
- **d.** The neck strength of players is extremely variable for high school players. Hence, the amount of damage to the neck varies considerably from player to player for exactly the same amount of shock transmitted by the helmet.

- **a.** The population of interest is the population of those who would vote in the 2004 senatorial campaign.
- **b.** The population from which the sample was selected is registered voters in this state.
- **c.** The sample will adequately represent the population, unless there is a difference between registered voters in the state and those who would vote in the 2004 senatorial campaign.
- **d.** The results from a second random sample of 5,000 registered voters will not be exactly the same as the results from the initial sample. Results vary from sample to sample. With either sample we hope that the results will be close to that of the views of the population of interest.

- **a.** The professor's population of interest is college freshmen at his university.
- **b.** The sampled population is all freshmen enrolled in HIST 101.
- **c.** Yes, there is a major difference in the two populations. Those enrolled in HIST 101 may not accurately reflect the population of all freshmen at his university. For example, they might be more interested in history.
- **d.** Had the professor lectured on the American Revolution, those students in HIST 101 would be more likely to know which country controlled the original 13 states prior to the American Revolution than other freshmen at the university.

Chapter 2

Using Surveys and Experimental Studies to Gather Data

2.1

- **a.** The explanatory variable is level of alcohol drinking. One possible confounding variable is smoking. Perhaps those who drink more often also tend to smoke more, which would impact incidence of lung cancer. To eliminate the effect of smoking, we could block the experiment into groups (e.g., nonsmokers, light smokers, heavy smokers).
- **b.** The explanatory variable is obesity. Two confounding variables are hypertension and diabetes. Both hypertension and diabetes contribute to coronary problems. To eliminate the effect of these two confounding variables, we could block the experiment into four groups (e.g., hypertension and diabetes, hypertension but no diabetes, diabetes but no hypertension, neither hypertension nor diabetes).

2.2

- **a.** The explanatory variable is the new blood clot medication. The confounding variable is the year in which patients were admitted to the hospital. Because those admitted to the hospital the previous year were not given the new blood clot medication, we cannot be sure that the medication is working or if something else is going on. We can eliminate the effects of this confounding by randomly assigning stroke patients to the new blood clot medication or a placebo.
- **b.** The explanatory variable is the software program. The confounding variable is whether students choose to stay after school for an hour to use the software on the school's computers. Those students who choose to stay after school to use the software on the school's computers may differ in some way from those students who do not choose to do so, and that difference may relate to their mathematical abilities. To eliminate the effect of the confounding variable, we could randomly assign some students to use the software on the school's computers during class time and the rest to stay in class and learn in a more traditional way.

2.3

Possible confounding factors include student-teacher ratios, expenditures per pupil, previous mathematics preparation, and access to technology in the inner city schools. Adding advanced mathematics courses to inner city schools will not solve the discrepancy between minority students and white students, since there are other factors at work.

2.4

There may be a difference in student-teacher ratios, expenditures per pupil, and previous preparation between the schools that have a foreign language requirement and schools that do not have a foreign language requirement.

The relative merits of the different types of sampling units depends on the availability of a sampling frame for individuals, the desired precision of the estimates from the sample to the population, and the budgetary and time constraints of the project.

2.6

She could conduct a stratified random sample in which the plants serve as the stratum. A simple random sample could then be selected within each plant. This would provide information concerning the differences between the plants along with the individual opinions of the employees.

2.7

The list of registered voters could be used as the sampling frame for selecting the persons to be included in the sample.

2.8

- **a.** No. The survey in which the interviewer showed the peanut butter should be more accurate because it does not rely on the respondent's memory of which brand was purchased.
- **b.** Both surveys may have survey nonresponse bias because an entire segment of the population (those not at home) cannot be contacted. Also, both surveys may have interviewer bias resulting from the way the question was posed (e.g., tone of voice). In the first survey, results may be biased by the respondent's ability to recall correctly which brand was purchased. The second survey may be biased by the respondent's unwillingness to show the interviewer the peanut butter jar (too intrusive), or by the respondent not recognizing that the peanut butter that was purchased was *low fat*.

2.9

- **a.** Alumni (men only?) who graduated from Yale in 1924.
- **b.** No. Alumni whose addresses were on file 25 years later would not necessarily be representative of their class.
- **c.** Alumni who <u>responded</u> to the mail survey would not necessarily be representative of those who were <u>sent</u> the questionnaires. Income figures may not be reported accurately (intentionally), or may be rounded off to the nearest \$5,000, say, in a self-administered questionnaire.
- **d.** Rounding income responses would make the figure \$25,111 unlikely. The fact that higher income respondents would be more likely to respond (bragging), and the fact that incomes are likely to be exaggerated, would tend to make the estimate too high.

- **a.** Simple random sampling.
- **b.** Stratified sampling.
- **c.** Cluster sampling.

- **a.** Simple random sampling.
- **b.** Stratified sampling.
- **c.** Cluster sampling.

2.13

- **a.** Stratified sampling. Stratify by job category and then take a random sample within each job category. Different job categories will use software applications differently, so this sampling strategy will allow us to investigate that.
- **b.** Systematic random sampling. Sample every tenth patient (starting from a randomly selected patient from the first ten patients). Provided that there is no relationship between the type of patient and the order that the patients come into the emergency room, this will give us a representative sample.

2.13

- **a.** Stratified sampling. We should stratify by type of degree and then sample 5% of the alumni within each degree type. This method will allow us to examine the employment status for each degree type and compare among them.
- **b.** Simple random sampling. Once we find 100 containers we will stop. Still it will be difficult to get a completely random sample. However, since we don't know the locations of the containers, it would be difficult to use either a stratified or cluster sample.

2.14

- **a.** Water temperature and Type of hardener
- **b.** Water temperature: 175 °F and 200 °F; Type of hardener: H_1 , H_2 , H_3
- **c.** Manufacturing plants
- d. Plastic pipe
- e. Location on Plastic pipe
- **f.** 2 pipes per treatment
- **g.** 6 treatments:

```
(175 \, {}^{\circ}\text{F}, H_1), (175 \, {}^{\circ}\text{F}, H_2), (175 \, {}^{\circ}\text{F}, H_3), (200 \, {}^{\circ}\text{F}, H_1), (200 \, {}^{\circ}\text{F}, H_2), (200 \, {}^{\circ}\text{F}, H_3)
```

2.15

a.

- Factors: Location in orchard, Location on tree, Time of year
- Factor levels: Location in orchard 8 sections; Location on tree top, middle, bottom; Time of year October, November, December, January, February, March, April, May
- Blocks: none
- Experimental units: Location on tree during one of the 8 months
- Measurement units: oranges
- Replications: For each section, time of year, and location on tree, there is one experimental unit, hence 1 replication.
- Treatments: 192 combinations of 8 sections, 8 months, and 3 locations on tree –
- (S_i, M_i, L_k) , for i = 1,...,8; j = 1,...,8; k = 1,2,3

b.

• Factors: Type of treatment

• Factor levels: T_1 , T_2

• Blocks: Hospitals

• Experimental units: Wards

• Measurement units: Patients

• Replications: 2 wards per treatment in each of the 8 hospitals

• Treatments: T_1 , T_2

c.

• Factors: Type of treatment

• Factor levels: T_1 , T_2

• Blocks: Hospitals, Wards

• Experimental units: Patients

• Measurement units: Patients

• Replications: 2 patients per treatment in each of the ward/hospital combinations

• Treatments: T_1 , T_2

d.

• Factors: Type of school

• Factor levels: Public; Private – non-parochial; Parochial

• Blocks: Geographic region

• Experimental units: Classrooms

• Measurement units: Students in classrooms

• Replications: 2 classrooms per each type of school in each of the city/region combinations

• Treatments: Public; Private – non-parochial; Parochial

2.16

a. Factors: Temperature, Type of seafood

b. Factor levels: Temperature (0 °C, 5 °C, 10 °C); Type of seafood (oysters, mussels)

c. Blocks: None

d. Experimental units: Package of seafood

e. Measurement units: Sample from package

f. Replications: 3 packages per temperature

g. Treatments: (0 °C, oysters), (5 °C, oysters), (10 °C, oysters), (0 °C, mussels), (5 °C, mussels), (10 °C, mussels)

- **a.** Randomized complete block design with blocking variable (5 farms) and 48 treatments in a 3 \times 4 \times 4 factorial structure.
- **b.** Completely randomized design with 10 treatments (software packages) and 3 replications of each treatment.
- **c.** Latin square design with blocking variables (position in kiln, day), each having 8 levels. The treatment structure is a 2×4 factorial structure (type of glaze, thickness).

- a. Design B. The experimental units are not homogeneous since one group of consumers gives uniformly low scores and another group gives uniformly high scores, no matter what recipe is used. Using design A, it is possible to have a group of consumers that gives mostly low scores randomly assigned to a particular recipe. This would bias this particular recipe. Using design B, the experimental error would be reduced since each consumer would evaluate each recipe. That is, each consumer is a block and each of the treatments (recipes) is observed in each block. This results in having each recipe subjected to consumers who give low scores and to consumers who give high scores.
- **b.** This would not be a problem for either design. In design A, each of the remaining 4 recipes would still be observed by 20 consumers. In design B, each consumer would still evaluate each of the 4 remaining recipes.

2.19

- **a.** "Employee" should refer to anyone who is eligible for *sick days*.
- **b.** Use payroll records. Stratify by employee categories (full-time, part-time, etc.), employment location (plant, city, etc.), or other relevant subgroup categories. Consider systematic selection within categories.
- c. Sex (women more likely to be care givers), age (younger workers less likely to have elderly relatives), whether or not they care for elderly relatives now or anticipate doing so in the near future, how many hours of care they (would) provide (to define "substantial"), etc. The company might want to explore alternative work arrangements, such as flex-time, offering employees 4 ten-hour days, cutting back to 3/4-time to allow more time to care for relatives, etc., or other options that might be mutually beneficial and provide alternatives to taking sick days.

- **a.** Each state agency and some federal agencies have records of licensed physicians, professional corporations, facility licenses, etc. Professional organizations such as the American Medical Association, American Hospital Administrators Association, etc., may have such lists, but they may not be as complete as licensing records.
- **b.** What nursing specialties are available at this time at the physician's offices or medical facilities? What medical specialties/facilities do they anticipate adding or expanding? What staffing requirements are unfilled at this time or may become available when expansion occurs? What is the growth/expansion time frame?
- **c.** Licensing boards may have this information. Many professional organizations have special categories for members who are unemployed, retired, working in fields not directly related to nursing, students who are continuing their education, etc.
- **d.** Population growth estimates may be available from the Census Bureau, university economic growth research, bank research studies (prevailing and anticipated load patterns), etc. Health risk factors and location information would be available from state health departments, the EPA, epidemiological studies, etc.
- **e.** Licensing information should be stratified by facility type, size, physician's specialty, etc., prior to sampling.

If phosphorous first: [P,N]

[10,40], [10,50], [10,60], then [20,60], [30,60]	or
[20,40], [20,50], [20,60], then [10,60], [30,60]	or
[30,40], [30,50], [30,60], then [10,60], [10,60]	

If nitrogen first: [N,P]

[40,10], [40,20], [40,30], then [50,30], [60,30]	or
[50,10], [50,20], [50,30], then [40,30], [60,30]	or
[60,10], [60,20], [60,30], then [40,30], [50,30]	

2.22

	Factor 2				
Factor 1	I	II	III		
A	25	45	65		
В	10	30	50		

2.23

a. Group dogs by sex and age:

Group	Dog
Young female	2, 7, 13, 14
Young male	3, 5, 6, 16
Old female	1, 9, 10, 11
Old male	4, 8, 12, 15

b. Generate a random permutation of the numbers 1 to 16:

Go through the list and the first two numbers that appear in each of the four groups receive treatment L_1 and the other two receive treatment L_2 .

Group	Dog-Treatment
Young female	$2-L_2$, $7-L_1$, 13 , $14-L_2$
Young male	$3-L_1$, $5-L_2$, $6-L_1$, $16-L_2$
Old female	$1-L_1$, $9-L_2$, $10-L_2$, $11-L_1$
Old male	$4-L_1$, $8-L_2$, $12-L_2$, $15-L_1$

2.24

a. Bake one cake from each recipe in the oven at the same time. Repeat this procedure *r* times. The baking period is a block with the four treatments (recipes) appearing once in each block. The four recipes should be randomly assigned to the four positions, one cake per position. Repeat this procedure *r* times.

b. If position in the oven is important, then position in the oven is a second blocking factor along with the baking period. Thus, we have a Latin square design. To have r = 4, we would need to have each recipe appear in each position exactly once within each of four baking periods. For example:

Period 1		Period 2		Period 3		Period 4	
R_1	R_2	R_4	R_1	R_3	R_4	R_2	R_3
R_3	R_4	R_2	R_3	R_{1}	R_2	R_4	$R_{_{1}}$

c. We now have an incompleteness in the blocking variable period since only four of the five recipes can be observed in each period. In order to achieve some level of balance in the design, we need to select enough periods in order that each recipe appears the same number of times in each period and the same total number of times in the complete experiment. For example, suppose we wanted to observe each recipe r = 4 times in the experiment. If would be necessary to have 5 periods in order to observe each recipe 4 times in each of the 4 positions with exactly 4 recipes observed in each of the 5 periods.

Period 1		Period	. 2	Period	13	Period	14	Period	15
R_1	R_2	R_5	$R_{\scriptscriptstyle 1}$	R_4	R_5	R_3	R_4	R_2	R_3
R_3	R_4	R_2	R_3	R_1	R_2	R_5	$R_{_1}$	R_4	R_5

2.25

Discussion question; answers will vary.

2.26

Discussion question; answers will vary.

2.27

Discussion question; answers will vary.

2.28

Discussion question; answers will vary.

2.29

Discussion question; answers will vary.

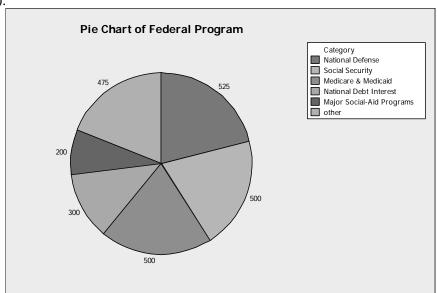
10

Chapter 3

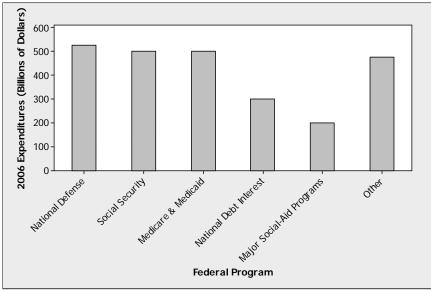
Data Description

3.1

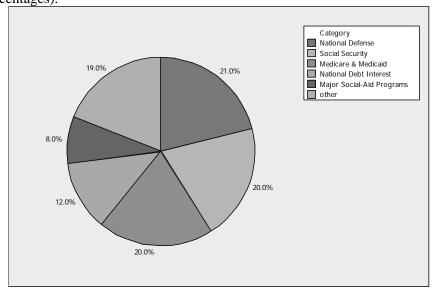
a. The following is a pie chart of the federal expenditures for the 2006 fiscal year (in billions of dollars).

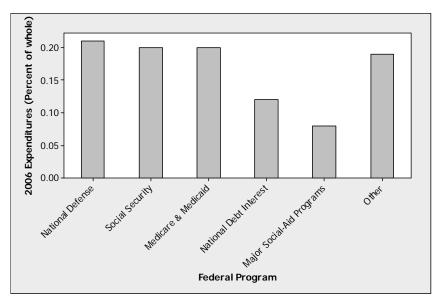


b. The following is a bar chart of the federal expenditures for the 2006 fiscal year (in billions of dollars).



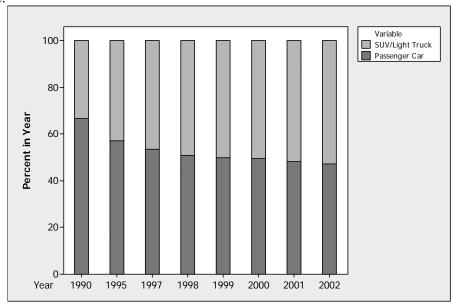
c. The following are a pie chart and bar chart of the federal expenditures for the 2006 fiscal year (in percentages).



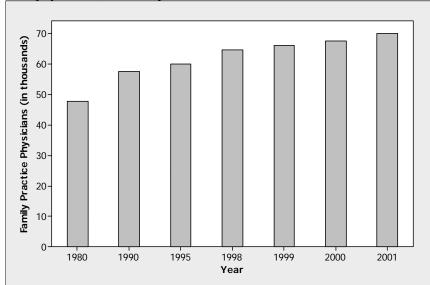


d. The pie chart using percentages is probably most informative to the tax-paying public. Here the tax-paying public can compare the percentages spent by the Federal government for domestic and defense programs as part of a whole.

- 3.2
 - **a.** Pie charts would not be appropriate to display these data. We would not be able to see trends over time.
 - **b.** The following bar chart shows the changes across the 12 years in the public's choice in vehicle.



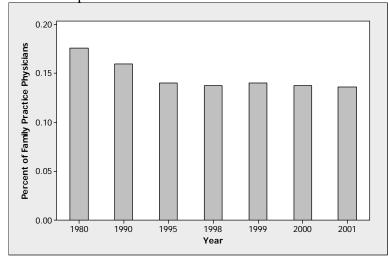
- **c.** It appears that the percentage of passenger cars has decreased over the period 1990-2002. If there was a substantial increase in gasoline prices, we would expect the percentage of passenger cars to increase.
- 3.3
- a. The following bar chart shows the increase in the number of family practice physicians (in thousands of physicians) over the period 1980-2001.



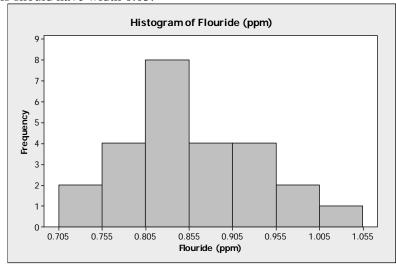
b. The percent of office-based physicians who are family practice physicians over the period 1980-2001 can be seen in the following table.

	1980	1990	1995	1998	1999	2000	2001
Percent Family Practice	17.6	16.0	14.0	13.8	14.0	13.8	13.6

The following bar chart shows the percent of office-based physicians who are family practice physicians over the period 1980-2001.

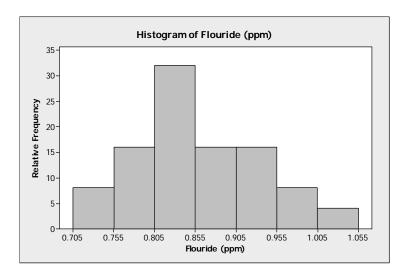


- **c.** While the number of family practice physicians increased over the period 1980-2001, the percent of total office-based physicians who are family practice physicians decreased over the same period.
- 3.4
- **a.** Range = 1.05 0.72 = 0.33
- **b.** The frequency histogram should be plotted with 7 classes ranging from 0.705 to 1.055. The intervals should have width 0.05.



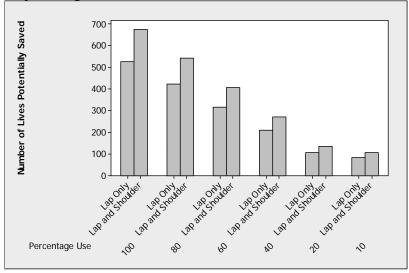
c.	Relative freq	uencies are	given	below.	Plot re	lative	frequencies	versus class	intervals.

Class	Class Interval	Frequency (f_i)	Relative Frequency ($f_i/25$)
1	0.705-0.755	2	0.08
2	0.755-0.805	4	0.16
3	0.805-0.855	8	0.32
4	0.855-0.905	4	0.16
5	0.905-0.955	4	0.16
6	0.955-1.005	2	0.08
7	1.005-1.055	1	0.04
Total		n = 25	1.00



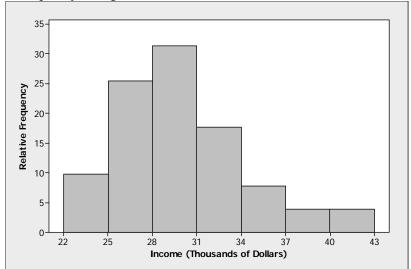
d. The probability is 7/25 = 0.28 that the fluoride reading would be greater than 0.90 ppm. Thus, we would predict that 28% of the days would have a reading greater than 0.90 ppm.

Two separate bar graphs could be plotted, one with Lap Belt Only and the other with Lap and Shoulder Belt. A single bar graph with the Lap Belt Only value plotted next to the Lap and Shoulder Belt for each value of Percentage of Use is probably the most effective plot. This plot would clearly demonstrate that the increase in the number of lives saved by using a shoulder belt increased considerably as the percentage use increased.



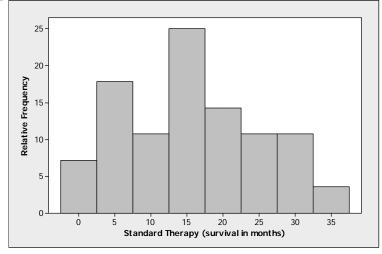
3.6

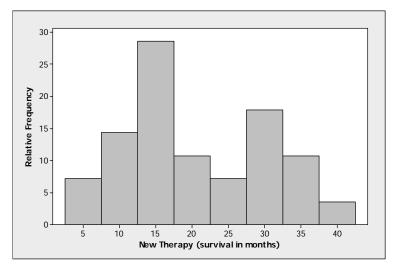
A relative frequency histogram for the income data is:



- **b.** The histogram is unimodal and skewed to the right. The median is in the \$31,000-\$33,900 bin. There are no outliers.
- Since there is a large spread in incomes (about \$21,000), the data are not homogeneous across the states.

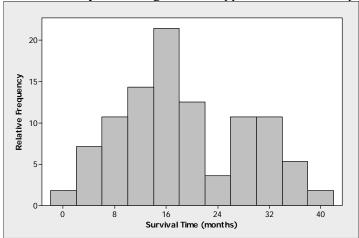
a. The separate relative frequency histograms for the two treatments appear below.





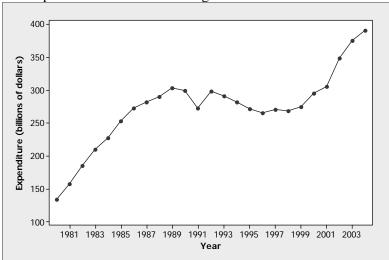
b. The histogram for the New Therapy begins and ends with bins that are slightly higher than the bins in the histogram for the Standard Therapy. This would indicate that the New Therapy generates a few more large values than the Standard Therapy. However, there is not convincing evidence that the New Therapy generates a longer survival time.

The following histogram of the data from the two therapies combined is bimodal and skewed to the right. Because of the bimodality, the histogram does appear to show two separate populations

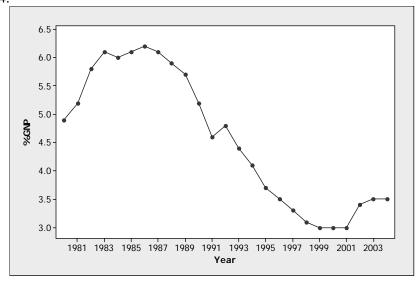


3.9

a. The time series plot shows an increase from 1980 through 1990, with a large dip at 1991. There is then a decrease in expenditures in billions of dollars over the period 1992 to 1999, and then a sharp increase from 1999 through 2004.



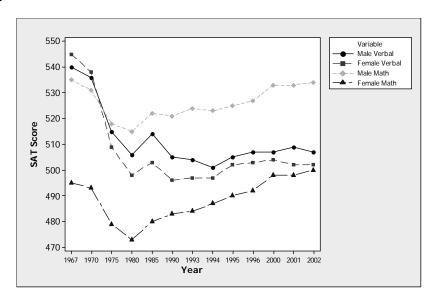
b. The time series plot of expenditures as a percent of GNP shows an increase from 1980 to 1983. It is then fairly steady from 1983 to 1987, decreases from 1987 to 2001 (with the exception of a spike in 1992), has a sharp increase 2001 to 2002, and then is fairly steady 2002 to 2004.



c. The plots do not show similar trends. The time series plot of expenditures supports the assertions of these members of Congress.

3.10

a. Both male and female math SAT scores decreased from 1967 to 1980 and then increased from 1980 to 2002, with the females increasing at a higher rate than the males. Both male and female verbal SAT scores decreased steeply from 1967 to 1980 and then remained fairly steady from 1980 to 2002.

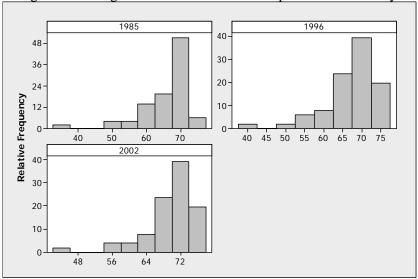


b. The gap between male and female math SAT scores does not appear to have changed over this period.

c. The gap between male and female verbal SAT scores has not changed much over this period, although there was a period (1994 to 2000) where the gap was smaller.

3.11

a. The following shows histograms for the homeownership data for all three years.



- **b.** All three histograms are unimodal and skewed to the left. It appears from the histograms that homeownership has increased over the years.
- **c.** Perhaps more people are able to afford to live in their own homes instead of renting homes from others.
- **d.** Congress could address the issues that a higher proportion of homes are owner-occupied. Congress could then develop tax laws for owner-occupied homes.

3.12 The following are stem-and-leaf plots for each of the three years 1985, 1996, and 2002.

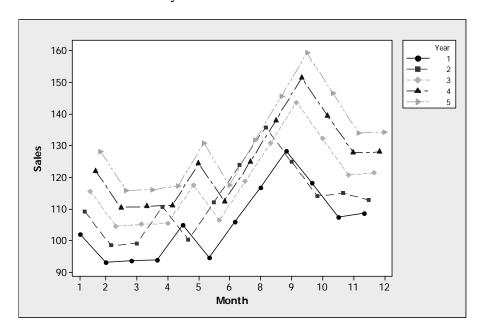
	1985		1996		2002
3	7	3		3	
4		4	0	4	4
4		4		4	
5	014	5	02	5	
5	7	5	56	5	5789
6	00011122334	6	1112233444	6	23
6	5566677777888888999999	6	5666777788888888999	6	55677778899999
7	000000111233	7	000111122233344	7	000000011122222333334444
7	5	7	56	7	556777

 $3 \mid 7 = 37\%$ (stems are tens, leaves are ones)

All three homeownership distributions are unimodal and skewed to the left (since they are skewed, they are not symmetric). The homeownership distributions each have outliers (37% in 1985, 40% in 1996, and 44% in 2002).

3.14

The plots show an upward trend from year 1 to year 5. There is a strong seasonal (cyclic) effect; the number of units sold increases dramatically in the late summer and fall months.



3.15

The mean is $\overline{y} = \frac{\sum_{i} y_{i}}{n} = \frac{55 + 85 + \dots + 31}{20} = \frac{1282}{20} = 64.1$. The median is the average of the 10th and 11th values when arranged in increasing order: median = 67.5. The mode is 90.

3.16

The new mean is $\overline{y} = \frac{1869}{20} = 93.45$. The median and mode remain unchanged at 67.5 and 90, respectively.

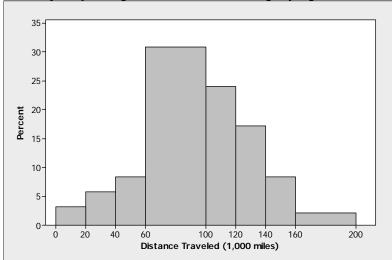
3.17

The 10% trimmed mean is 52.2 for both data sets, since the lowest two and highest two values are deleted before calculating the mean. The 5% trimmed mean is different for both data sets (58.1 for the original data set and 69.85 for the altered data set) since the extreme value of 345 impacts the mean of the altered data set.

The mean is 9.97, the median is 9.5, and the mode is 9.5.

3.19

a. The relative frequency histogram is unimodal and slightly right skewed.



b. The following table is used to calculate the summary statistics:

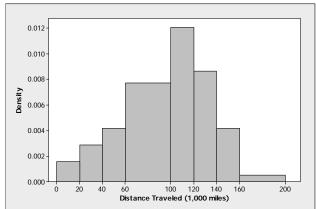
Class Interval	Frequency (f_i)	Midpoint (y_i)	$f_i y_i$
0-20.0	6	10	60
20.1-40.0	11	30	330
40.1-60.0	16	50	800
60.1-100.0	59	80	4720
100.1-120.0	46	110	5060
120.1-140.0	33	130	4290
140.1-160.0	16	150	2400
160.1-200.0	4	180	720
Total	191		18,380

Mean $\approx 18380/191 = 96.2$

Mode \approx 80

Median
$$\approx L + \frac{w}{f_m}(0.5n - cf_b) = 100.1 + \frac{20}{46}[(0.5)(191) - 92] = 101.6$$

c. Since the median is larger than the mean, it would indicate that the plot is somewhat left skewed. This contradiction between what is indicated in the relative frequency histogram and what is indicated by the summary statistics is due to the fact that the class intervals are a different width. The correct plot would have relative frequency divided by class width on the vertical axis. This would then produce a left skewed histogram with mode at approximately 110.



d. The median is more informative since the distribution is somewhat skewed to the left which produces a mean somewhat less than the middle of the distribution. The median distance traveled would at least represent a value such that half of the buses traveled less and half greater than 101,600 miles.

3.20

 a. The mean cannot be approximated since we do not know the endpoint for the last interval, hence cannot compute the midpoint for this interval.
 Mode ≈ 240.

median
$$\approx L + \frac{w}{f_{m}}(0.5n - cf_{b}) = 200 + \frac{19.9}{88}[(0.5)(408) - 119] = 219.2$$

b. The median since it would give an indication of the cholesterol level for which half of the men in this group have a greater or lesser value.

3.21

- **a.** Mean = 8.04, Median = 1.54
- **b.** Terrestrial: Mean = 15.01, Median = 6.03

Aquatic: Mean = 0.38, Median = 0.375

- **c.** The mean is more sensitive to extreme values than is the median.
- **d.** Terrestrial: Median, because the two large values (76.50 and 41.70) result in a mean that is larger than 82% of the values in the data set.

Aquatic: Mean or median since the data set is relatively symmetric.

- **a.** After removing the survival times of the two individuals who left the study, we obtain Mean = 35.22 days. The median can be calculated for all 11 patients, since we know that the values for the two individuals who left the study were greater than the list values 57 and 60 which would place them in the upper half of the survival times. Thus, the Median = 29 days.
- **b.** The median would be unchanged but the mean would increase since these two values will be greater than the mean calculated from the nine observed values.

3.23

- a. If we use all 14 failure times, we obtain Mean > 173.7 days and Median = 154 days. In fact, we know that the mean is greater than 173.7 days since the failure times for two of the engines are greater than the reported times of 300 days.
- **b.** The median would be unchanged if we replace the failure times of 300 with the true failure times for the two engines that did not fail. However, the mean would be increased.

3.24

a. The values are given below:

Group	Mean	Median	Mode
I	2.923	2.805	no mode
II	1.592	1.565	1.55, 1.57
III	0.797	0.755	0.70

- **b.** Mean = 1.7707; Median = 1.565; Modes = 0.70, 1.55, 1.57
- c. If we were to use one summary for the combined group, then the median would be most appropriate because the three groups are substantially different. If separate summaries are computer for each group, then the mean and median are both appropriate since the three groups have relatively symmetric distributions.

3.25

Mean = 1.7707, Median = 1.7083, Mode = 1.273

The average of the three net group means and the mean of the complete set of measurements are the same. This will be true whenever the groups have the same number of measurements, but it is not true if the groups have different sample sizes. However, the average of the group medians and modes are different from the overall median and mode.

- **a.** $\overline{y} = \frac{\sum y_i}{n} = \frac{40}{8} = 5$ years. This value does appear to adequately represent the data set.
- **b.** $\sum_{i=1}^{8} (y_i 5)^2 = (6 5)^2 + (3 5)^2 + (10 5)^2 + (4 5)^2 + (4 5)^2 + (2 5)^2 + (4 5)^2 + (7 5)^2$ = 1 + 4 + 25 + 1 + 1 + 9 + 1 + 4 = 46
- c. $s^2 = \frac{46}{8-1} = 6.57 \Rightarrow s = 2.56$ years. The $CV = 100 \frac{s}{y} \% = 100 \frac{2.56}{5} \% = 51\%$. The standard deviation is 51% of the mean.

- **a.** s = 7.95 years
- **b.** Because the magnitude of the racers' ages is larger than that of their experience.

3.28

a. Racers' age:
$$CV = 100 \frac{s}{y}\% = 100 \frac{7.95}{29.875}\% = 26.6\%$$

Years of experience:
$$CV = 100 \frac{s}{y} \% = 100 \frac{2.56}{5} \% = 51\%$$

The CV for the racers' age is about half of the CV for their years of experience.

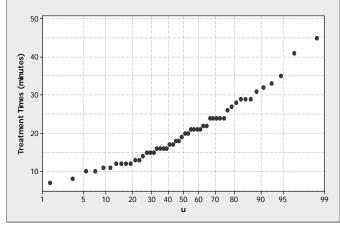
b. Estimated SD for racers' age: (39-18)/4 = 5.25

Estimated SD for years of experience: (10-2)/4=2

The estimate standard deviation for the years of experience is off by about 0.5 years, while the estimated standard deviation for the racers' ages is off by about 2.7 years.

3.29

The quantile plot is given below.

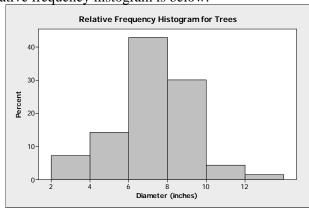


- **a.** The 25th percentile is the value associated with u = 0.25 on the graph, which is 14 minutes. Also, by definition 14 minutes is the 25th percentile since 25% of the times are less than or equal to 14 minute and 75% of the times are greater than or equal to 14 minutes.
- **b.** Yes; the 90th percentile is 31.5 minutes. This means that 90% of the patients have a treatment time less than or equal to 31.5 minutes (which is less than 40 minutes).

a. The frequency table is given here.

Class Intervals	Frequency	Relative Frequency
2.0-4.0	5	0.0714
4.1-6.0	10	0.1429
6.1-8.0	30	0.4286
8.1-10	21	0.3000
10.1-12.0	3	0.0429
12.1-14.0	1	0.0142
Total	70	1.0000

The relative frequency histogram is below:



- **b.** $\overline{y} = 541/70 = 7.7286$
- **c.** s = 1.985

 $\overline{y} \pm s$ yields (5.744, 9.713); 50 / 70 = 71.43%

 $\overline{y} \pm 2s$ yields (3.759, 11.698); 68/70 = 95.71%

 $\overline{y} \pm 3s$ yields (1.774, 13.683); 70/70 = 100%

The percentages are very close to the percentages given by the Empirical Rule: 68%, 95%, and 99.7%.

3.31

a. Luxury: $\overline{y} = 145.0$, s = 27.6

Budget: $\overline{y} = 46.1$, s = 5.13

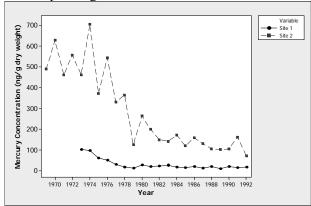
b. Luxury: CV = 19%

Budget: CV = 11%

c. Luxury hotels vary in quality, location, and price, whereas budget hotels are more competitive for the low-end market so prices tend to be similar.

d. The CV would be better because it takes into account the larger difference in the means between the two types of hotels.

a. The time series plot is given here.



For Site 1, there is a steady decrease until 1980, after which the level is fairly constant but at a much lower level than the values for Site 2. The concentrations at Site 2 are very erratic from 1969 to 1980, with alternating rises and falls. From 1980 through 1992, there is a fairly steady decline in mercury concentration.

- **b.** Site 1: Median = 18.25 ng/g; Mean = 29.18 ng/g
 - Site 2: Median = 184.1 ng/g; Mean = 287.1 ng/g

Both distributions are right skewed, thus the median is a more appropriate measure of center than is the mean. Site 1 has a considerably lower center than that of Site 2.

- **c.** Site 1: s = 26.95 ng/g; CV = 92%
 - Site 2: s = 194.7 ng/g; CV = 68%

Comparing the standard deviations can be misleading because 26.95 is smaller in magnitude than 194.7. However, the data values for Site 1 are considerably smaller in magnitude as well. Therefore, it is more informative to compare the CV values. Based on CV values, the concentrations from Site 1 are relatively more variable than those from Site 2.

d. No, Site 1 does not have values for these years.

3.33

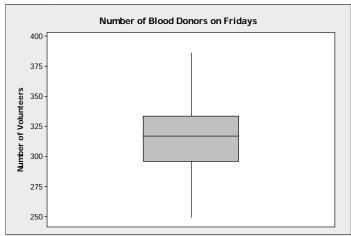
Median = 15;
$$Q_1 = 10$$
; $Q_3 = 21$

3.34

a. A stem-and-leaf plot is below:

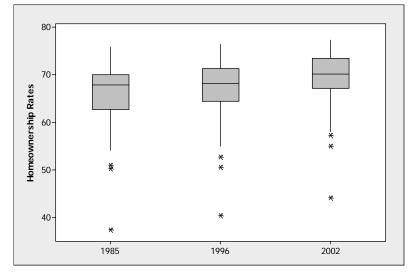
Stem	Leaf
2	5
2	677
2	9
3	00111
3	223333
3	5
3 3 3	67
3	8

b. Min = 250, $Q_1 = (295 + 301) / 2 = 298$, $Q_2 = \text{Median} = (315 + 320) / 2 = 317.5$, $Q_3 = (334 + 334) / 2 = 334$, Max = 386



There are no outliers because $Q_1 - (1.5)IQR = 244 < 250$ and $Q_3 + (1.5)IQR = 388 > 386$. The distribution is approximately symmetric with no outliers.

- **a.** CAN: $Q_1 \approx 1.45$, $Q_2 = \text{Median} \approx 1.65$, $Q_3 \approx 2.4$
 - DRY: $Q_1 \approx 0.55$, $Q_2 = \text{Median} \approx 0.60$, $Q_3 \approx 0.70$
- **b.** Canned dog food is more expensive (median much greater than that for dry dog food), highly skewed to the right with a few large outliers. Dry dog food is slightly left skewed with a considerably less degree of variability than canned dog food.
- **3.36** The following shows comparative boxplots for homeownership rates for the years 1985, 1996, and 2002:



- **a.** The distributions for each of the three years (1985, 1996, and 2002) are each left skewed with a few low outliers. The medians appear to be greater, and the distributions as a whole tend to be higher, in subsequent years.
- **b.** The descriptions in part (a) agree with our description of the distributions in Exercise 3.11.

a. 1985: mean = 65.876; median = 67.90 1996: mean = 66.843; median = 68.20 2002: mean = 69.449; median = 70.20

Since the distributions are left skewed, it is better to use the median for each of the three years.

b. 1985: s = 6.734; CV = 10.2% 1996: s = 6.688; CV = 10.0% 2002: s = 6.163; CV = 8.9%

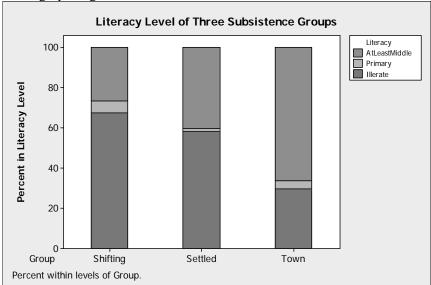
The coefficients of variation are decreasing over the three years.

3.38

- **a.** See the boxplot in Exercise 3.36. Median homeownership rate is increasing over the three years.
- **b.** See the boxplot in Exercise 3.36. The variation in homeownership rate is decreasing over the three years.
- **c.** The District of Columbia, Hawaii, and New York have extremely low homeownership rates in each of the three years. These states are indicated as outliers in the side-by-side boxplot.
- **d.** No states have extremely high homeownership rates in each of the three years.

3.39

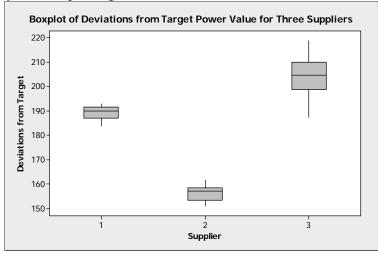
a. A stacked bar graph is given below.



a. The means and standard deviations are given below:

Supplier	\overline{y}	S
1	189.23	2.96
2	156.28	3.30
3	203.94	8.98

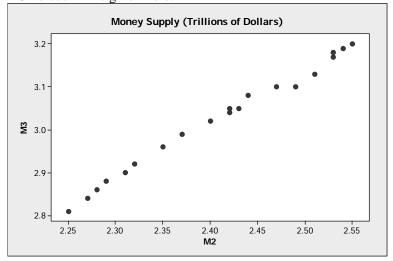
b. A side-by-side boxplot is given below:



The three distributions are relatively symmetric, but supplier 3 is considerably more variable and is shifted about supplier 1's values, which in turn are shifted above supplier 2's values.

- **c.** Supplier 3 not only has the largest mean but also the largest standard deviation. Suppliers 1 and 2 have similar degrees of variability, but supplier 1 has a greater mean than supplier 2.
- **d.** Supplier 2 because it has the smallest mean and deviated with essentially the same degree of variability as supplier 1.

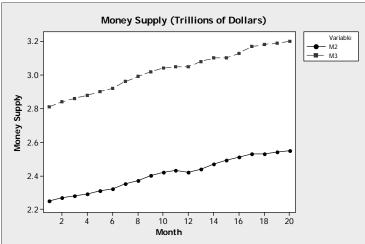
3.41 A scatterplot of M3 versus M2 is given here.



- **a.** Yes, it would since we want to determine the relative changes in the two over the 20-month period of time.
- **b.** See scatterplot. The two measures follow an increasing, approximately linear relationship.

3.42

A time series plot with M2 and M3 values on the vertical axis and months on the horizontal axis is given here.



This is a more informative plot than the scatterplot because it shows the relative changes of the two measures of money supply across the 20 months.

- **a.** Mean = 57.5; Median = 34.0
- **b.** Median since the data has a few very large values which results in the mean being larger than all but a few of the data values.

- **c.** Range = 273; s = 70.2
- **d.** Using the approximation, $s \approx \text{range} / 4 = 273 / 4 = 68.3$. The approximation is fairly accurate.
- **e.** $\overline{y} \pm s \Rightarrow (-12.7,127.7)$; yields 82% $\overline{y} \pm 2s \Rightarrow (-82.9,197.0)$; yields 94% $\overline{y} \pm 3s \Rightarrow (-153.1,268.1)$; yields 97%
- **f.** The Empirical Rule applies to data sets with roughly a "mound-shaped" histogram. The distribution of this data set is highly skewed right.

- **a.** Mode = 5 hours; Median = 15 hours; Mean = 15.96 hours
- **b.** Range = 34 4 = 30 hours; $s \approx 30 / 4 = 7.5$ hours
- c. s = 8.5 hours
- **d.** No, the histogram for the data set is skewed to the right and hence is not mound-shaped.

3.45

- **a.** Price per roll: Mean = 0.9196, s = 0.4233Price per sheet: Mean = 0.01091, s = 0.0059
- **b.** Price per roll: $CV = 100 \frac{0.4233}{0.9196} \% = 46.03\%$

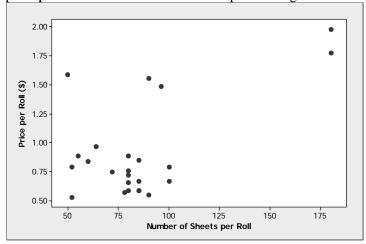
Price per sheet:
$$CV = 100 \frac{0.0059}{0.01091} \% = 54.13\%$$

The price per sheet is more variable relative to its mean.

c. CV; The CV is unit free, whereas the standard deviation also reflects the relative magnitude of the data values.

3.46

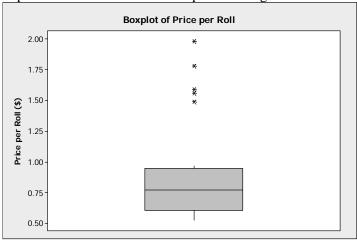
A scatterplot of price per roll versus number of sheets per roll is given here.

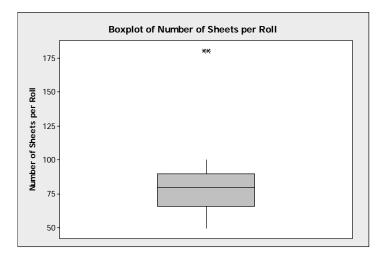


a. No.

- **b.** No, as the number of sheets increases from 50 to 100, there is just a scatter of points, no real pattern. The price per roll jumps dramatically for the two brands having the largest number of sheets.
- c. Paper towel sheets vary in thickness and size, both of which will affect the price.

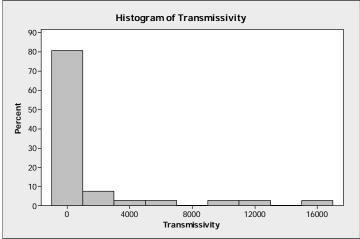
3.47 Boxplots for price per roll and number of sheets per roll are given here.





From the two boxplots, there a 5 unusual brands with regard to price per roll: \$1.49, \$1.56, \$1.59, \$1.78, and \$1.98. There are 2 unusual brands with respect to number of sheets per roll: 180 and 180.

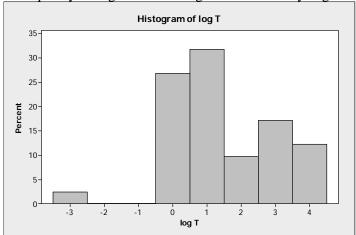
a. A relative frequency histogram for transmissivity is given here.



- **b.** The distribution is unimodal and highly skewed to the right.
- c. $\overline{y} \pm s = 1424 \pm 3488 \Rightarrow (-2063, 4912)$ contains 34/71 = 90.2% $\overline{y} \pm 2s = 1424 \pm (2)3488 \Rightarrow (-5551, 8400)$ contains 38/41 = 92.7% $\overline{y} \pm 3s = 1424 \pm (3)3488 \Rightarrow (-9039, 11888)$ contains 39/41 = 95.1%

These values do not match the values from the Empirical Rule: 68%, 95%, and 99.7%.

d. A relative frequency histogram for the log of transmissivity is given here.



The shape is more mound-shaped than the original data, although it appears somewhat skewed to the right with a low outlier.

 $\overline{y} \pm s = 1.48 \pm 1.54 \Rightarrow (-0.06, 3.02)$ contains 31/41 = 75.6%

 $\overline{y} \pm 2s = 1.48 \pm (2)1.54 \Rightarrow (-1.60, 4.56)$ contains 40/41 = 97.6%

 $\overline{y} \pm 3s = 1.48 \pm (3)1.54 \Rightarrow (-3.14, 6.10)$ contains 41/41 = 100%

These values more closely match the percentages from the Empirical Rule.

3.49

A relative frequency histogram for murder rate is given here.



3.50

a. Mode = 2.5; Median
$$\approx L + \frac{w}{f_m} [0.5n - cf_b] = 5.5 + \frac{2}{13} [0.5(90) - 35] = 7.04$$

b. Mean
$$\approx \frac{1}{n} \sum_{i=1}^{13} f_i y_i = 747 / 90 = 8.3$$

c. Since the distribution is skewed to the right, the median provides a better measure of the center of the distribution.

a. 75th percentile
$$\approx L + \frac{w}{f_m} [0.75n - cf_b] = 13.5 + \frac{2}{9} [0.75(90) - 72] = 12.5$$

25th percentile $\approx L + \frac{w}{f_m} [0.25n - cf_b] = 3.5 + \frac{2}{15} [0.25(90) - 20] = 3.83$
IQR $\approx 12.5 - 3.83 = 8.67$

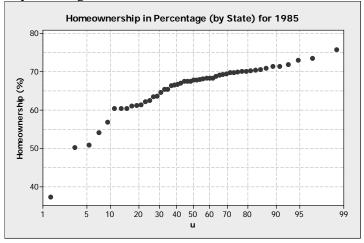
b.
$$s^2 \approx \frac{1}{n-1} \sum_{i=1}^{13} f_i (y_i - 8.3)^2$$

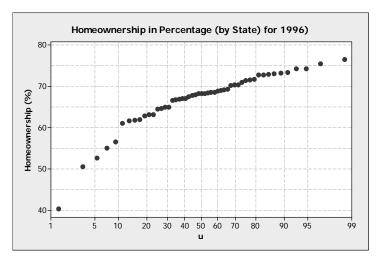
= $\frac{1}{89} \Big[2(0.5 - 8.3)^2 + 18(2.5 - 8.3)^2 + \dots + 1(24.5 - 8.3)^2 \Big] = 29.0382$
Thus, $s \approx \sqrt{29.0382} = 5.389$.

3.52

The years referred to in the questions are 1985 and 1996, so the quantile plots for homeownership

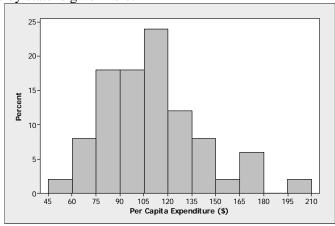
rates for those two years are given here.





- The 20th percentile for 1996 is given by reading the vertical value on the graph for $u = 0.20 \Rightarrow 63\%$. Thus, approximately 20% of the 1996 homeownership percentages are less than or equal to 63%.
- **b.** The upper 10th percentile would correspond to the states having the 5 largest percentages: Michigan, Indiana, West Virginia, Minnesota, and Maine.
- In 1985, the states falling in the upper 10th percentile are Pennsylvania, South Carolina, Wyoming, Maine, and West Virginia. There are only two states that fall into both groups.

a. A relative frequency histogram of per capita expenditures (dollars) for health and hospital services by state is given here.



b. Mean $\approx \frac{1}{n} \sum_{i=1}^{11} f_i y_i = 5480 / 50 = 109.6$

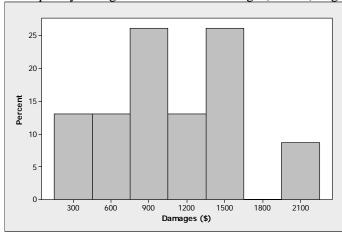
$$s^2 \approx \frac{1}{n-1} \sum_{i=1}^{11} f_i (y_i - 109.6)^2$$

$$=\frac{1}{49}[47412]=967.592$$

$$s \approx \sqrt{967.592} = 31.106$$

3.54

a. A relative frequency histogram for the total damage (dollars) is given here.



- **b.** Mean $\approx 1100
- **c.** Mean = 25115 / 23 = \$1091.96; Median = \$1039
- **d.** Because the mean is slightly larger than the median, it is likely that the distribution is slightly skewed to the right.

38

The means and medians are given here:

Number of Members	Mean	Median
1	93.75	78.5
2	98.652	95
3	113.3125	100
4	124.857	112.5
5+	131.90	128.50

3.56

a. Mean = 9024/83 = 108.7228

b. Yes, we can use the following method: [20(93.75) + 23(98.652) + 16(113.3125) + 14(124.857) + 10(131.90)]/83 = 9023.994/83 = 108.7228

c. Median = 104

d. No.

3.57

a. The sample mean will be distorted by several large values which skew the distribution. State 5 and State 11 have more than 10 times as many plants destroyed as any other state; for arrests, States 1, 2, 8, and 12 exceed the other arrest figures substantially.

b. Plants: $\overline{y} = 10,166,919/15 = 677,794.60$

Arrests: $\overline{y} = 1425 / 15 = 95$

10% trimmed mean:

Plants: $\overline{y} = 1,565,604/11 = 142,327.64$

Arrests: $\overline{y} = 657/11 = 59.7$

20% trimmed mean:

Plants: $\overline{y} = 1{,}197{,}354/9 = 133{,}039.33$

Arrests: $\overline{y} = 372 / 9 = 41.30$

For plants, the 10% trimmed mean works well since it eliminates the effect of States 5 and 11. For arrests, the means differ because each takes some of the high values out of the calculation. It appears that the distribution is not skewed, but rather separated into at least two parts: states with high numbers of arrests and states with low numbers. By trimming the mean, we may be losing important information.

3.58

a. Mean = 1437.93/30 = \$47.93

b. Range = 108.58 - 21.14 = \$87.44

c. DJIA =
$$\frac{\sum_{i=1}^{30} y_i}{C} = \frac{1437.93}{0.1409017} = $10,205.20$$

d. The DJIA does provide information about a population, the population of all companies. However, the sample is not a random sample.

a. The job-history percentages within each source are given here:

Source						
Job History	Within Firm	Related Business	Unrelated Business			
Promoted	22.80	19.05	23.81			
Same position	56.14	38.10	42.86			
Resigned	15.80	28.57	23.81			
Dismissed	5.26	14.28	9.52			
Total	100(n = 57)	100(n = 21)	100(n = 42)			

b. If for each source we compute the percentages combined over the promoted and same position categories, we find that they are 78.94% for within firms, 57.15% for related business, and 66.67% for unrelated business. This ordering by source also holds for every job history category except the "promoted" one in which the three sources are nearly equal. It appears that a company does best when it selects its middle managers from within its own firm and worst when it takes its choices from a related firm.

3.60

- **a.** The value 62 reflects the number of respondents in coal producing states who preferred a national energy policy that encourages coal production. The value 32.8 is of those who favored a coal policy, 32.8% came from major coal producing states. The value 41.3 tells us that 41.3% of those from coal states favored a coal policy. And the value 7.8 tells us that 7.8% of all the responses come from residents of major coal producing states who were in favor of a national energy policy that encouraged coal production.
- **b.** The column percentages because they displayed the distribution of opinions within each of the three types of states.
- **c.** Yes. For both the coal and oil-gas states, the largest percentage of responses favored the type of energy produced in their own state.

3.61

Arbitration seems to win the largest wage increases. If we assume that the Empirical Rule holds for these data, then a standard error for the mean of the arbitration figures would be $s/\sqrt{n}=0.25$. Thus the mean increase after arbitration is $(9.42-8.40)/0.25\approx 4$ standard errors above the next largest mean, that for poststrike. Management, on the other hand, should favor negotiation. It has the smallest mean percentage wage increase and the smallest variance, or least risk.

3.62

a. At each of the clinic visits, there appears to be one person receiving the treatment who has a higher seizure count than others in the study (see the black square in the upper right corner of each panel on the left). There also appears to be one person receiving the treatment who has a high seizure count for his or her age (see the topmost black square around age 22 of each panel on the right). This person appears different than the rest of the patients in some way. Other than that, it is difficult to tell any difference between the placebo and treatment groups at any of the clinic visits when number of seizures is compared to baseline and compared to age.

b. For the fourth clinic visit, it appears that the treatment group has fewer seizures when compared to the baseline number of seizures. Additionally, by the fourth clinic visit, it appears that the treatment is working better than the placebo for people of all ages.

3.63

The treatment group has a stronger linear relationship with the baseline for each of the four clinic visits than does the placebo group, especially for the third clinic visit (Y_3). The placebo group has a weak positive relationship with age for all four clinic visits, while the treatment group has a weak negative relationship with age for all four clinic visits.

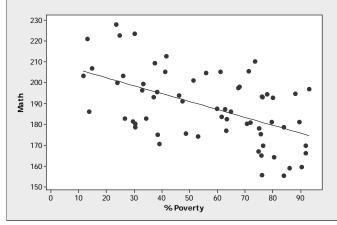
3.64

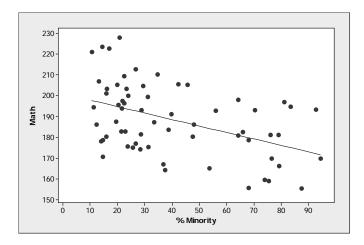
- **a.** When this patient is removed from the data, the correlations should be weaker than they are with this patient included.
- **b.** The correlations for the treatment group with patient ID 207 removed from the data are

	Treatment Group				
	Y_1	Y_2	Y_3	Y_4	Base
Y_2	0.453				
Y_3	0.626	0.699			
Y_4	0.773	0.722	0.826		
Base	0.545	0.513	0.501	0.627	
Age	0.041	-0.218	-0.109	-0.128	-0.355

The correlations are all weaker with patient ID 207 removed than they were when that patient was included.

3.65 Scatterplots for math versus %poverty and %minority are given here:

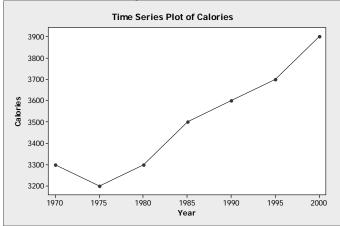




- **a.** The scatterplots for math versus %poverty and %minority are very similar to the scatterplots for reading versus the same two variables (as shown on page 115 of the text). Thus, there is support for the same conclusions for the math scores as for the reading scores.
- **b.** The conclusions are not different.

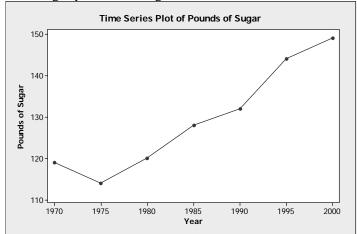
- **a.** If we concluded that large relative values for %minority and %poverty in a school results in lower reading scores for children in these social classes, we would be assuming that %minority and %poverty *cause* the lower reading scores.
- **b.** Some variables might be: expenditures per pupil, student-teacher ratios, teacher training, teacher motivation, prior preparation of the students.

3.67 A time-series plot for calorie intake is given here:



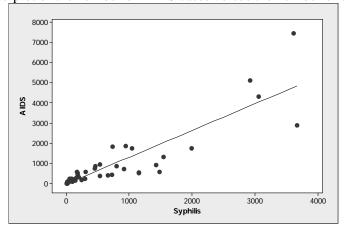
- **a.** Aside from a slight dip at 1975, calorie intake has increased over the 30 years.
- **b.** If we assume that calories would continue increasing through 2005, we would predict that the calorie intake would be about 4000 calories.

3.68 A time-series plot for sugar production is given here:



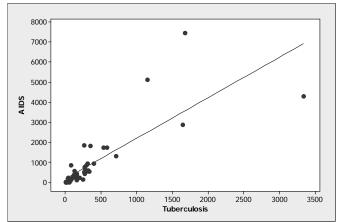
- **a.** Aside from a slight dip at 1975, sugar production has increased over the 30 years.
- **b.** The correlation coefficient between calorie intake and sugar production is 0.986. The correlation coefficient only suggests that there is a strong linear relationship between calorie intake and sugar production. It does not suggest that sugar production is causing the increased calorie intake by the U.S. public.

a. A scatterplot of the number of AIDS cases versus the number of syphilis cases is given here:



- **b.** The correlation between the number of AIDS cases and the number of syphilis cases is 0.883.
- **c.** The scatterplot shows a relatively strong positive linear association, and the correlation coefficient reflects this.
- **d.** Both diseases are sexually-transmitted diseases (STDs). It seems reasonable to believe that a person who puts himself at risk for one STD also puts himself at risk for another STD and that incidences of STDs might occur together.

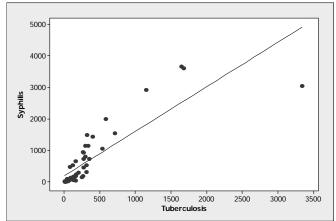
a. A scatterplot of the number of AIDS cases versus the number of tuberculosis cases is given here:



- **b.** The correlation between the number of AIDS cases and the number of tuberculosis cases is 0.806.
- **c.** Both diseases may be associated with drug abusers. It seems reasonable to believe that a person who puts himself at risk for one disease through drug abuse also puts himself at risk for another disease and that incidences of such diseases might occur together.

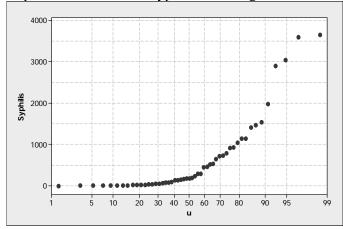
3.71

a. A scatterplot of the number of syphilis cases versus the number of tuberculosis cases is given here:



- **b.** The correlation between the number of syphilis cases and the number of tuberculosis cases is 0.849.
- **c.** Both diseases may be associated with drug abusers. It seems reasonable to believe that a person who puts himself at risk for one on disease through drug abuse also puts himself at risk for another disease and that incidences of such diseases might occur together.

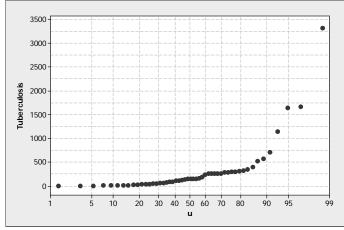
a. A quantile plot of the number of syphilis cases is given here:



- **b.** The 90th percentile is the value associated with u = 0.90 on the graph, which is 1500 cases.
- **c.** There are 5 states that have a number of syphilis cases above the 90th percentile: Georgia, Florida, California, New York, and Texas.

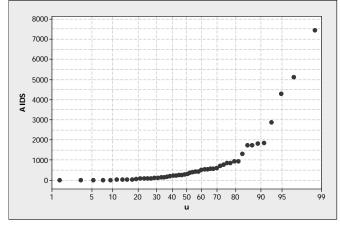
3.73

a. A quantile plot of the number of tuberculosis cases is given here:



- **b.** The 90th percentile is the value associated with u = 0.90 on the graph, which is 750 cases.
- **c.** There are 5 states that have a number of tuberculosis cases above the 90th percentile: Illinois, Florida, Texas, New York, and California.

a. A quantile plot of the number of AIDS cases is given here:

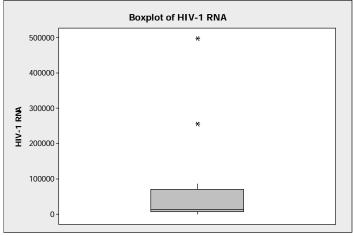


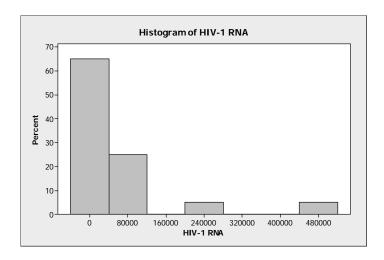
- **b.** The 90th percentile is the value associated with u = 0.90 on the graph, which is 2000 cases.
- **c.** There are 5 states that have a number of AIDS cases above the 90th percentile: Maryland, Texas, California, Florida, and New York.

3.75

- **a.** There were 4 states that had number of AIDS, tuberculosis, and syphilis cases all above the 90th percentiles.
- **b.** These four states were California, Florida, New York, and Texas. All four states are large states, but they are also associated with larger immigrant populations than other states.
- **c.** The U.S. government should educate all people in all states about sexually-transmitted diseases. However, it appears that larger states and states with large immigrant populations need more education.

- **a.** Mean = 61,667.95; Median = 13,956.5; s = 117,539.3
- **b.** 25th percentile = 8,914; 50th percentile = median = 13,956.5; 75th percentile = 63,554.5
- **c.** A boxplot and a histogram of the HIV-1 RNA levels are given here:

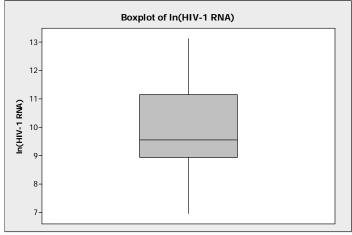


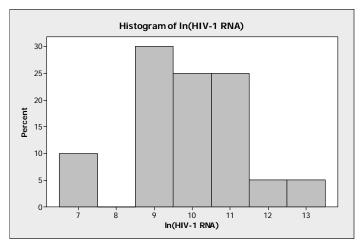


d. The distribution of HIV-1 RNA levels is unimodal and highly skewed to the right with two high outliers (256,440 and 496,433).

- **a.** Mean = 9.905; Median = 9.544; s = 1.550
- **b.** 25th percentile = 8.946; 50th percentile = median = 9.544; 75th percentile = 11.143

c. A boxplot and a histogram of the natural log of HIV-1 RNA levels are given here:

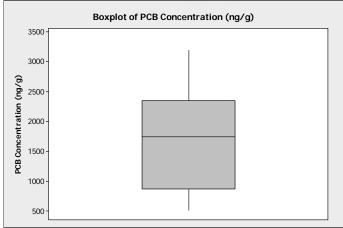




d. The natural logarithm transformation did result in a somewhat symmetric distribution, certainly more symmetric than the original distribution.

- **a.** Mean = 1709.3 ng/g; Median = 1750 ng/g; s = 824.8 ng/g
- **b.** 25th percentile = 878 ng/g; 50th percentile = median = 1750 ng/g; 75th percentile = 2350 ng/g

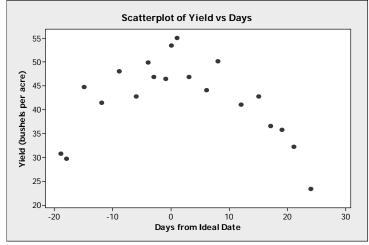
c. A boxplot of the PCB concentrations (ng/g) is given here:



d. $\overline{y} \pm s = 1709.3 \pm 824.8 \Rightarrow (884.5, 2534.1)$ contains 9/14 = 64.3% $\overline{y} \pm 2s = 1709.3 \pm (2)824.8 \Rightarrow (59.7, 3358.9)$ contains 14/14 = 100% $\overline{y} \pm 3s = 1709.3 \pm (3)824.8 \Rightarrow (-765.1, 4183.7)$ contains 14/14 = 100% These values match the values from the Empirical Rule fairly well: 68%, 95%, and 99.7%.

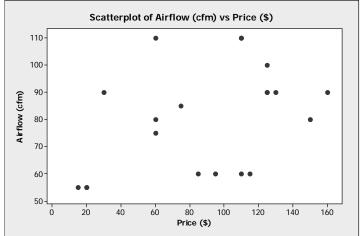
3.79

a. A scatterplot of yield (bushels per acre) versus number of days from the ideal planting date is given here:



- **b.** The relationship between yield (bushels per acre) and number of days from the ideal planting date is curved, increasing from the negative number of days to the ideal planting date and decreasing after the ideal planting date.
- **c.** The correlation coefficient between yield (bushels per acre) and number of days from the ideal planting date is -0.226.
- **d.** The correlation coefficient is relatively small for this data set because of the curved relationship between the two variables. The correlation coefficient is meant to measure the direction and strength of the straight line relationship between two variables, but these variables appear to have a quadratic relationship.

a. A scatterplot of airflow (cfm) versus price (dollars) is given here:



The relationship between airflow and price of the fans appears to be fairly scattered with no pattern emerging.

- **b.** The correlation coefficient between airflow (cfm) and price (dollars) is 0.413. There is a relatively weak relationship between price and airflow of the fans.
- **c.** Yes, we would expect a relatively weak correlation coefficient when we don't see much of a linear relationship in the scatterplot.
- **d.** It would not be reasonable to conclude that higher priced fans generate greater airflow. First, this implies that there is a causal relationship between the two variables. Second, there is not much relationship between the two variables to speak of.