

112 學年度第 1 學期

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

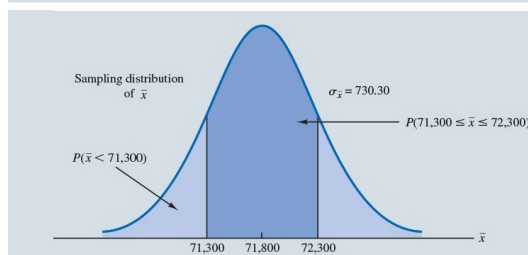
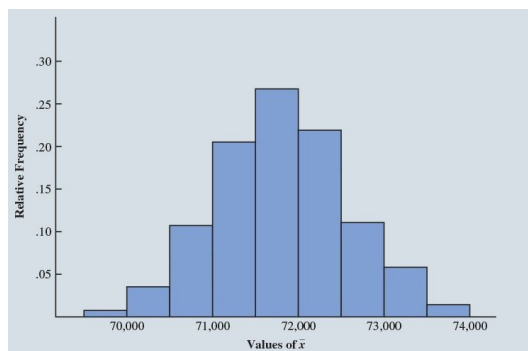
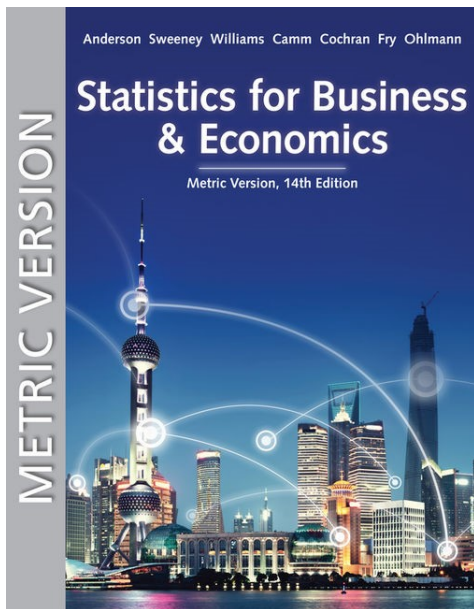
授課教師: 吳漢銘 國立政治大學統計學系

開課單位: 統計系

科目代碼: 000359221

教學網站: <http://www.hmwu.idv.tw>

系級: _____ 學號: _____ 姓名: _____



目錄

- Ch 1. Data and Statistics
 - Ch 2. Descriptive Statistics: Tabular and Graphical Displays
 - Ch 3. Descriptive Statistics: Numerical Measures
 - Ch 4. Introduction to Probability
 - Ch 5. Discrete Probability Distributions
 - Ch 6. Continuous Probability Distributions
 - Ch 7. Sampling and Sampling Distributions
 - Ch 8. Interval Estimation
 - Ch 9. Hypothesis Tests
 - Ch 10. Inference about Means and Proportions with Two Populations (Optional)
- 附錄：111-1 學年第 1 學期小考題、期中考題、期末考題。

叮嚀

- A. 平常就要唸書，做習題。
- B. 考過的題目，要主動訂正。
- C. 上課以「互相尊重」為最高原則並盡到「告知老師」的義務。
- D. 上課可小聲討論、上廁所安靜去回、不鼓勵飲食。(請一定要維護教室整潔)
- E. 四不一要: 「上課不聊天，睡覺不趴著，手機不要滑，考試不作弊，要認真。」

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 1: Data and Statistics

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

Overview: What is Statistics?

1. The term statistics can refer to numerical facts such as averages, medians, percentages, and maximums that help us understand a variety of business and economic situations.
2. Wikipedia: Statistics is the discipline (art and science) that concerns the collection, organization, analysis, interpretation, and presentation of data.

1.1 Applications in Business and Economics

Accounting

Public accounting firms use statistical sampling procedures when conducting audits for their clients.

Economics

Economists use statistical information in making forecasts about the future of the economy or some aspect of it.

Finance

Financial advisors use price-earnings ratios (本益比) and dividend yields (現金殖利率) to guide their investment advice.

Marketing

Electronic point-of-sale scanners at retail checkout counters are used to collect data for a variety of marketing research applications.

Production

A variety of statistical quality control charts are used to monitor the output of a production process.

Information Systems

A variety of statistical information helps administrators assess the performance of computer networks.

1.2 Data

1. Data are the facts and figures collected, analyzed, and summarized for presentation and interpretation.
2. All the data collected in a particular study are referred to as the data set for the study.

Elements, Variables, and Observations

1. Elements are the entities on which data are collected.
2. A variable is a characteristic of interest for the elements.
3. The set of measurements obtained for a particular element is called an observation.
4. A data set with n elements contains n observations.

5. The total number of data values in a complete data set is the number of elements multiplied by the number of variables.
6. Example WTO, 惠譽國際 (Fitch Group) (Per Capita 人均)

Nation	WTO Status	Per Capita GDP (\$)	Fitch Rating	Fitch Outlook
Armenia	Member	3,615	BB-	Stable
Australia	Member	49,755	AAA	Stable
Austria	Member	44,758	AAA	Stable
Azerbaijan	Observer	3,879	BBB-	Stable
Bahrain	Member	22,579	BBB	Stable
Belgium	Member	41,271	AA	Stable
Brazil	Member	8,650	BBB	Stable
Bulgaria	Member	7,469	BBB-	Stable
Canada	Member	42,349	AAA	Stable
Cape Verde	Member	2,998	B+	Stable
Chile	Member	13,793	A+	Stable
China	Member	8,123	A+	Stable
Colombia	Member	5,806	BBB-	Stable
Costa Rica	Member	11,825	BB+	Stable
Croatia	Member	12,149	BBB-	Negative

Scales of Measurement

- Scales of measurement include: Nominal, Ordinal, Interval, Ratio.
- The scale determines the amount of information contained in the data.
- The scale indicates the data stigmatization (污名化) and statistical analyses that are most appropriate.
- Nominal scale**
 - Data are labels or names used to identify an attribute of the element.
 - A non-numeric label or numeric code may be used.
 - Example: Students of a university are classified by the school in which they are enrolled using a nonnumeric label such as Business, Humanities, Education, and so on. Alternatively, a numeric code could be used for the school variable (e.g., 1 denotes Business, 2 denotes Humanities, 3 denotes Education, and so on).

5. Ordinal scale

- (a) The data have the properties of nominal data and the order or rank of the data is meaningful.
- (b) A non-numeric label or numeric code may be used.
- (c) Example: Students of a university are classified by their class standing using a non-numeric label such as Freshman, Sophomore, Junior, or Senior. Alternatively, a numeric code could be used for the class standing variable (e.g., 1 denotes Freshman, 2 denotes Sophomore, and so on).

6. Interval scale

- (a) The data have the properties of ordinal data, and the interval between observations is expressed in terms of a fixed unit of measure.
- (b) Interval data are always numeric.
- (c) Example: Melissa has an SAT score of 1985, while Kevin has an SAT score of 1880. Melissa scored 105 points more than Kevin.

7. Ratio scale

- (a) Data have all the properties of interval data and the ratio of two values is meaningful.
- (b) Ratio data are always numerical.
- (c) Zero value is included in the scale.
- (d) Example: Price of a book at a retail store is \$200, while the price of the same book sold online is \$100. The ratio property shows that retail stores charge twice the online price.

Categorical and Quantitative Data

1. Data can be further classified as being categorical or quantitative.
2. The statistical analysis that is appropriate depends on whether the data for the variable are categorical or quantitative.

3. Categorical Data

- (a) Labels or names are used to identify an attribute of each element.
- (b) Often referred to as qualitative data.
- (c) Use either the nominal or ordinal scale of measurement.
- (d) Can be either numeric or nonnumeric.
- (e) Appropriate statistical analyses are rather limited.

4. Quantitative Data

- (a) Quantitative data indicate how many or how much.
- (b) Quantitative data are always numeric.
- (c) Ordinary arithmetic operations are meaningful for quantitative data.

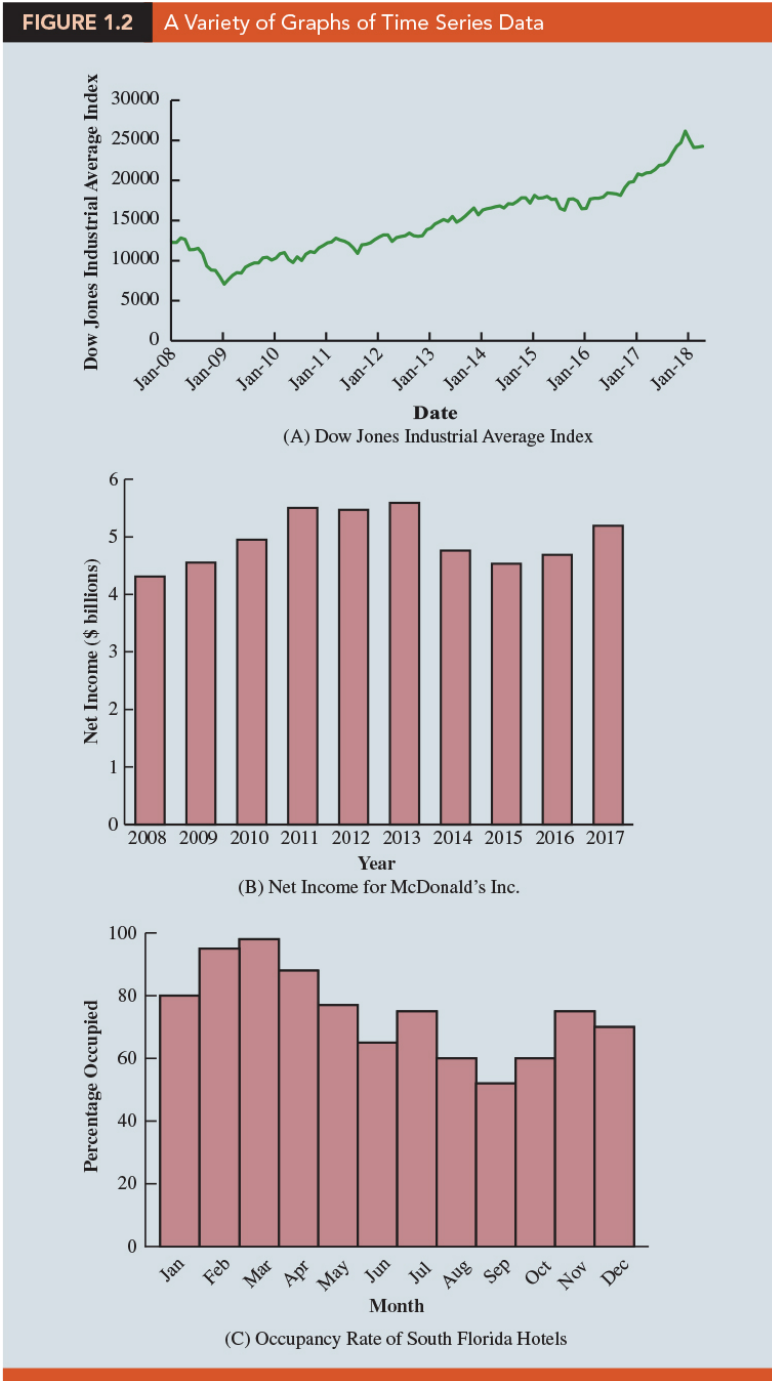
Cross-Sectional and Time Series Data

1. Cross-sectional data

- (a) Cross-sectional data are collected at the same or approximately the same point in time.
- (b) Example: The data in Table 1.1 are cross-sectional because they describe the five variables for the 60 World Trade Organization nations at the same point in time.

2. Time Series Data

- (a) Time series data are collected over several time periods.
- (b) Example: Data detailing the number of building permits issued in Lucas County, Ohio in each of the last 36 months.
- (c) Graphs of time series data help analysts understand: (i) what happened in the past, (ii) identify any trends over time, and (iii) predict future levels for the time series.



1.3 Data Sources

Existing Sources

1. Internal company records – almost any department
2. Business database services – Dow Jones & Co.
3. Government agencies – U.S. Department of Labor
4. Industry associations – Travel Industry Association of America
5. Special-interest organizations – Graduate Management Admission Council (GMAT)
6. Internet –e.g., 政府資料開放平臺: <https://data.gov.tw>

Observational Study

1. In observational (nonexperimental) studies no attempt is made to control or influence the variables of interest.
2. Example: Survey and public opinion, e.g., studies of smokers and nonsmokers are observational studies because researchers do not determine or control who will smoke and who will not smoke.

Experiment

1. In experimental studies the variable of interest is first identified. Then one or more other variables are identified and controlled so that data can be obtained about how they influence the variable of interest.
2. Example: The largest experimental study ever conducted is believed to be the 1954 Public Health Service experiment for the Salk polio vaccine (小兒麻痺疫苗). Nearly two million U.S. children (grades 1- 3) were selected.

Time and Cost Issues

1. Searching for information can be time consuming.
2. Information may no longer be useful by the time it is available.

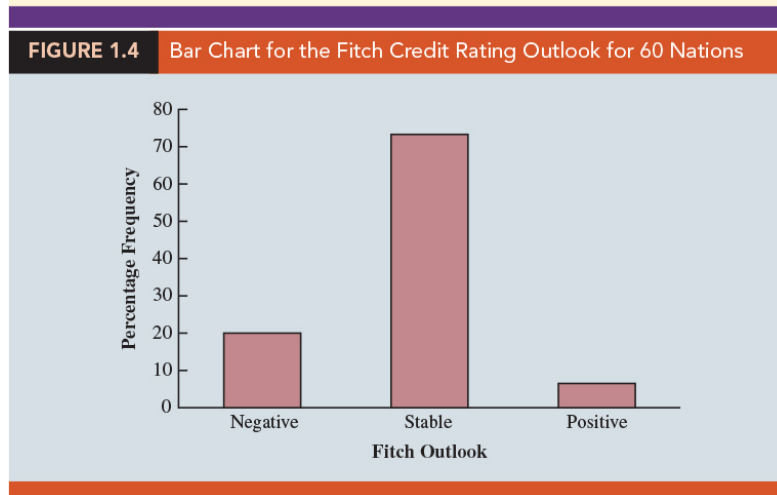
Data Acquisition Errors

1. Organizations often charge for information even when it is not their primary business activity.
2. Using any data that happen to be available or were acquired with little care can lead to misleading information.

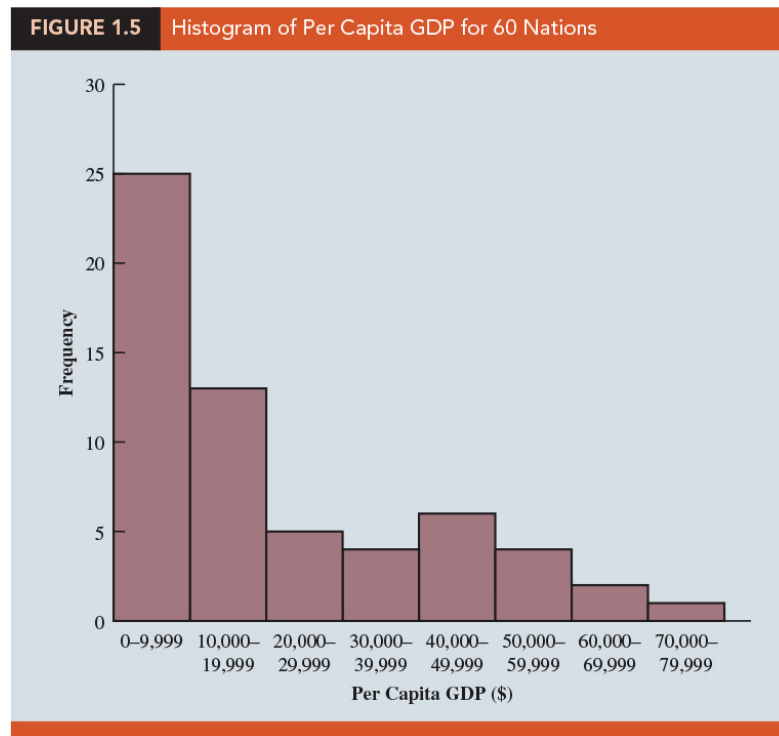
1.4 Descriptive Statistics

1. Most of the statistical information in the media, company reports, and other publications consists of data that are summarized and presented in a form that is easy for the reader to understand. Such summaries of data, which may be tabular, graphical, or numerical, are referred to as descriptive statistics.
2. **Example** A tabular summary of the data showing the number of nations with each of the Fitch Outlook ratings (Table 1.4). A graphical summary of the same data, called a bar chart (Figure 1.4). We can see that the majority of Fitch Outlook credit ratings are stable, with 73.3% of the nations having this rating. More nations have a negative outlook (20%) than a positive outlook (6.7%).

Fitch Outlook	Frequency	Percent Frequency (%)
Positive	4	6.7
Stable	44	73.2
Negative	12	20.0



3. **Example** A graphical summary of the data for the quantitative variable Per Capita GDP in Table 1.1, called a histogram (Figure 1.5). Using the histogram, it is easy to see that Per Capita GDP for the 60 nations ranges from \$0 to \$80,000, with the highest concentration between \$0 and \$10,000. Only one nation had a Per Capita GDP exceeding \$70,000.



4. The most common numerical descriptive statistic is the mean (or average).
5. The mean demonstrates a measure of the central tendency, or central location of the data for a variable.

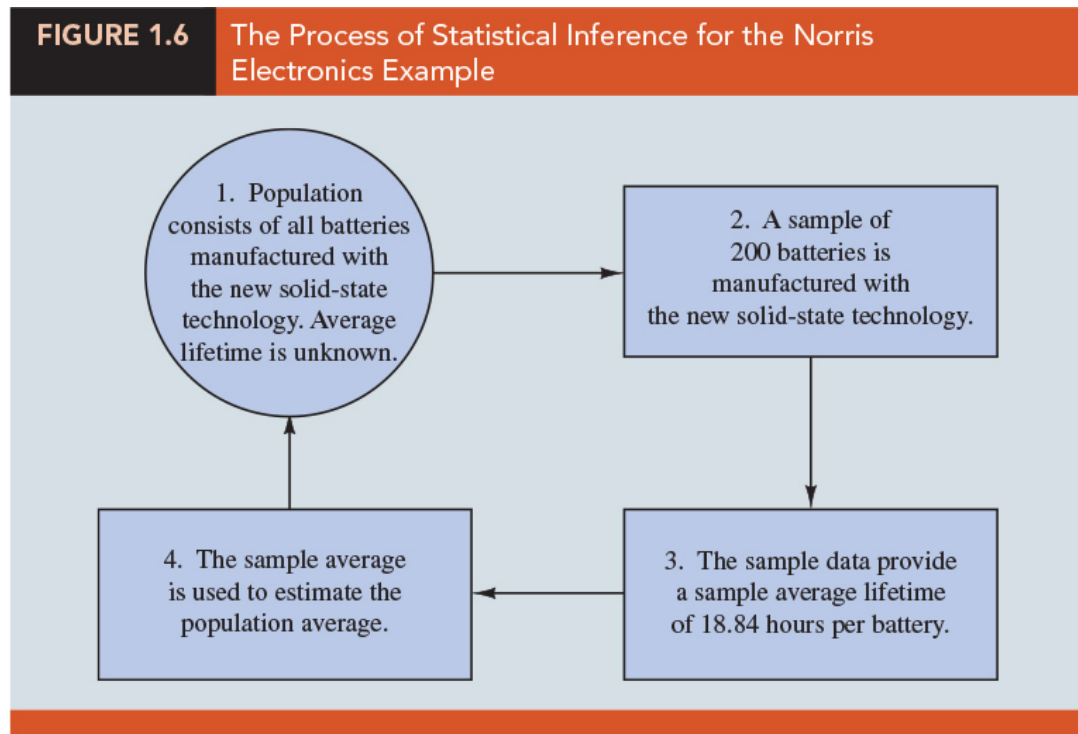
1.5 Statistical Inference

- Population:** The set of all elements of interest in a particular study.
- Sample:** A subset of the population.
- Statistical inference:** The process of using data obtained from a sample to make estimates and test hypotheses about the characteristics of a population.

4. **Census:** Collecting data for the entire population.
5. **Sample survey:** Collecting data for a sample.
6. **Example** Rogers Industries manufactures lithium batteries used in rechargeable electronics such as laptop computers and tablets. In an attempt to increase battery life for its products, Rogers has developed a new solid-state lithium battery that should last longer and be safer to use.
 - (a) population: all lithium batteries that could be produced using the new solid-state technology.
 - (b) sample (data): 200 batteries manufactured with the new solid-state technology were tested. (Table 1.5) the number of hours each battery lasted before needing to be recharged under controlled conditions.

TABLE 1.5 Hours Until Recharge for a Sample of 200 Batteries for the Rogers Industries Example									
Battery Life (hours)									
19.49	18.18	18.65	19.45	19.89	18.94	17.72	18.35	18.66	18.23
19.08	19.92	19.01	18.84	17.73	19.70	18.37	18.69	19.98	18.80
19.11	18.26	19.05	17.89	19.61	18.52	18.10	19.08	18.27	18.29
19.55	18.81	18.68	17.43	20.34	17.73	17.66	18.52	19.90	19.33
18.81	19.12	18.39	19.27	19.43	19.29	19.11	18.96	19.65	18.20
19.18	20.07	18.54	18.37	18.13	18.29	19.11	20.22	18.07	18.91

- (c) Suppose Rogers wants to use the sample data to make an inference about the average hours of battery life for the population of all batteries that could be produced with the new solid-state technology.
- (d) The sample average battery life: 18.84 hours. We can use this sample result to estimate that the average lifetime for the batteries in the population is 18.84 hours.
- (e) (Figure 1.6) a graphical summary of the statistical inference process for Rogers Industries.



1.6 Analytics

1. Analytics is the scientific process of transforming data into insight for making better decisions.
2. Three categories of techniques:
 - (a) **Descriptive analytics**: This describes what has happened in the past.
 - (b) **Predictive analytics**: Use models constructed from past data to predict the future or to assess the impact of one variable on another.
 - (c) **Prescriptive analytics**: The set of analytical techniques that yield a best course of action.

1.7 Big Data and Data Mining

1. **Big data**: Large and complex data set. Three V's of Big data:
 - (a) Volume: Amount of available data.
 - (b) Velocity: Speed at which data is collected and processed.

- (c) Variety : Different data types.

2. Data Mining

- (a) Methods for developing useful decision-making information from large databases .
- (b) Using a combination of procedures from statistics, mathematics, and computer science, analysts "mine the data" to convert it into useful information.
- (c) The most effective data mining systems use automated procedures to discover relationships in the data and predict future outcomes prompted by general and even vague queries by the user.
- (d) Statistical methodology such as multiple regression, logistic regression, and correlation are heavily used.
- (e) Also needed are computer science technologies involving artificial intelligence and machine learning .
- (f) With the enormous amount of data available, the data set can be partitioned into a training set (for model development) and a test set (for validating the model).
- (g) Careful interpretation of results and extensive testing is important.

1.8 Computers and Statistical Analysis

1. Computer: laptop, PC, Server, ...
2. OS: Windows, Mac OS, Linux, ...
3. Software for statistical analysis: Microsoft Excel, JMP, SPSS, SAS, R, Python

1.9 Ethical Guidelines for Statistical Practice

1. In a statistical study, unethical behavior can take a variety of forms including: improper sampling , inappropriate analysis of the data, development of misleading graphs , use of inappropriate summary statistics , biased interpretation of the statistical results.

2. One should strive to be fair, thorough, objective, and neutral (中性的) as you collect, analyze, and present data.
3. As a consumer of statistics, one should also be aware of the possibility of unethical behavior by others.
4. The American Statistical Association (ASA) developed the report “Ethical Guidelines for Statistical Practice” . It contains 67 guidelines organized into 8 topic areas:
 - (a) Professionalism.
 - (b) Responsibilities to Funders, Clients, and Employers.
 - (c) Responsibilities in Publications and Testimony.
 - (d) Responsibilities to Research Subjects.
 - (e) Responsibilities to Research Team Colleagues.
 - (f) Responsibilities to Other Statisticians/Practitioners.
 - (g) Responsibilities Regarding Allegations of Misconduct.
 - (h) Responsibilities of Employers Including Organizations, Individuals, Attorneys, or Other Clients.

☺ SUPPLEMENTARY EXERCISES

3, 4, 5, 9, 13, 14, 19, 25

“不要努力想成為一個成功者，要努力成為一個有價值的人”

“Try not to become a man of success, but rather try to become a man of value”

— *Albert Einstein (March 14, 1879 – April 18, 1955)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 2: Descriptive Statistics: Tabular and Graphical Displays

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

2.1 Summarizing Data for a Categorical Variable

Frequency Distribution

1. **Frequency distribution:** A frequency distribution is a tabular summary of data showing the number (frequency) of observations in each of several non-overlapping categories or classes.
2. **Example** Soft Drinks Data (Table 2.1)
 - (a) Coca-Cola, Diet Coke, Dr. Pepper, Pepsi, and Sprite are five popular soft drinks. The data show below the soft drink selected in a sample of 50 soft drink purchases.

Coca-Cola	Coca-Cola	Coca-Cola	Sprite	Coca-Cola
Diet Coke	Dr. Pepper	Diet Coke	Dr. Pepper	Diet Coke
Pepsi	Sprite	Coca-Cola	Pepsi	Pepsi
Diet Coke	Coca-Cola	Sprite	Diet Coke	Pepsi
Coca-Cola	Diet Coke	Pepsi	Pepsi	Pepsi
Coca-Cola	Coca-Cola	Coca-Cola	Coca-Cola	Pepsi
Dr. Pepper	Coca-Cola	Coca-Cola	Coca-Cola	Coca-Cola
Diet Coke	Sprite	Coca-Cola	Coca-Cola	Dr. Pepper
Pepsi	Coca-Cola	Pepsi	Pepsi	Pepsi
Pepsi	Diet Coke	Coca-Cola	Dr. Pepper	Sprite

- (b) The number of times each soft drink (frequency distribution) summarizes information about the popularity of the five soft drinks.

Soft Drink	Frequency
Coca-Cola	19
Diet Coke	8
Dr. Pepper	5
Pepsi	13
Sprite	5
Total	50

Relative Frequency and Percent Frequency Distribution

- The relative frequency of a class equals the fraction or proportion of observations belonging to a class.
- For a data set with n observations, the relative frequency of each class is:

$$\text{Relative frequency of a class} = \frac{\text{Frequency of the class}}{n}$$

- A relative frequency distribution gives a tabular summary of data showing the relative frequency for each class.
- The percent frequency of a class is the relative frequency multiplied by 100. A percent frequency distribution summarizes the percent frequency of the data for each class.

5. **Example** Soft Drinks Data (Table 2.3)

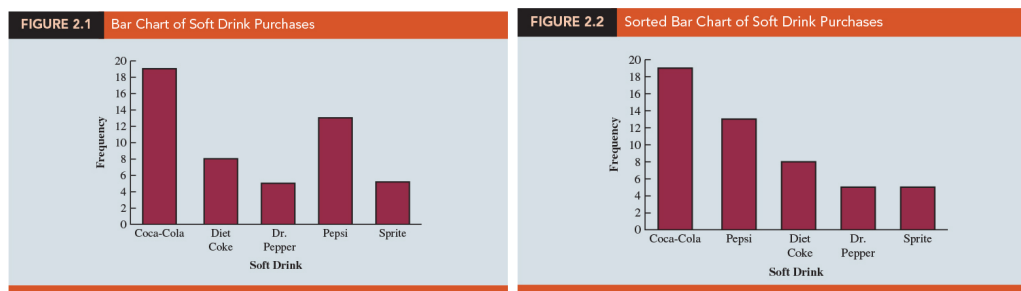
- The relative frequency for Coca-Cola is $19/50 = 0.38$, the relative frequency for Diet Coke is $8/50 = 0.16$, and so on.
- From the percent frequency distribution, we see that 38% of the purchases were Coca-Cola, 16% of the purchases were Diet Coke, and so on.
- We can also note that $38\% + 26\% + 16\% = 80\%$ of the purchases were for the top three soft drinks.

TABLE 2.3 Relative Frequency and Percent Frequency Distributions of Soft Drink Purchases

Soft Drink	Relative Frequency	Percent Frequency
Coca-Cola	.38	38
Diet Coke	.16	16
Dr. Pepper	.10	10
Pepsi	.26	26
Sprite	.10	10
Total	1.00	100

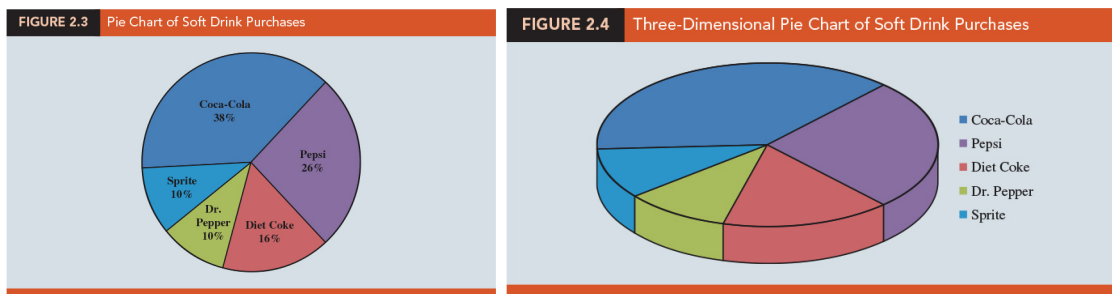
Bar Chart

1. A bar chart is a graphical display for depicting categorical (qualitative) data summarized in a frequency, relative frequency, or percent frequency distribution.
2. On one axis of the chart (usually the horizontal axis), we specify the labels that are used for the classes (categories). A frequency, relative frequency, or percent frequency scale can be used for the other axis of the chart (usually the vertical axis).
3. Using a bar of fixed width drawn above each class label, we extend the length of the bar until we reach the frequency, relative frequency, or percent frequency of the class.
4. For categorical data, the bars should be separated to emphasize the fact that each category is separate.
5. Example Soft Drinks Data (Figure 2.1)(Figure 2.2)



Pie Chart

1. The pie chart is a commonly used graphical display for presenting relative frequency and percent frequency distributions for categorical data.
2. First draw a circle, then use the relative frequencies to subdivide the circle into sectors that correspond to the relative frequency for each class.
3. Because there are 360 degrees in a circle, a class with a relative frequency of 0.25 would consume $0.25(360) = 90$ degrees of the circle.
4. In most cases, a bar chart is superior to a pie chart for displaying categorical data.
5. Numerous options involving the use of colors, shading, legends, text font, and three-dimensional perspectives are available to enhance the visual appearance of bar and pie charts.
6. In general, pie charts are not the best way to present percentages for comparison.
7. **Example** Soft Drinks Data (Figure 2.3)(Figure 2.4)



2.2 Summarizing Data for a Quantitative Variable

Frequency Distribution

1. The three steps necessary to define the classes for a frequency distribution with quantitative data are
 - (a) **Determine the number of nonoverlapping classes:** Classes are formed by specifying ranges that will be used to group the data. The goal is to use

enough classes to show the variation in the data, but not so many classes that some contain only a few data items.

- (b) **Determine the width of each class:** As a general guideline, we recommend that the width be the same for each class:

$$\text{Approximate class width} = \frac{\text{Largest data value} - \text{Smallest data value}}{\text{Number of classes}}$$

The approximate class width can be rounded to a more convenient value based on the preference of the person developing the frequency distribution.

- (c) **Determine the class limits:** Class limits must be chosen so that each data item belongs to one and only one class.

2. With the number of classes, class width, and class limits determined, a frequency distribution can be obtained by counting the number of data values belonging to each class.

3. **Example** Audit Time Data (Table 2.4.)(Table 2.5.)

12	14	19	18
15	15	18	17
20	27	22	23
22	21	33	28
14	18	16	13

- (a) **Number of Classes:** The number of data items relatively small ($n = 20$), we chose to develop a frequency distribution with five classes.
- (b) **Width of the Classes:** After deciding to use five classes, the width for each class is five days.
- (c) **Class limits:** A total of five classes: 10–14, 15–19, 20–24, 25–29, and 30–34.

Audit Time (days)	Frequency
10–14	4
15–19	8
20–24	5
25–29	2
30–34	1
Total	20

- (d) The class midpoint is the value halfway between the lower and upper class limits. For the audit time data, the five class midpoints are 12, 17, 22, 27, and 32.

Relative Frequency and Percent Frequency Distributions

- The relative frequency is the proportion of the observations belonging to a class. With n observations,

$$\text{Relative frequency of the class} = \frac{\text{Frequency of class}}{n}$$

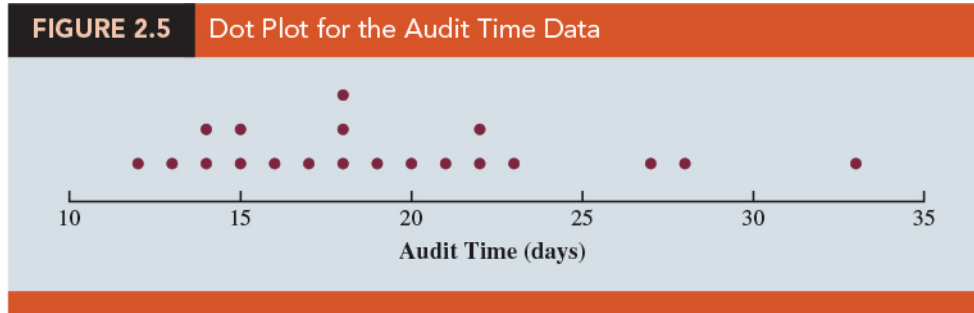
- The percent frequency of a class is the relative frequency multiplied by 100.
- Example** Audit Time Data (Table 2.6.)

Audit Time (days)	Relative Frequency	Percent Frequency
10–14	.20	20
15–19	.40	40
20–24	.25	25
25–29	.10	10
30–34	.05	5
Total	1.00	100

Dot Plot

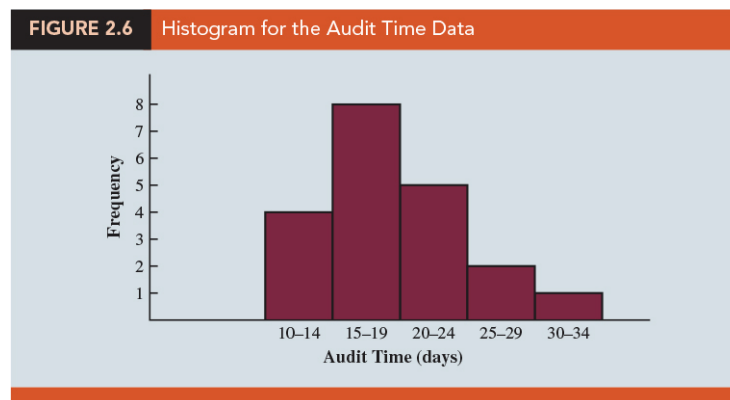
- A horizontal axis of a dot plot shows the range of data values.

- Then each data value is represented by a dot placed above the axis. A vertical axis shows the number of data values.
- Example** Audit Time Data (Figure 2.5.)



Histogram

- A histogram is constructed by placing the variable of interest on the horizontal axis and the frequency, relative frequency, or percent frequency on the vertical axis.
- The frequency, relative frequency, or percent frequency of each class is shown by drawing a rectangle whose base is determined by the class limits on the horizontal axis and whose height is the corresponding frequency, relative frequency, or percent frequency.
- Example** Audit Time Data (Figure 2.6.)
Note that the class with the greatest frequency is shown by the rectangle appearing above the class of 15–19 days. The height of the rectangle shows that the frequency of this class is 8.



4. The adjacent rectangles of a histogram touch one another. Unlike a bar graph, a histogram contains no natural separation between the rectangles of adjacent classes.
5. (Figure 2.7.) A histogram provides information about the shape, or form, of a distribution.
- (a) **Panel A:** A histogram is said to be skewed to the left if its tail extends farther to the left. E.g., exam scores; with no scores above 100%, most of the scores above 70%, and only a few really low scores.
- (b) **Panel B:** A histogram is said to be skewed to the right if its tail extends farther to the right. E.g., housing prices; a few expensive houses create the skewness in the right tail.
- (c) **Panel C:** In a symmetric histogram, the left tail mirrors the shape of the right tail. Data for SAT scores, heights and weights of people, and so on lead to histograms that are roughly symmetric.
- (d) **Panel D:** a histogram highly skewed to the right. E.g., data on housing prices, salaries, purchase amounts, and so on often result in histograms skewed to the right.



Cumulative Distributions

- The cumulative frequency distribution shows the number of data items with values less than or equal to the upper class limit of each class.

- Example** Audit Time Data (Table 2.7.)

The cumulative frequency for the class ("less than or equal to 24") is simply the sum of the frequencies for all classes with data values less than or equal to 24.

Audit Time (days)	Cumulative Frequency	Cumulative Relative Frequency	Cumulative Percent Frequency
Less than or equal to 14	4	.20	20
Less than or equal to 19	12	.60	60
Less than or equal to 24	17	.85	85
Less than or equal to 29	19	.95	95
Less than or equal to 34	20	1.00	100

Stem-and-Leaf Display

- A stem-and-leaf display is a graphical display used to show simultaneously the rank order and shape of a distribution of data.

- Example** Aptitude Test Data (Table 2.8.)

These data result from a 150-question aptitude test given to 50 individuals recently interviewed for a position at Haskens Manufacturing. The data indicate the number of questions answered correctly.

112	72	69	97	107
73	92	76	86	73
126	128	118	127	124
82	104	132	134	83
92	108	96	100	92
115	76	91	102	81
95	141	81	80	106
84	119	113	98	75
68	98	115	106	95
100	85	94	106	119

6	8	9									
7	2	3	3	5	6	6					
8	0	1	1	2	3	4	5	6			
9	1	2	2	2	4	5	5	6	7	8	8
10	0	0	2	4	6	6	6	7	8		
11	2	3	5	5	8	9	9				
12	4	6	7	8							
13	2	4									
14	1										

3. **Stretched Stem-and-Leaf Display:** If we believe the original stem-and-leaf display has condensed the data too much, we can stretch the display vertically by using two stems for each leading digit(s).

6		8	9				
7		2	3	3			
7		5	6	6			
8		0	1	1	2	3	4
8		5	6				
9		1	2	2	2	4	
9		5	5	6	7	8	8
10		0	0	2	4		
10		6	6	6	7	8	
11		2	3				
11		5	5	8	9	9	
12		4					
12		6	7	8			
13		2	4				
13							
14		1					

2.3 Summarizing Data for Two Variables Using Tables

Crosstabulation

1. A crosstabulation is a tabular summary of data for two variables. Both variables can be either categorical or quantitative. One variable is categorical and the other variable is quantitative are just as common.
2. **Example** Quality Rating and Meal Price Data for 300 Los Angeles Restaurants (Table 2.9.) (the first 10 restaurants)

Restaurant	Quality Rating	Meal Price (\$)
1	Good	18
2	Very Good	22
3	Good	28
4	Excellent	38
5	Very Good	33
6	Good	28
7	Very Good	19
8	Very Good	11
9	Very Good	23
10	Good	13
.	.	.
.	.	.
.	.	.

- (a) Quality rating is a categorical variable with rating categories of good, very good, and excellent. Meal price is a quantitative variable that ranges from \$10 to \$49.

TABLE 2.10 Crosstabulation of Quality Rating and Meal Price Data for 300 Los Angeles Restaurants

Quality Rating	Meal Price				Total
	\$10–19	\$20–29	\$30–39	\$40–49	
Good	42	40	2	0	84
Very Good	34	64	46	6	150
Excellent	2	14	28	22	66
Total	78	118	76	28	300

- (b) The right and bottom margins of the crosstabulation provide the frequency distributions for quality rating and meal price separately.
- (c) Note that the values in the relative frequency column do not add exactly to 1.00 and the values in the percent frequency distribution do not add exactly to 100; the reason is that the values being summed are rounded.
- (d) Restaurants with higher meal prices received higher quality ratings than restaurants with lower meal prices.
3. When quantitative variables are used, however, we must first create classes for the values of the variable.
4. **Example** Quality Rating and Meal Price Data (Table 2.11.)
For row percentages, the results of dividing each frequency in Table 2.10 by its corresponding row total are shown in Table 2.11. Each row of Table 2.11 is a percent frequency distribution of meal price for one of the quality rating categories.

TABLE 2.11 Row Percentages for Each Quality Rating Category

Quality Rating	Meal Price				Total
	\$10–19	\$20–29	\$30–39	\$40–49	
Good	50.0	47.6	2.4	.0	100
Very Good	22.7	42.7	30.6	4.0	100
Excellent	3.0	21.2	42.4	33.4	100

Simpson's Paradox

1. The data in two or more crosstabulations are often combined or aggregated to produce a summary crosstabulation showing how two variables are related. In such

cases, conclusions drawn from two or more separate crosstabulations can be reversed when the data are aggregated into a single crosstabulation. The reversal of conclusions based on aggregate and unaggregated data is called Simpson's paradox.

2. **Example** Two Judges in Two different Courts

- (a) Judges Ron Luckett and Dennis Kendall presided over cases in Common Pleas Court (民事訴訟法庭) and Municipal Court (市政法庭) during the past three years. Some of the verdicts they rendered were appealed. In most of these cases the appeals (上訴) court upheld (堅持) the original verdicts, but in some cases those verdicts were reversed.
- (b) For each judge a crosstabulation was developed based upon two variables: Verdict (upheld or reversed) and Type of Court (Common Pleas and Municipal).

Judge			
Verdict	Luckett	Kendall	Total
Upheld	129 (86%)	110 (88%)	239
Reversed	21 (14%)	15 (12%)	36
Total (%)	150 (100%)	125 (100%)	275

- (c) This crosstabulation shows the number of appeals in which the verdict was upheld and the number in which the verdict was reversed for both judges. A review of the column percentages shows that 86% of the verdicts were upheld for Judge Luckett, while 88% of the verdicts were upheld for Judge Kendall.
- (d) From this aggregated crosstabulation, we conclude that Judge Kendall is doing the better job because a greater percentage of Judge Kendall's verdicts are being upheld.
- (e) When we unaggregate the data, we see that Judge Luckett has a better record because a greater percentage of Judge Luckett's verdicts are being upheld in both courts.

Judge Luckett			
Verdict	Common Pleas	Municipal Court	Total
Upheld	29 (91%)	100 (85%)	129
Reversed	3 (9%)	18 (15%)	21
Total (%)	32 (100%)	118 (100%)	150

Judge Kendall			
Verdict	Common Pleas	Municipal Court	Total
Upheld	90 (90%)	20 (80%)	110
Reversed	10 (10%)	5 (20%)	15
Total (%)	100 (100%)	25 (100%)	125

- (f) Note that for both judges the percentage of appeals that resulted in reversals was much higher in Municipal Court than in Common Pleas Court. Because Judge Luckett tried a much higher percentage of his cases in Municipal Court, the aggregated data favored Judge Kendall. When we look at the crosstabulations for the two courts separately, however, Judge Luckett shows the better record.
- (g) Thus, for the original crosstabulation, we see that the type of court is a hidden variable that cannot be ignored when evaluating the records of the two judges.
3. Before drawing a conclusion, you may want to investigate whether the aggregated or unaggregated form of the crosstabulation provides the better insight and conclusion.

2.4 Summarizing Data for Two Variables Using Graphical Displays

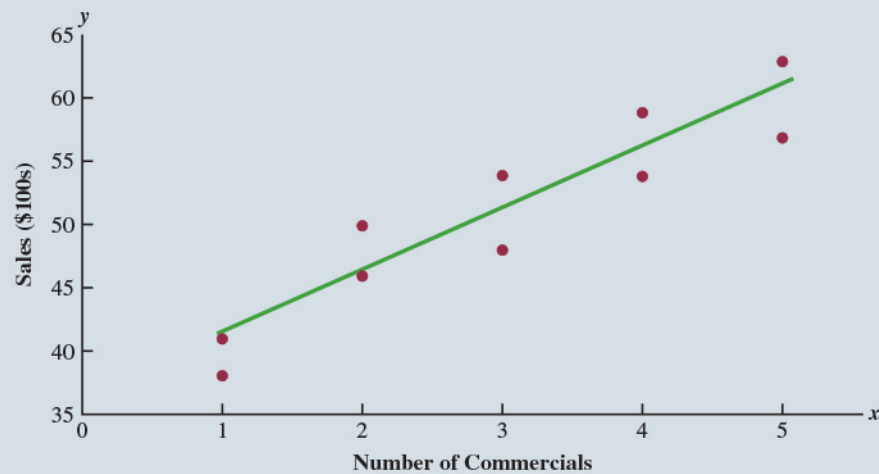
Scatter Diagram and Trendline

1. A scatter diagram (scatter plot) is a graphical presentation of the overall relationship between two quantitative variables.
2. One variable is shown on the horizontal axis and the other variable is shown on the vertical axis.
3. A trendline provides an approximation of the relationship.
4. Example Electronics Store in San Francisco (Table 2.14.)

- (a) Consider the advertising/sales relationship for an electronics store in San Francisco. On 10 occasions during the past three months, the store used weekend television commercials to promote sales at its stores. The managers want to investigate whether a relationship exists between the number of commercials shown and sales at the store during the following week. Sample data for the 10 weeks with sales in hundreds of dollars are shown in Table 2.14.

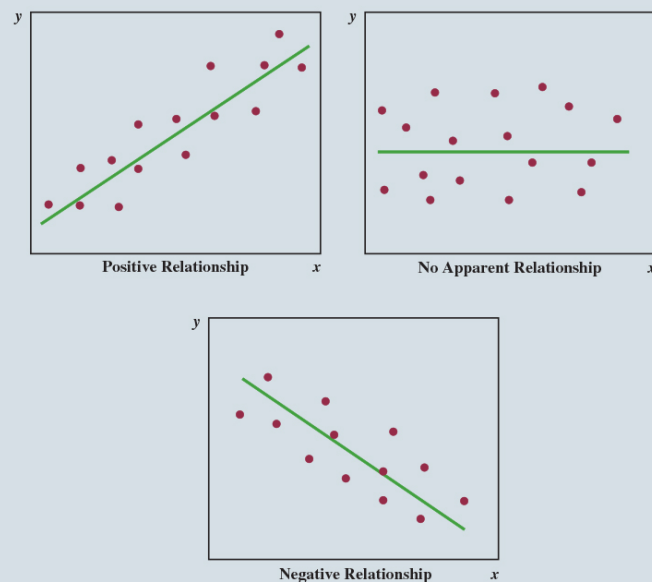
Week	Number of Commercials <i>x</i>	Sales (\$100s) <i>y</i>
1	2	50
2	5	57
3	1	41
4	3	54
5	4	54
6	1	38
7	5	63
8	3	48
9	4	59
10	2	46

- (b) The scatter diagram in Figure 2.8 indicates a positive relationship between the number of commercials and sales. Higher sales are associated with a higher number of commercials.

FIGURE 2.8 Scatter Diagram and Trendline for the San Francisco Electronics Store

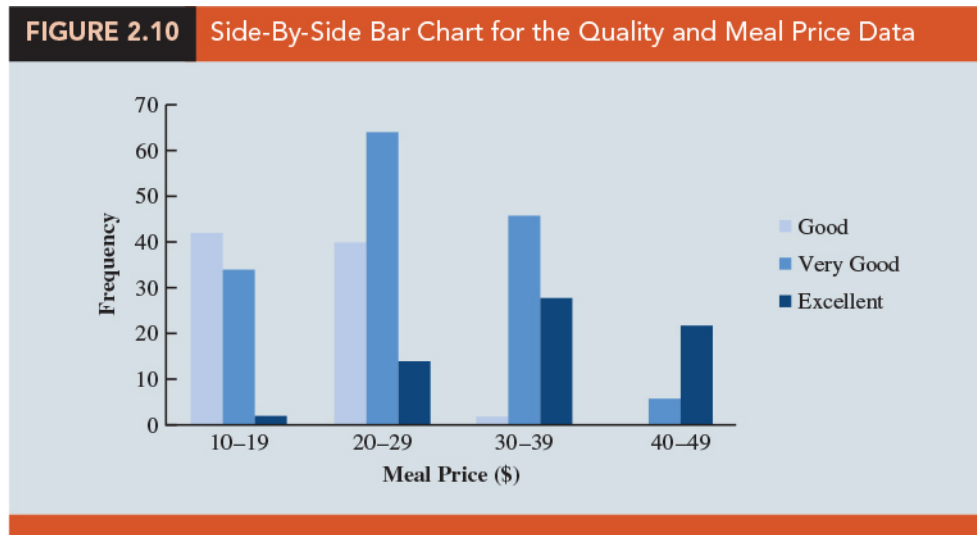
(c) The relationship is not perfect in that all points are not on a straight line. However, the general pattern of the points and the trendline suggest that the overall relationship is positive.

5. Some general scatter diagram patterns and the types of relationships they suggest are shown in Figure 2.9.

FIGURE 2.9 Types of Relationships Depicted by Scatter Diagrams

Side-by-Side Bar Chart

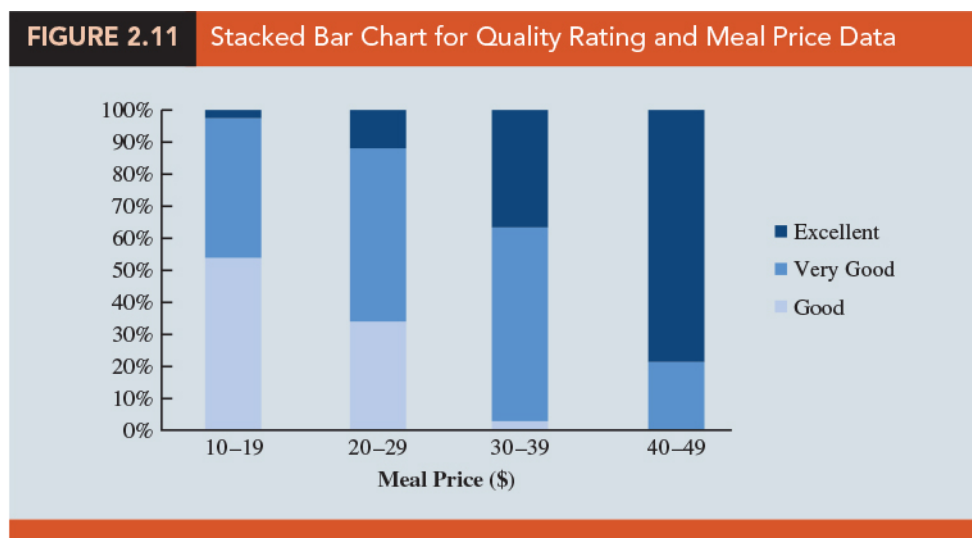
1. A side-by-side bar chart is a graphical display for depicting multiple bar charts on the same display.
2. Each cluster of bars represents one value of the first variable.
3. Each bar within a cluster represents one value of the second variable.
4. **Example** Quality Rating and Meal Price Data (Figure 2.10.)
 - (a) We see that the lowest meal price category (\$10–\$19) received mostly good and very good ratings, but very few excellent ratings. The highest price category (\$40–\$49) received mostly excellent ratings, some very good ratings, but no good ratings.
 - (b) Notice that as the price increases (left to right), the height of the light blue bars decreases and the height of the dark blue bars generally increases. This indicates that as price increases, the quality rating tends to be better.
 - (c) The very good rating, as expected, tends to be more prominent in the middle price categories as indicated by the dominance of the middle bar in the moderate price ranges of the chart.



Stacked Bar Chart

1. A stacked bar chart is a bar chart in which each bar is broken into rectangular segments of a different color.
2. If percentage frequencies are displayed, all bars will be of the same height (or length), extending to the 100% mark.
3. **Example** Quality Rating and Meal Price Data (Table 2.15.)(Figure 2.11.)
 - (a) Table 2.15 shows the column percentages for each meal price category.
 - (b) Figure 2.11 shows even more clearly than Figure 2.10 the relationship between the variables. As we move from the low price category (\$10–19) to the high price category (\$40–49), the length of the light blue bars decreases and the length of the dark blue bars increases.

Quality Rating	Meal Price			
	\$10–19	\$20–29	\$30–39	\$40–49
Good	53.8%	33.9%	2.6%	.0%
Very Good	43.6	54.2	60.5	21.4
Excellent	2.6	11.9	36.8	78.6
Total	100.0%	100.0%	100.0%	100.0%



2.5 Data Visualization: Best Practices in Creating Effective Graphical Displays

1. Data visualization is a term used to describe the use of graphical displays to summarize and present information about a data set.
2. The goal of data visualization is to communicate as effectively and clearly as possible, the key information about the data.

Creating Effective Graphical Displays

1. Give the display a clear and concise title.
2. Keep the display simple. Do not use three dimensions when two dimensions are sufficient.
3. Clearly label each axis and provide the units of measure.
4. If color is used to distinguish categories, make sure the colors are distinct.
5. If multiple colors or line types are used, use a legend to define how they are used and place the legend close to the representation of the data.

Choosing the Type of Graphical Display

1. **Displays used to show the distribution of data:**
 - (a) Bar Chart to show the frequency distribution and relative frequency distribution for categorical data
 - (b) Pie Chart to show the relative frequency and percent frequency for categorical data
 - (c) Dot Plot to show the distribution for quantitative data over the entire range of the data
 - (d) Histogram to show the frequency distribution for quantitative data over a set of class intervals

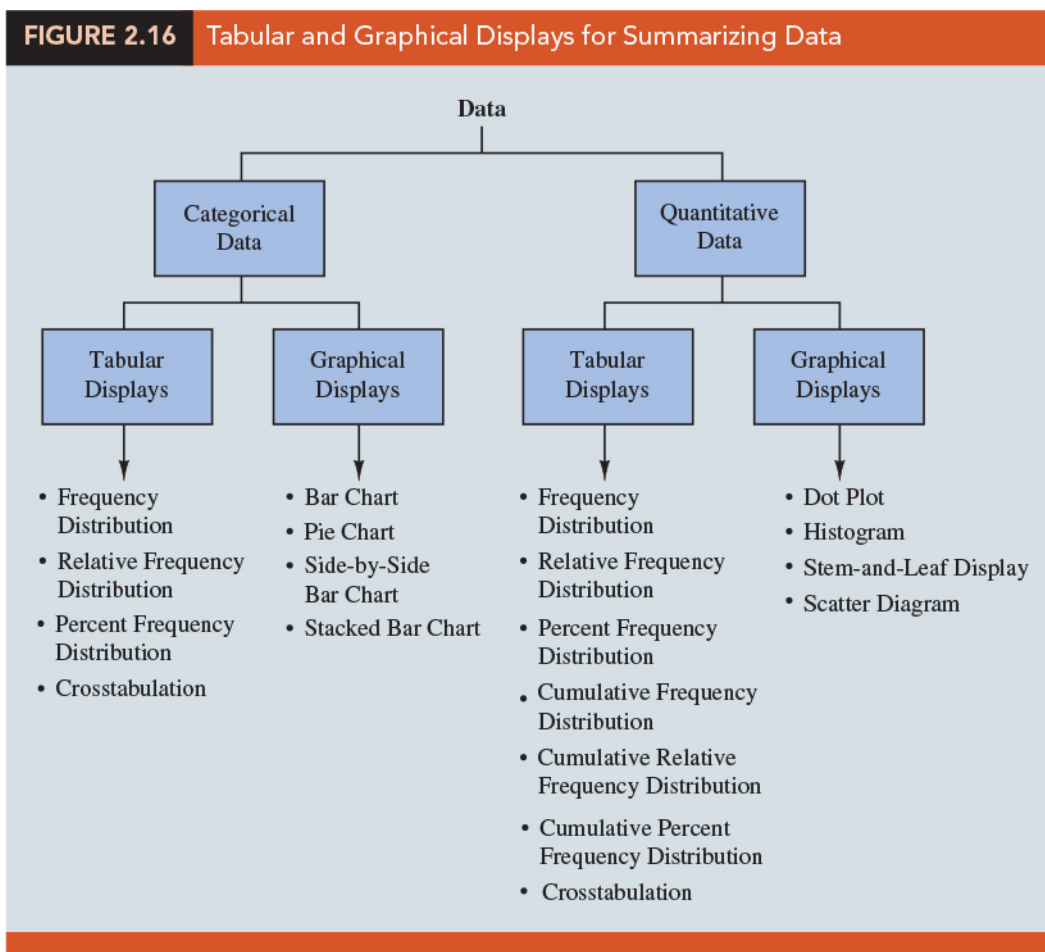
- (e) Stem-and-Leaf Display to show both the rank order and shape of the distribution for quantitative data

2. Displays used to make comparisons:

- (a) Side-by-Side Bar Chart to compare two variables
- (b) Stacked Bar Chart to compare the relative frequency or percent frequency of two categorical variables

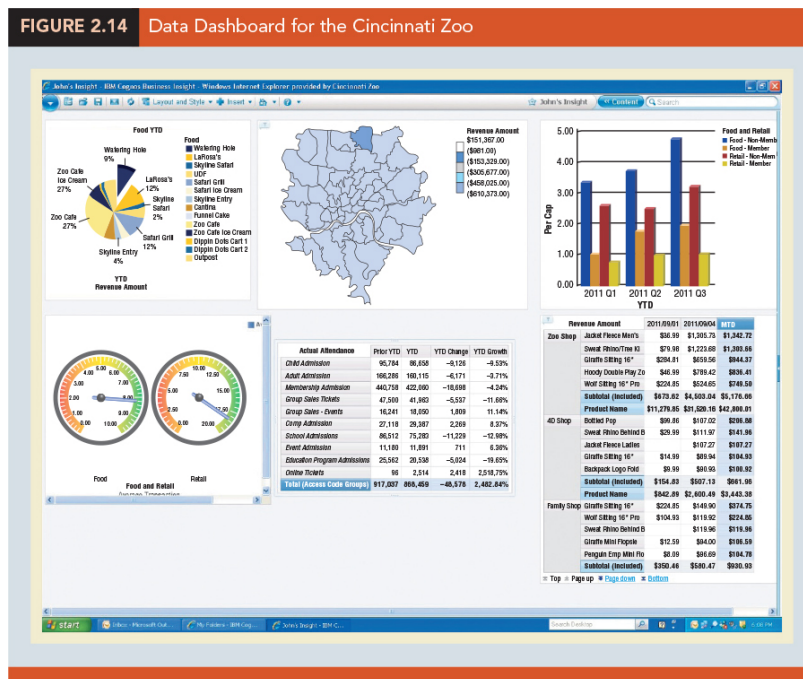
3. Displays used to show relationships:

- (a) Scatter Diagram to show the relationship between two quantitative variables
- (b) Trendline to approximate the relationship of data in a scatter diagram



Data Dashboards

1. A data dashboard is a widely used data visualization tool.
2. It organizes and presents key performance indicators (KPIs) used to monitor an organization or process.
3. It provides timely summary information that is easy to read, understand, and interpret.
4. Some additional guidelines include:
 - (a) Minimize the need for screen scrolling.
 - (b) Avoid unnecessary use of color or 3D displays.
 - (c) Use borders between charts to improve readability.



NOTE : use "misleading graphs" as the search keyword on Google or YouTube.

☺ EXERCISES

2.1 : 3, 4, 5, 10

2.2 : 11, 17, 20, 25

2.3 : 31, 33, 34

2.4 : 39, 41, 43

2.5 : 49, 54, 55, 56

“因為每個人在某些方面都會發生失敗，所以我能夠接受失敗，但我不能接受沒有嘗試”

“I can accept failure, everyone fails at something. But I can't accept not trying”

— *Michael Jordan (February 17, 1963 –)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 3: Descriptive Statistics: Numerical Measures

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

Overview

1. Develop numerical summary measures for data sets consisting of a single variable: location, variability, and distribution shape.
2. Other topics: relative location, detecting outliers, five-number summaries and box plots, the measures of association between two variables,
3. Sample statistics : if the measures are computed for data from a sample.
4. Population parameters : if the measures are computed for data from a population.
5. A sample statistic is referred to as the point estimator of the corresponding population parameter.

3.1 Measures of Location

Mean

1. The mean of a data set is the average of all the data values (x_1, x_2, \dots, x_n) .
2. The sample mean $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ is the point estimator of the population mean $\mu = \frac{\sum_{i=1}^N X_i}{N}$.
3. The mean provides a measure of central location.

Weighted Mean

1. In some instances, the mean is computed by giving each observation a weight $(w_i, i = 1, \dots, n)$ that reflects its relative importance.
2. The choice of weights depends on the application.
3. The weighted sample mean

$$\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$$

4. **Example** The computation of a grade point average (GPA) for college students.
 - (a) The data values generally used are 4 for an A grade, 3 for a B grade, 2 for a C grade, 1 for a D grade, and 0 for an F grade.
 - (b) The weights are the number of credit hours earned for each grade.

Median

1. The median of a data set is the value in the middle when the data items are arranged in ascending order (smallest value to largest value).
 - (a) For an odd number of observations, the median is the middle value.
 - (b) For an even number of observations, the median is the average of the two middle values.
2. Whenever a data set has extreme values, the median is the preferred measure of central location than mean.
3. The median is the measure of location most often reported for annual income and property value data.


Geometric Mean

1. The geometric mean is calculated by finding the n th root of the product of n values.

$$\bar{x}_g = \sqrt[n]{x_1 x_2 \cdots x_n}$$

2. It is often used in analyzing growth rates in financial data (where using the arithmetic mean will provide misleading results).

3. The sample mean is only appropriate for an additive process. For a multiplicative process, such as applications involving growth rates, the geometric mean is the appropriate measure of location.
4. While the applications of the geometric mean to problems in finance, investments, and banking are particularly common, the geometric mean should be applied any time you want to determine the mean rate of change over several successive periods.

 **Question** (p113)

Consider Table 3.2, which shows the percentage annual returns, or growth rates, for a mutual fund over the past 10 years. (a) Compute how much \$100 invested in the fund at the beginning of year 1 would be worth at the end of year 10. (b) What was the mean percentage annual return or mean rate of growth for this investment over the 10-year period?

Year	Return (%)	Growth Factor
1	-22.1	.779
2	28.7	1.287
3	10.9	1.109
4	4.9	1.049
5	15.8	1.158
6	5.5	1.055
7	-37.0	.630
8	26.5	1.265
9	15.1	1.151
10	2.1	1.021

sol:

Mode


1. The mode of a data set is the value that occurs with greatest frequency.
2. The greatest frequency can occur at two or more different values.
3. If the data have exactly two modes, the data are bimodal.
4. If the data have more than two modes, the data are multimodal.

Percentiles

1. A percentile provides information about how the data are spread over the interval from the smallest value to the largest value.
2. The p th percentile of a data set is a value such that at least p percent of the items take on this value or less and at least $(100 - p)\%$ of the items take on this value or more.
3. To calculate the p th percentile for a data set containing n observations:
 - (a) Arrange the data in ascending order.
 - (b) Compute L_p , the location of the p th percentile

$$L_p = \frac{p}{100}(n + 1)$$

4. The 50th percentile is also the median.
5. Admission test scores for colleges and universities are frequently reported in terms of percentiles.

 Question (p116)


Compute and interpret the 80th percentile for the starting salary data in Table 3.1.

TABLE 3.1 Monthly Starting Salaries for a Sample of 12 Business School Graduates			
Graduate	Monthly Starting Salary (\$)	Graduate	Monthly Starting Salary (\$)
1	5850	7	5890
2	5950	8	6130
3	6050	9	5940
4	5880	10	6325
5	5755	11	5920
6	5710	12	5880

sol:

Quartiles

1. Quartiles are specific percentiles .
2. Q_1 : First Quartile = 25th Percentile
3. Q_2 : Second Quartile = 50th Percentile = Median
4. Q_3 : Third Quartile = 75th Percentile

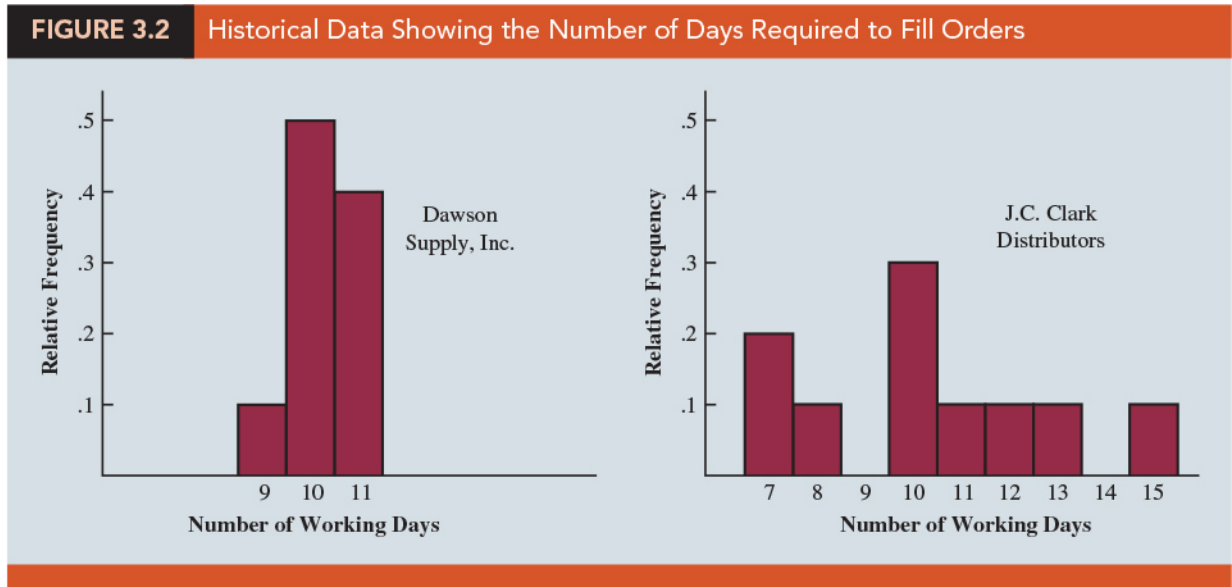
 Question (p116)

Compute the first and third quartiles for the starting salary data in Table 3.1.

sol:

3.2 Measures of Variability

1. It is often desirable to consider measures of variability (dispersion), as well as measures of location.
2. For example, in choosing supplier A or supplier B we might consider not only the average delivery time for each, but also the variability in delivery time for each.



Range

1. The range of a data set is the difference between the largest and smallest data value.

$$\text{Range} = \underline{\text{largest value} - \text{smallest value}}$$

2. It is very sensitive to the smallest and largest data values.

Interquartile Range (IQR)

1. The interquartile range of a data set is the difference between the third quartile and the first quartile: $IQR = Q_3 - Q_1$.
2. It is the range for the middle 50% of the data.
3. It overcomes the sensitivity to extreme data values.

Variance

1. The variance is a measure of variability that utilizes all the data.
2. deviation: the difference between the value of each observation (x_i) and the mean (\bar{x} for a sample, μ for a population): $(x_i - \bar{x}, x_i - \mu)$.
3. For any data set, the sum of the deviations about the mean will always equal zero: $\sum(x_i - \bar{x}) = 0$.
4. The variance is the average of the squared deviations between each data value and the mean.
5. The variance for a population is: $\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N}$.
6. The variance of a sample is: $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$.

Standard Deviation

1. The standard deviation of a data set is the positive square root of the variance.
2. It is measured in the same units as the original data, making it more easily interpreted than the variance.
3. The standard deviation of a sample is: $s = \sqrt{s^2}$.
4. The standard deviation of a population is: $\sigma = \sqrt{\sigma^2}$.

Coefficient of Variation

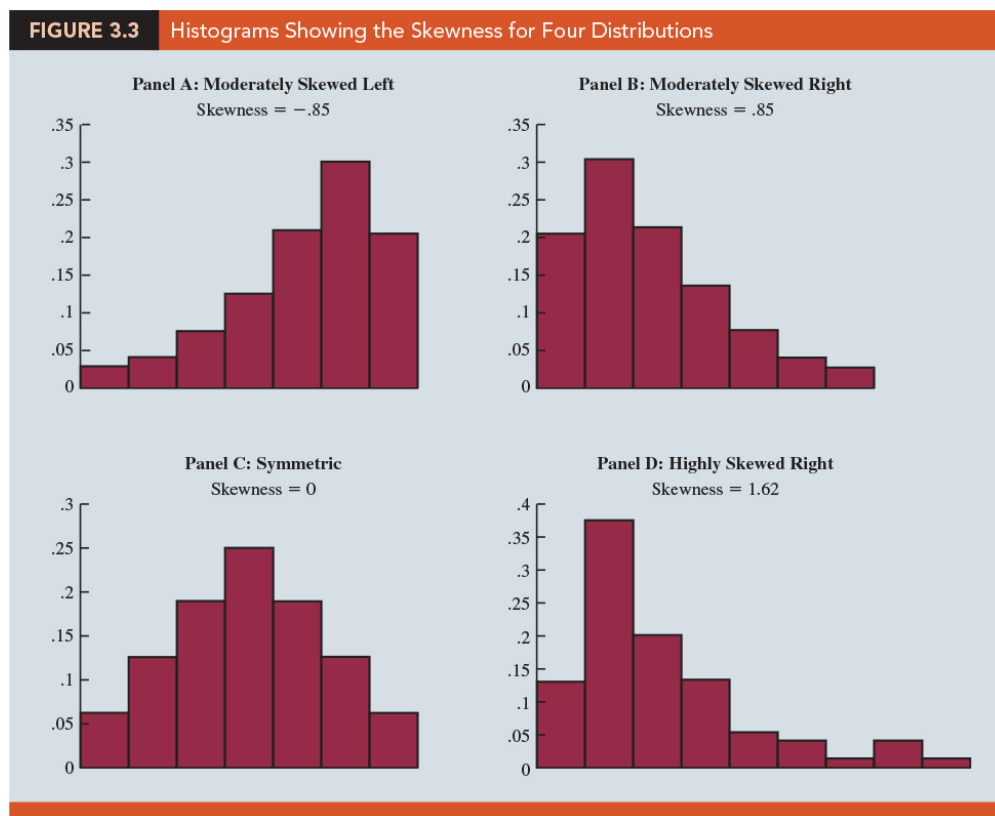
1. The coefficient of variation indicates how large the standard deviation is relative to the mean.
2. The coefficient of variation = $\frac{\text{standard deviation}}{\text{mean}} \times 100\%$.
3. In general, the coefficient of variation is a useful statistic for comparing the variability of variables that have different standard deviations and different means.

3.3 Measures of Distribution Shape, Relative Location, and Detecting Outliers

Distribution Shape

1. A histogram provides a graphical display showing the shape of a distribution.
2. An important numerical measure of the shape of a distribution is called skewness.
3. The formula for the skewness of sample data is

$$\text{Skewness} = \frac{n}{(n-1)(n-2)} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^3$$



4. For data skewed to the left, the skewness is negative; for data skewed to the right, the skewness is positive. If the data are symmetric, the skewness is zero.

5. For a symmetric distribution, the mean and the median are equal. When the data are positively (negatively) skewed, the mean will usually be greater (less) than the median.

***z*-Scores**

1. The *z*-score is often called the standardized value: $z_i = \frac{x_i - \bar{x}}{s}$.
2. It denotes the number of standard deviations a data value x_i is from the mean.
3. An observation's *z*-score is a measure of the relative location of the observation in a data set.
4. A data value equal to the sample mean will have a *z*-score of zero.
5. A data value greater (less) than the sample mean will have a *z*-score greater (less) than zero.

Chebyshev's Theorem

1. At least $(1-1/z^2)$ of the data values must be within z standard deviations of the mean, where z is any value greater than 1.
2. Chebyshev's theorem enables us to make statements about the proportion of data values that must be within a specified number of standard deviations of the mean.
3. Chebyshev's theorem requires $z > 1$; but z need not be an integer.
4. At least $(75\%, 89\%, 94\%)$ of the data values must be within $(z = 2, z = 3, z = 4)$ standard deviations of the mean.

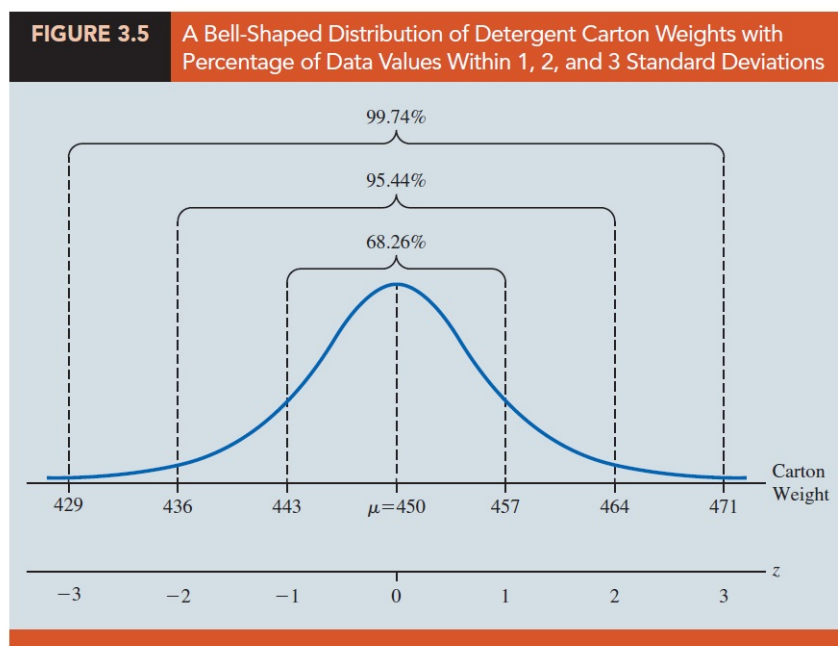
 **Question** (p132)

Suppose that the midterm test scores for 100 students in a college business statistics course had a mean of 70 and a standard deviation of 5. How many students had test scores between 60 and 80? How many students had test scores between 58 and 82? (Using Chebyshev's theorem)

sol:

Empirical Rule

1. When the data are believed to approximate a bell-shaped distribution (e.g., normal distribution), the empirical rule can be used to determine the percentage of data values that must be within a specified number of standard deviations of the mean.
2. Approximately (68%, 95%, 99%) of the data values will be within (one, two, three) standard deviation of the mean.



 Question (p134)

The liquid detergent cartons are filled automatically on a production line. Filling weights frequently have a bell-shaped distribution. If the mean filling weight is 16 ounces and the standard deviation is .25 ounces, how many filled cartons will (a) weights between 16 and 16.25 ounces? (b) weights between 15.50 and 16 ounces? (c) weights less than 15.50 ounces? (d) weights between 15.50 and 16.25 ounces?

sol:

Detecting Outliers

1. An outlier is an unusually small or unusually large value in a data set.
2. A data value with a z-score less than -3 or greater than +3 might be considered an outlier.
3. Another approach to identifying outliers is based upon the values of the first and third quartiles (Q_1 and Q_3) and the interquartile range (IQR).


$$\text{Lower Limit} = \underline{Q_1 - 1.5(IQR)}, \quad \text{Upper Limit} = \underline{Q_3 + 1.5(IQR)}$$

4. An observation is classified as an outlier if its value is less than the Lower limit or greater than the Upper limit.

3.4 Five-Number Summaries and Box Plots

Five-Number Summary

1. Five-Number Summary: smallest value, first quartile (Q_1), median, third quartile (Q_3), and largest value.

 **Question** (p138)

The monthly starting salary data is shown in ascending order below,

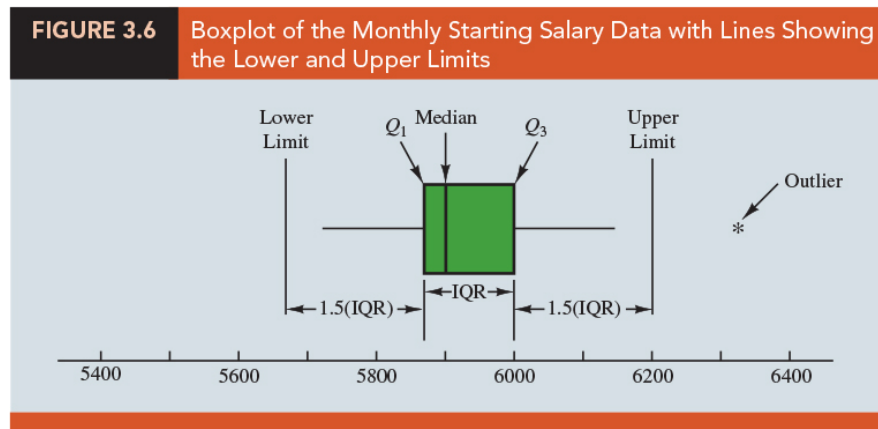
5710 5755 5850 5880 5880 5890 5920 5940 5950 6050 6130 6325

Find the five-number summary.

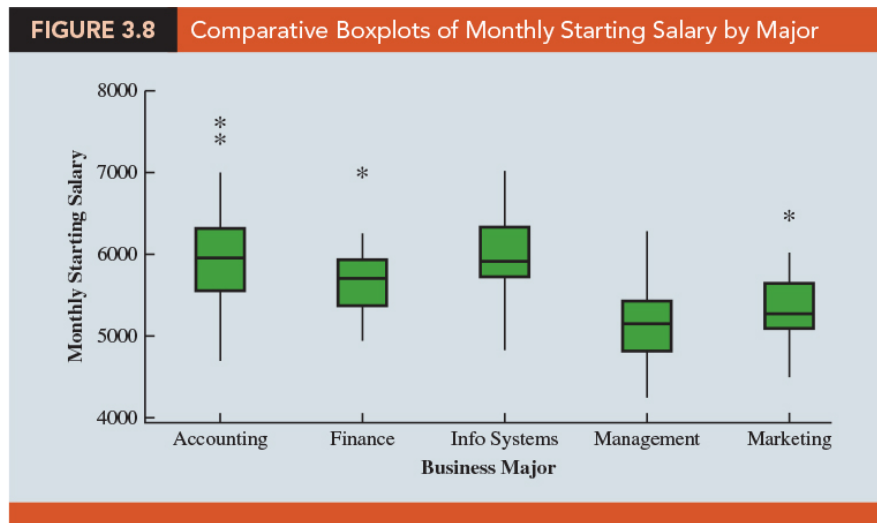
sol:

Box Plot

1. A box plot is a graphical display of data that is based on a five-number summary.



2. Box plots provide another way to identify outliers.
- (a) A box is drawn with its ends located at the first and third quartiles.
 - (b) A vertical line is drawn in the box at the location of the median (second quartile).
 - (c) Box limits are located using the interquartile range (IQR).
 - (d) Data outside the box $\pm 1.5IQR$ are considered outliers (shown with the symbol *).
 - (e) The horizontal lines extending from each end of the box in Figure 3.6 are called whiskers. The whiskers are drawn from the ends of the box to the smallest and largest values inside the limits.
3. Boxplots can also be used to provide a graphical summary of two or more groups and facilitate visual comparisons among the groups.



3.5 Measures of Association Between Two Variables

Covariance

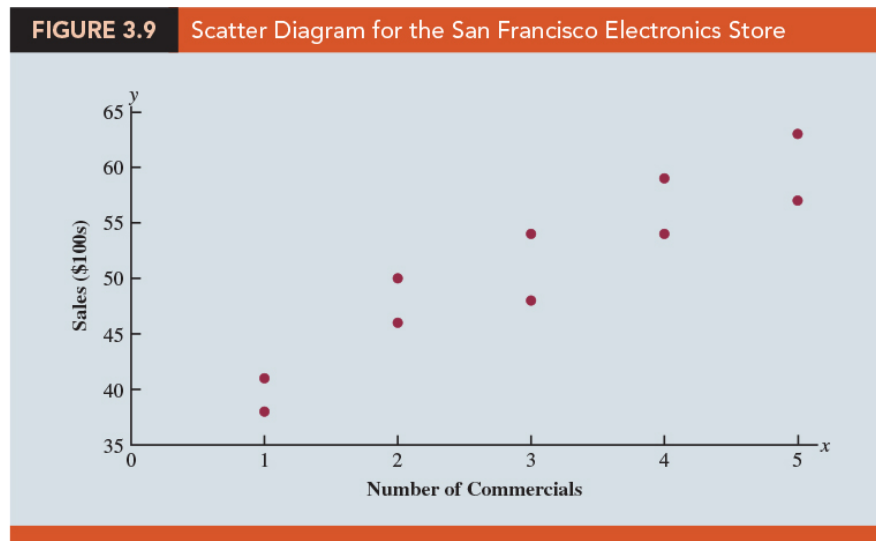
1. The covariance is a measure of the linear association between two variables. For a sample of size n with the observations $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, the sample covariance is defined as follows:

$$S_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

2. Let μ_x be the population mean of the variable x and μ_y be the population mean of the variable y . The population covariance σ_{xy} is defined for a population of size N :

$$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$

3. Positive (Negative) values of covariance indicate a positive (negative) relationship.



 Question (p143)

Compute the sample covariance to measure the strength of the linear relationship between the number of commercials x and the sales volume y in the San Francisco electronics store problem.

sol:

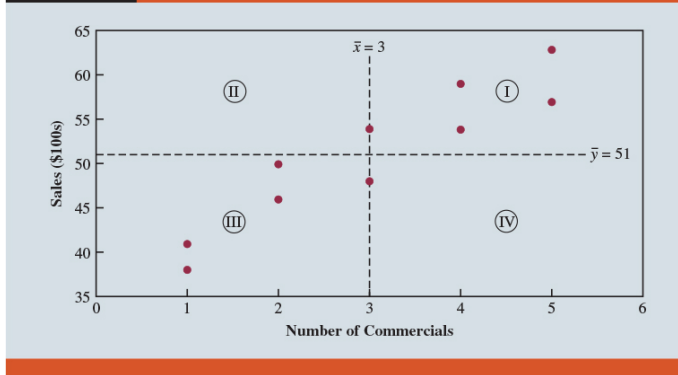
	x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
	2	50	-1	-1	1
	5	57	2	6	12
	1	41	-2	-10	20
	3	54	0	3	0
	4	54	1	3	3
	1	38	-2	-13	26
	5	63	2	12	24
	3	48	0	-3	0
	4	59	1	8	8
	2	46	-1	-5	5
Totals	$\overline{30}$	$\overline{510}$	$\overline{0}$	$\overline{0}$	$\overline{99}$

$$s_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n - 1} = \frac{99}{10 - 1} = 11$$

Interpretation of the Covariance

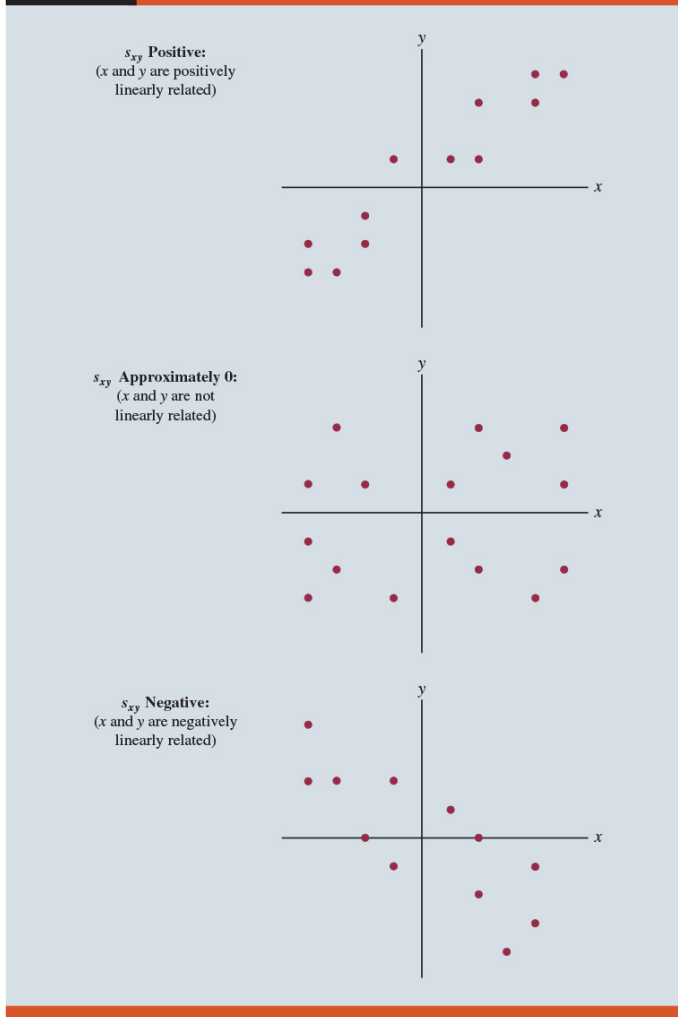
- (Figure 3.10) Points in quadrant I correspond to x_i greater than \bar{x} and y_i greater than \bar{y} .
- If the value of s_{xy} is positive (negative), the points with the greatest influence on s_{xy} must be in quadrants I and III (II and IV). Hence, a positive (negative) value for s_{xy} indicates a positive (negative) linear association between x and y ; that is, as the value of x increases (decreases), the value of y increases (decreases).
- The points are evenly distributed across all four quadrants, the value of s_{xy} will be close to zero, indicating no linear association between x and y .

FIGURE 3.10 Partitioned Scatter Diagram for the San Francisco Electronics Store



4. (Figure 3.11) the values of s_{xy} can be expected with three different types of scatter diagrams.

FIGURE 3.11 Interpretation of Sample Covariance



- However, one problem with using covariance as a measure of the strength of the linear relationship is that the value of the covariance depends on the units of measurement for x and y .
- A measure of the relationship between two variables that is not affected by the units of measurement for x and y is the correlation coefficient.

Correlation Coefficient

- Correlation is a measure of linear association and not necessarily causation.
- The sample Pearson product moment correlation coefficient:

$$r_{xy} = \frac{s_{xy}}{s_x s_y},$$

where s_x (s_y) is the sample standard deviation of x (y).

- The population Pearson product moment correlation coefficient:

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y},$$

where σ_x (σ_y) is the population standard deviation of x (y).

- The coefficient can take on values between -1 and +1.
- Values near -1 (+1) indicate a strong negative (positive) linear relationship.
- The closer the correlation is to zero, the weaker the relationship.
- NOTE: A high correlation between two variables does not mean that changes in one variable will cause changes in the other variable.

 **Question** (p147)

Let $(x_i, y_i), i = 1, 2, 3$ be the $(5, 10), (10, 30), (15, 50)$. Compute the sample correlation between x and y .

sol:

TABLE 3.8 Computations Used in Calculating the Sample Correlation Coefficient

	x_i	y_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$y_i - \bar{y}$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
	5	10	-5	25	-20	400	100
	10	30	0	0	0	0	0
	15	50	5	25	20	400	100
Totals	$\underline{30}$	$\underline{90}$	$\underline{0}$	$\underline{50}$	$\underline{0}$	$\underline{800}$	$\underline{200}$
	$\bar{x} = 10$	$\bar{y} = 30$					

3.6 Data Dashboards: Adding Numerical Measures to Improve Effectiveness*

☺ EXERCISES

3.1 : 5, 6, 10

3.2 : 25, 32, 34

3.3 : 38, 41, 43, 45

3.4 : 46, 47, 50

3.5 : 55, 59, 60

SUP : 64, 67, 71

“學習知識要善於思考，思考，再思考，我就是這個方法成為科學家的”

“Try not to become a man of success, but rather try to become a man of value”

— *Albert Einstein (March 14, 1879 – April 18, 1955)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 4: Introduction to Probability

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

Overview

1. **Example** Managers often base their decisions on an analysis of uncertainties such as the following:
 - (a) What are the chances that the sales will decrease if we increase prices?
 - (b) What is the likelihood a new assembly method will increase productivity?
 - (c) How likely is it that the project will be finished on time?
 - (d) What is the chance that a new investment will be profitable?
2. Probability:
 - (a) Probability is a numerical measure of the likelihood that an event will occur.
 - (b) Probability values are always assigned on a scale from 0 to 1.
 - (c) A probability near zero indicates an event is quite unlikely to occur.
 - (d) A probability near one indicates an event is almost certain to occur.

4.1 Random Experiments, Counting Rules, and Assigning Probabilities

1. **Random experiments:** a random experiment is a process that generates well-defined

experimental outcomes. On any single repetition or trial, the outcome that occurs is determined completely by chance.

2. **Sample space S** : the sample space for a random experiment is the set of all experimental outcomes.

3. An experimental outcome is also called a sample point to identify it as an element of the sample space.

4. **Example** Tossing a Coin

(a) Tossing the coin the upward face will be either a head (H) or a tail (T). There are two possible experimental outcomes: head or tail (sample points).

(b) Each time we toss the coin we will either observe a head or a tail. And, the outcome that occurs on any trial is determined solely by chance or random variability.

(c) By specifying all the possible experimental outcomes, we identify the sample space $S = \{Head, Tail\}$ or $S = \{H, T\}$ for a random experiment.

5. **Example** Rolling a Die

(a) The experimental outcomes, defined as the number of dots appearing on the face of the die, are the six sample points in the sample space $S = \{1, 2, 3, 4, 5, 6\}$ for this random experiment.

Counting Rules, Combinations, and Permutations

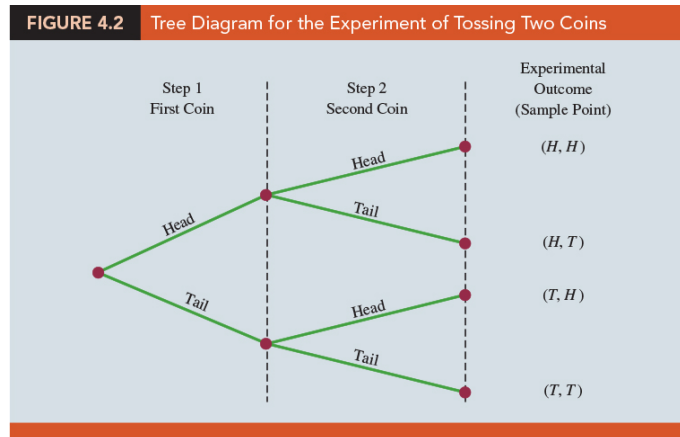
1. Being able to identify and count the experimental outcomes is a necessary step in assigning probabilities.

2. **Counting Rule for Multiple-Step Experiments**: If an experiment can be described as a sequence of k steps with n_1 possible outcomes on the first step, n_2 possible outcomes on the second step, and so on, then the total number of experimental outcomes is given by $(n_1)(n_2) \cdots (n_k)$.

3. **Tree Diagram**: a graphical representation that helps in visualizing a multiple-step experiment.

4. Example Tossing two Coins

- (a) The experiment of tossing two coins can be thought of as a two-step experiment in which step 1 is the tossing of the first coin $(n_1 = 2)$ and step 2 is the tossing of the second coin $(n_2 = 2)$.
- (b) The sample space for this coin-tossing experiment $S = \{(H, H), (H, T), (T, H), (T, T)\}$, 4 distinct experimental outcomes are possible.



Question (p182)

The Kentucky Power & Light Company (KP&L) is starting a project designed to increase the generating capacity of one of its plants in northern Kentucky. The project is divided into two sequential stages or steps: stage 1 (design) and stage 2 (construction). An analysis of similar construction projects revealed possible completion times for the design stage of 2, 3, or 4 months and possible completion times for the construction stage of 6, 7, or 8 months. Management sets a goal of 10 months for the completion of the entire project. Summarizes the experimental outcomes to complete the entire project for the KP&L problem. Draw the tree diagram to show how the outcomes (sample points) occur.

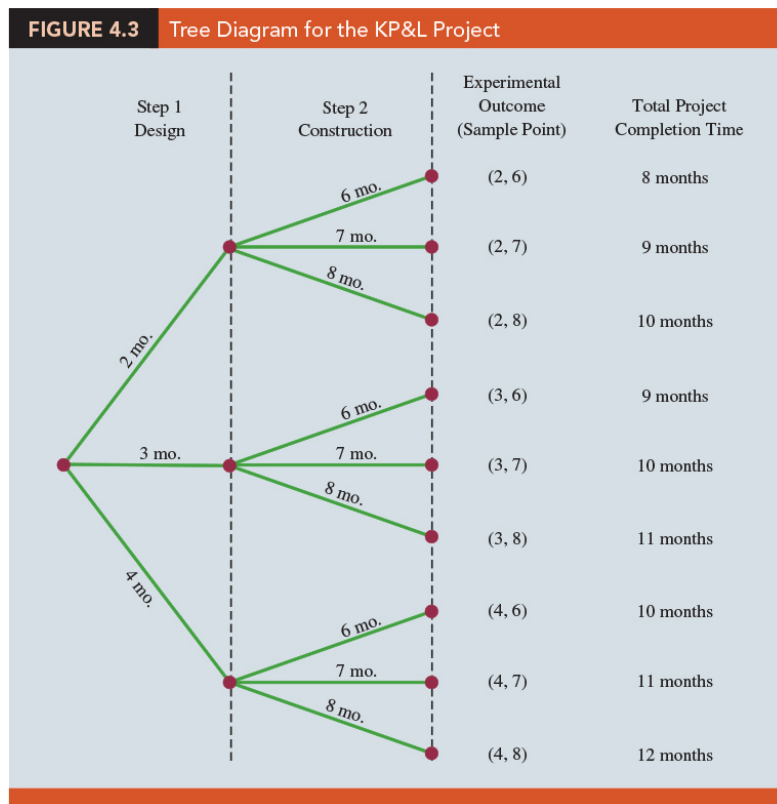
sol:

Notation: let (a, b) indicates that the design stage is completed in a months and the construction stage is completed in b months. Table 4.1: summarizes the experimental outcomes.

TABLE 4.1 Experimental Outcomes (Sample Points) for the KP&L Project

Completion Time (months)		Notation for Experimental Outcome	Total Project Completion Time (months)
Stage 1 Design	Stage 2 Construction		
2	6	(2, 6)	8
2	7	(2, 7)	9
2	8	(2, 8)	10
3	6	(3, 6)	9
3	7	(3, 7)	10
3	8	(3, 8)	11
4	6	(4, 6)	10
4	7	(4, 7)	11
4	8	(4, 8)	12

Figure 4.3: the project will be completed in 8 to 12 months, with six of the nine experimental outcomes providing the desired completion time of 10 months or less.



5. **Counting Rule for Combinations:** count the number of experimental outcomes when the experiment involves selecting n objects from a set of N objects.
6. The number of combinations of N objects taken n at a time is

$$C_n^N = \binom{N}{n} = \frac{N!}{n!(N-n)!} ,$$

where $N! = N(N-1)(N-2)\cdots(2)(1)$, $n! = n(n-1)(n-2)\cdots(2)(1)$, and $0! = 1$.

7. **Example** Consider a quality control procedure in which an inspector randomly selects two ($n = 2$) of five parts ($N = 5$) to test for defects. In a group of five parts, there are 10 combinations of two parts can be selected:

$$C_2^5 = \binom{5}{2} = \frac{5!}{2!(5-2)!} = 10.$$

If we label the five parts as A, B, C, D, and E, the 10 combinations or experimental outcomes can be identified as AB, AC, AD, AE, BC, BD, BE, CD, CE, and DE.

8. **Counting Rule for Permutations:** the number of experimental outcomes when n objects are to be selected from a set of N objects where the order of selection is important. The same n objects selected in a different order are considered a different experimental outcome.
9. The number of permutations of N objects taken n at a time is given by

$$P_n^N = n! \binom{N}{n} = \frac{N!}{(N-n)!}$$

10. **Example** Consider again the quality control process in which an inspector selects two ($n = 2$) of five parts ($N = 5$) to inspect for defects. There are 20 permutations may be selected:

$$P_2^5 = \frac{5!}{(5-2)!} = \frac{120}{6} = 20.$$

If we label the parts A, B, C, D, and E, the 20 permutations are AB, BA, AC, CA, AD, DA, AE, EA, BC, CB, BD, DB, BE, EB, CD, DC, CE, EC, DE, and ED.

Assigning Probabilities

1. **Basic Requirements for Assigning Probabilities:** the probability assigned to each experimental outcome must be between 0 and 1, inclusively.

(a) Let E_i is the i th experimental outcome and $P(E_i)$ is its probability

$$\underline{0 \leq P(E_i) \leq 1}, \quad \text{for all } i$$

(b) The sum of the probabilities for all experimental outcomes must equal 1.

$$\underline{\sum_{i=1}^n P(E_i) = 1}$$

where n is the number of experimental outcomes.

2. **Classical Method:** assigning probabilities based on the assumption of equally likely outcomes.
3. **Example** Consider the experiment of tossing a fair coin; the two experimental outcomes—head and tail—are equally likely. Because one of the two equally likely outcomes is a head, the probability of observing a head is $1/2$. Similarly, the probability of observing a tail is also $1/2$.
4. **Relative Frequency Method:** assigning probabilities based on experimental or historical data. That is when data are available to estimate the proportion of the time the experimental outcome will occur if the experiment is repeated a large number of times.
5. **Example** Consider a study of waiting times in the X-ray department for a local hospital. A clerk recorded the number of patients waiting for service at 9:00 a.m. on 20 successive days and obtained the following results.

Number Waiting: 0, 1, 2, 3, 4

Number of Days Outcome Occured: 2, 5, 6, 4, 3

These data show that on 2 of the 20 days, zero patients were waiting for service; on 5 of the days, one patient was waiting for service; and so on. Using the relative frequency method, we would assign a probability of $2/20 = 0.10$ ($5/20 = 0.25$) ($6/20 = 0.30$) ($4/20 = 0.20$) ($3/20 = 0.15$) to the experimental outcome of zero (one) (two) (three) (four) patients waiting for service. As with the classical method, using the relative frequency method automatically satisfies the two basic requirements.

6. **Subjective Method:** assigning probability based on judgment. That is, a probability value that expresses our degree of belief (on a scale from 0 to 1) that the experimental outcome will occur is specified.

7. **Example** Consider the case in which Tom and Judy make an offer to purchase a house. Two outcomes are possible:

E_1 = their offer is accepted

E_2 = their offer is rejected

Judy believes that the probability their offer will be accepted is 0.8; thus, Judy would set $P(E_1) = 0.8$ and $P(E_2) = 0.2$. Tom, however, believes that the probability that their offer will be accepted is 0.6; hence, Tom would set $P(E_1) = 0.6$ and $P(E_2) = 0.4$. Note that Tom's probability estimate for E_1 reflects a greater pessimism that their offer will be accepted. Both Judy and Tom assigned probabilities that satisfy the two basic requirements. The fact that their probability estimates are different emphasizes the personal nature of the subjective method.

8. The best probability estimates often are obtained by combining the estimates from the classical or relative frequency approach with the subjective estimate.

 **Question** (p113)

To perform further analysis on the KP&L project, we must develop probabilities for each of the nine experimental outcomes listed in Table 4.1. On the basis of experience and judgment, management concluded that the experimental outcomes were not equally likely. Hence, the classical method of assigning probabilities could not be used. Management then decided to conduct a study of the completion times for similar projects undertaken by KP&L over the past three years. The results of a study of 40 similar projects are summarized in Table 4.2. After reviewing the results of the study, management decided to employ the relative frequency method of assigning probabilities. Management could have provided subjective probability estimates but felt that the current project was quite similar to the 40 previous projects. Thus, the relative frequency method was judged best.

Completion Time (months)			Number of Past Projects Having These Completion Times
Stage 1 Design	Stage 2 Construction	Sample Point	
2	6	(2, 6)	6
2	7	(2, 7)	6
2	8	(2, 8)	2
3	6	(3, 6)	4
3	7	(3, 7)	8
3	8	(3, 8)	2
4	6	(4, 6)	2
4	7	(4, 7)	4
4	8	(4, 8)	6
Total			40

In using the data in Table 4.2 to compute probabilities for each outcome.

sol:

Let $P(a, b)$ represents the probability of the sample point (a, b) . We note that outcome $(2, 6)$ - stage 1 completed in 2 months and stage 2 completed in 6 months - occurred six times in the 40 projects. Use the relative frequency method to assign a probability of $6/40 = 0.15$ to this outcome. Continuing in this manner, we obtain the probability assignments for the sample points of the KP&L project shown in Table 4.3.

Sample Point	Project Completion Time	Probability of Sample Point
(2, 6)	8 months	$P(2, 6) = 6/40 = .15$
(2, 7)	9 months	$P(2, 7) = 6/40 = .15$
(2, 8)	10 months	$P(2, 8) = 2/40 = .05$
(3, 6)	9 months	$P(3, 6) = 4/40 = .10$
(3, 7)	10 months	$P(3, 7) = 8/40 = .20$
(3, 8)	11 months	$P(3, 8) = 2/40 = .05$
(4, 6)	10 months	$P(4, 6) = 2/40 = .05$
(4, 7)	11 months	$P(4, 7) = 4/40 = .10$
(4, 8)	12 months	$P(4, 8) = 6/40 = .15$
Total		1.00

4.2 Events and Their Probabilities

1. **Event:** an event is a collection (set) of sample points.
2. **Probability of an Event:** the probability of any event is equal to the sum of the probabilities of the sample points in the event.
3. If we can identify all the sample points of an experiment and assign a probability to each, we can compute the probability of an event.
4. Example KP&L project

(a) Assume that the project manager is interested in the event that the entire project can be completed in 10 months or less. Referring to Table 4.3, we see that six sample points - (2, 6), (2, 7), (2, 8), (3, 6), (3, 7), and (4, 6) - provide a project completion time of 10 months or less.

(b) Let C denote the event that the project is completed in 10 months or less:

$$C = \underline{\{(2, 6), (2, 7), (2, 8), (3, 6), (3, 7), (4, 6)\}}$$

Event C is said to occur if any one of these six sample points appears as the experimental outcome.

(c) Other events that might be of interest:


L = The event that the project is completed in less than 10 months

M = The event that the project is completed in more than 10 months

Using the information in Table 4.3, we see that these events consist of the following sample points.

$$L = \underline{\{(2, 6), (2, 7), (3, 6)\}}$$

$$M = \underline{\{(3, 8), (4, 7), (4, 8)\}}$$

 **Question** (p190)

In the KP&L project problem, compute the probability that the project will take 10 months or less to complete. Compute the probability of the event that the project is completed in more than 10 months.

sol: Using these probability results, we can now tell KP&L management that:
The probability that the project will be completed in 10 months or less:

$$\underline{P(C) = 0.7.}$$

The probability that the project will be completed in less than 10 months:

$$\underline{P(L) = 0.40.}$$

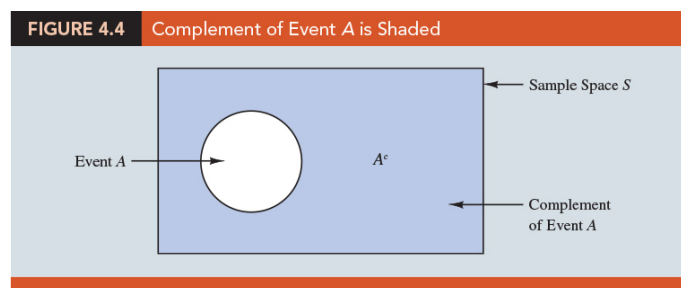
The probability that the project will be completed in more than 10 months:

$$\underline{P(M) = 0.30.}$$

4.3 Some Basic Relationships of Probability

Complement of an Event

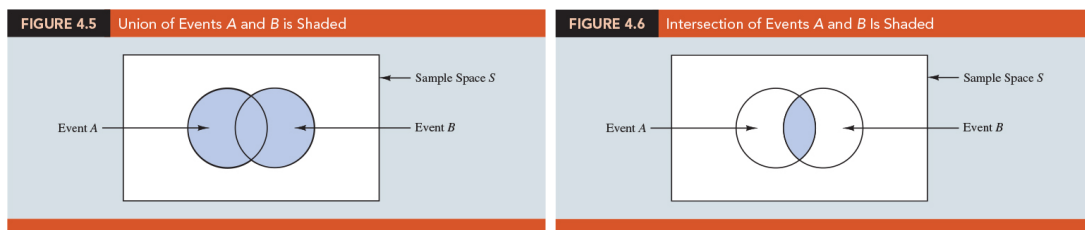
1. The complement of event A is defined to be the event consisting of all sample points that are not in A .
2. The complement of A is denoted by A^c .
3. Computing probability using the complement: $P(A) = 1 - P(A^c)$.
4. (Figure 4.4): Venn diagram illustrates the concept of a complement.



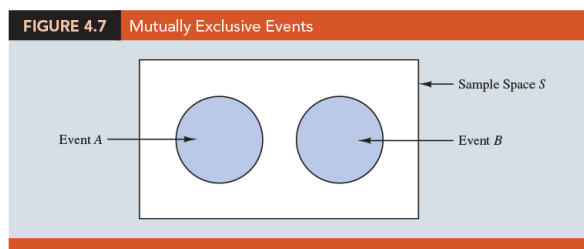
Addition Law

- Union:** given two events A and B , the union of A and B (denoted by $A \cup B$) is the event containing all sample points belonging to A or B or both.
- Intersection:** The intersection of events A and B (denoted by $A \cap B$) is the set of all sample points that are in both A and B .
- Addition Law:** a way to compute the probability of event A , or B , or both A and B occurring:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



- Two events are said to be mutually exclusive if the events have no sample points in common. Two events are mutually exclusive if, when one event occurs, the other cannot occur.
- Addition Law for Mutually Exclusive Events:** $P(A \cup B) = P(A) + P(B)$



Question (p195)

Consider the case of a small assembly plant with 50 employees. Each worker is expected to complete work assignments on time and in such a way that the assembled product will pass a final inspection. On occasion, some of the workers fail to meet the performance standards by completing work late or assembling a defective product. At the end of a performance evaluation period, the production manager

found that 5 of the 50 workers completed work late, 6 of the 50 workers assembled a defective product, and 2 of the 50 workers both completed work late and assembled a defective product. Let

L = the event that the work is completed late

D = the event that the assembled product is defective

After reviewing the performance data, the production manager decided to assign a poor performance rating to any employee whose work was either late or defective; thus the event of interest is $L \cup D$. What is the probability that the production manager assigned an employee a poor performance rating?

sol:

The relative frequency information:

$$\underline{P(L) = \frac{5}{50} = 0.10, P(D) = \frac{6}{50} = 0.12, P(L \cap D) = \frac{2}{50} = 0.04} .$$

To compute $P(L \cup D)$:

$$\begin{aligned} P(L \cup D) &= \frac{P(L) + P(D) - P(L \cap D)}{=} \\ &= \underline{0.10 + 0.12 - 0.04 = 0.18} \end{aligned}$$

There is a 0.18 probability that a randomly selected employee received a poor performance rating.

 Question (p196)

Consider a recent study conducted by the personnel manager of a major computer software company. The study showed that 30% of the employees who left the firm within two years did so primarily because they were dissatisfied with their salary, 20% left because they were dissatisfied with their work assignments, and 12% of the former employees indicated dissatisfaction with both their salary and their work assignments. What is the probability that an employee who leaves within two years does so because of dissatisfaction with salary, dissatisfaction with the work assignment, or both?

sol:

Let

$S =$ the event that the employee leaves because of salary

$W =$ the event that the employee leaves because of work assignment

We have $P(S) = 0.30$, $P(W) = 0.20$, and $P(S \cap W) = 0.12$.

By the addition law, we have

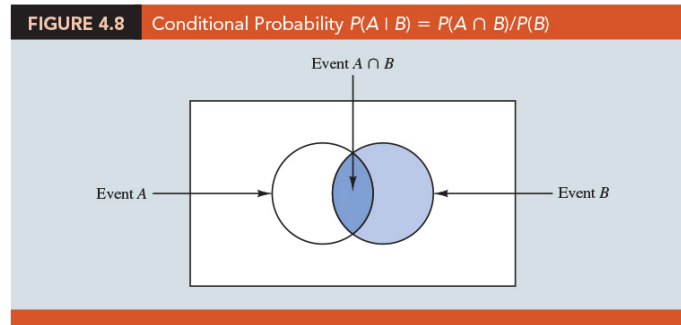
$$\underline{P(S \cup W) = P(S) + P(W) - P(S \cap W) = 0.30 + 0.20 - 0.12 = 0.38}.$$

We find a 0.38 probability that an employee leaves for salary or work assignment reasons.

4.4 Conditional Probability

1. **Conditional probability:** the probability of an event A given that another event B has occurred is called a conditional probability (denoted by $P(A|B)$):

$$P(A|B) = \underline{\frac{P(A \cap B)}{P(B)}}$$



2. **Joint probability:** The probability of the intersection of two events is called joint probability.
3. **Marginal probabilities:** The values in the margins of the joint probability table provide the probabilities of each event separately. The marginal probabilities are found by summing the joint probabilities in the corresponding row or column of the joint probability table.
4. **Example** Promotion status of male and female officers
 - (a) Consider the situation of the promotion status of male and female officers of a major metropolitan police force in the eastern United States. The police force consists of 1200 officers, 960 men and 240 women. Over the past two years, 324 officers on the police force received promotions. The specific breakdown of promotions for male and female officers is shown in Table 4.4.

TABLE 4.4 Promotion Status of Police Officers Over the Past Two Years

	Men	Women	Total
Promoted	288	36	324
Not Promoted	672	204	876
Total	960	240	1200

- (b) After reviewing the promotion record, a committee of female officers raised a discrimination case on the basis that 288 male officers had received promotions, but only 36 female officers had received promotions. The police administration argued that the relatively low number of promotions for female officers was due not to discrimination, but to the fact that relatively few females are members of the police force.

- (c) Let us show how conditional probability could be used to analyze the discrimination charge. Let

M = event an officer is a man, W = event an officer is a woman

A = event an officer is promoted, A^c = event an officer is not promoted

Obtain the following joint probability table:

TABLE 4.5 Joint Probability Table for Promotions

	Men (M)	Women (W)	Total
Promoted (A)	.24	.03	.27
Not Promoted (A^c)	.56	.17	.73
Total	.80	.20	1.00

Joint probabilities appear in the body of the table.

Marginal probabilities appear in the margins of the table.

- probability that a randomly selected officer is a man and is promoted:
$$\underline{P(M \cap A) = 288/1200 = 0.24}$$
 - probability that a randomly selected officer is a man and is not promoted:
$$\underline{P(M \cap A^c) = 672/1200 = 0.56}$$
 - probability that a randomly selected officer is a woman and is promoted:
$$\underline{P(W \cap A) = 36/1200 = 0.03}$$
 - probability that a randomly selected officer is a woman and is not promoted:
$$\underline{P(W \cap A^c) = 204/1200 = 0.17}$$
- (d) The marginal probabilities: $\underline{P(M) = 0.80, P(W) = 0.20, P(A) = 0.27}$, and $\underline{P(A^c) = 0.73}$. We see that 80% of the force is male, 20% of the force is female, 27% of all officers received promotions, and 73% were not promoted.
- (e) The conditional probability of being promoted given that the officer is $P(A|M)$:

$$\underline{P(A|M) = \frac{P(A \cap M)}{P(M)} = \frac{0.24}{0.80} = 0.30}$$

- (f) Similarly, the probability that an officer is promoted given that the officer is a woman; that is, $P(A|W)$:

$$\underline{P(A|W) = \frac{P(A \cap W)}{P(W)} = \frac{0.03}{0.20} = 0.15}$$

- (g) The marginal probability in row 1 of Table 4.5 shows that the probability of promotion of an officer is $P(A) = 0.27$ (regardless of whether that officer is male or female). However, the critical issue in the discrimination case involves the two conditional probabilities $P(A|M)$ and $P(A|W)$. That is, what is the probability of a promotion given that the officer is a man, and what is the probability of a promotion given that the officer is a woman? If these two probabilities are equal, a discrimination argument has no basis because the chances of a promotion are the same for male and female officers. However, a difference in the two conditional probabilities will support the position that male and female officers are treated differently in promotion decisions.
- (h) What conclusion do you draw?
- The probability of a promotion given that the officer is a man is 0.30, twice the 0.15 probability of a promotion given that the officer is a woman.
 - Although the use of conditional probability does not in itself prove that discrimination exists in this case, the conditional probability values support the argument presented by the female officers.

Independent Events

1. Two events A and B are independent if

$$\underline{P(A|B) = P(A)} \quad \text{or} \quad \underline{P(B|A) = P(B)}$$

Otherwise, the events are dependent.

2. Example Promotion status of male and female officers

- $P(A) = 0.27$, $P(A|M) = 0.30$, and $P(A|W) = 0.15$. We see that the probability of a promotion (event A) is affected or influenced by whether the officer is a man or a woman.
- Particularly, because $P(A|M) \neq P(A)$, we would say that events A and M are dependent events. That is, the probability of event A (promotion) is altered or affected by knowing that event M (the officer is a man) exists.
- Similarly, with $P(A|W) \neq P(A)$, we would say that events A and W are dependent events.

Multiplication Law

1. **Multiplication law:** a way to compute the probability of the intersection of two events:

$$\underline{P(A \cap B) = P(B)P(A|B)} \quad \text{and} \quad \underline{P(A \cap B) = P(A)P(B|A)}$$

2. **Example** Telecommunications company

- (a) Consider a telecommunications company that offers services such as high-speed Internet, cable television, and telephone services. For a particular city, it is known that 84% of the households subscribe to high-speed Internet service. If we let H denote the event that a household subscribes to high-speed Internet service, $P(H) = 0.84$. In addition, it is known that the probability that a household that already subscribes to high-speed Internet service also subscribes to cable television service (event C) is 0.75; that is, $P(C|H) = 0.75$. What is the probability that a household subscribes to both high-speed Internet and cable television services?
- (b) Using the multiplication law, we compute the desired $P(C \cap H)$ as

$$P(C \cap H) = \underline{P(H)P(C|H) = 0.84(0.75) = 0.63}$$

We now know that 63% of the households subscribe to both high-speed Internet and cable television services.

3. **Multiplication law for two independent events:** events A and B are independent whenever $P(A|B) = P(A)$ or $P(B|A) = P(B)$, the multiplication law is

$$\underline{P(A \cap B) = P(A)P(B)}$$

4. If $P(A \cap B) = P(A)P(B)$, then A and B are independent; if $\underline{P(A \cap B) \neq P(A)P(B)}$, then A and B are dependent.

5. **Example** Service station manager

- (a) Consider the situation of a service station manager who knows from past experience that 80% of the customers use a credit card when they purchase gasoline. What is the probability that the next two customers purchasing gasoline will each use a credit card?

(b) Let

A = the event that the first customer uses a credit card

B = the event that the second customer uses a credit card

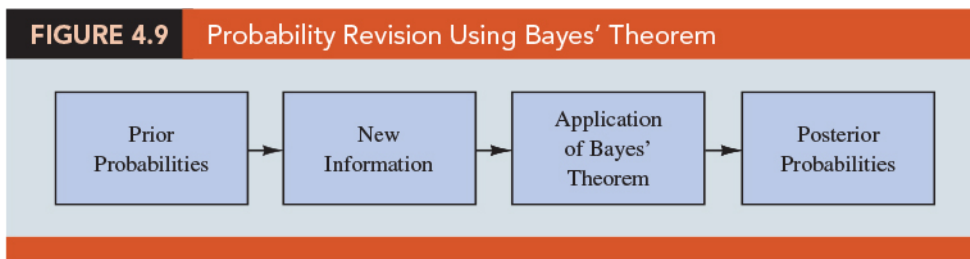
then the event of interest is $A \cap B$. Assume that A and B are independent events. Thus,

$$\underline{P(A \cap B) = P(A)P(B) = (0.80)(0.80) = 0.64}$$

4.5 Bayes' Theorem

1. Motivation:

- (a) Often we begin probability analysis with initial or prior probabilities.
- (b) Then, from a sample, special report, or a product test, we obtain some additional information.
- (c) Given this information, we calculate revised or posterior probabilities.
- (d) Bayes' theorem provides the means for revising the prior probabilities.



2. The Bayes' theorem for the case of two events:

$$P(A_1|B) = \frac{P(A_1)P(B|A_1)}{P(A_1)P(B|A_1) + P(A_2)P(B|A_2)}$$

$$P(A_2|B) = \frac{P(A_2)P(B|A_2)}{P(A_2)P(B|A_2) + P(A_1)P(B|A_1)}$$

3. **Example** Manufacturing firm receives shipments from two suppliers
- (a) Consider a manufacturing firm that receives shipments of parts from two different suppliers. Let A_1 (A_2) denote the event that a part is from supplier 1 (2).
- (b) Currently, 65% (35%) of the parts purchased by the company are from supplier 1 (2). Hence, if a part is selected at random, we would assign the prior probabilities $P(A_1) = 0.65$ and $P(A_2) = 0.35$.
- (c) The quality of the purchased parts varies with the source of supply. Historical data suggest that the quality ratings of the two suppliers are as shown in Table 4.6.

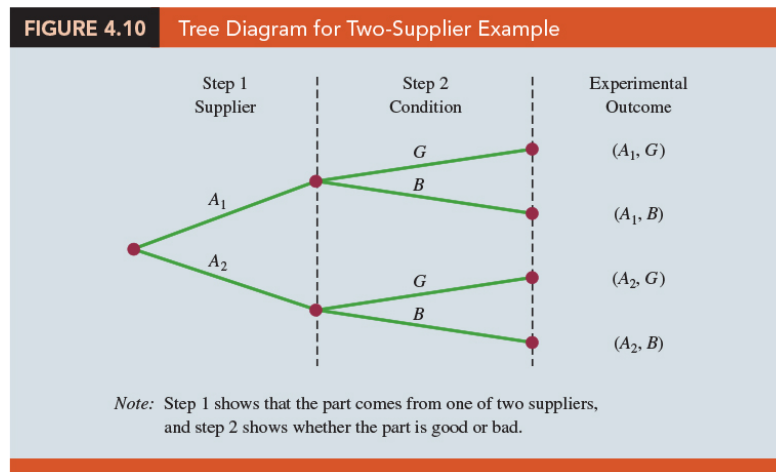
	Percentage Good Parts	Percentage Bad Parts
Supplier 1	98	2
Supplier 2	95	5

- (d) Let G (B) denote the event that a part is good (bad), the information in Table 4.6 provides the following conditional probability values.

$$P(G|A_1) = 0.98, \quad P(B|A_1) = 0.02$$

$$P(G|A_2) = 0.95, \quad P(B|A_2) = 0.05$$

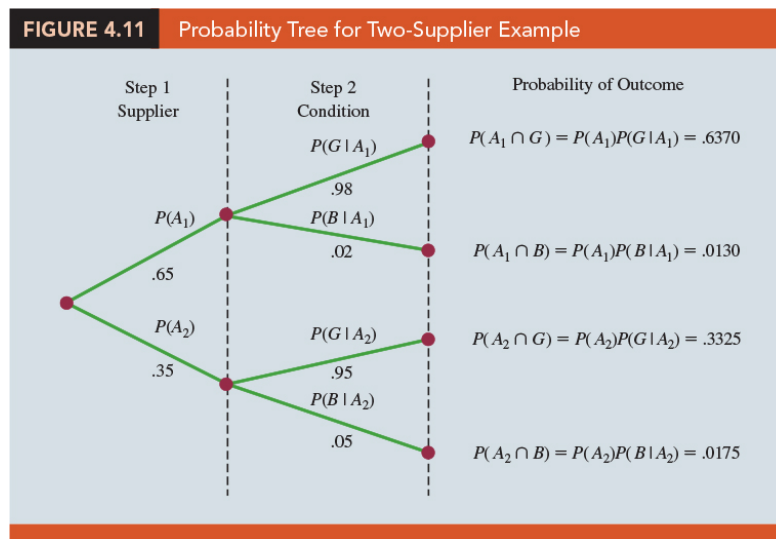
- (e) (Figure 4.10) The tree diagram depicts the process of the firm receiving a part from one of the two suppliers and then discovering that the part is good or bad as a two-step experiment. We see that four experimental outcomes are possible; two correspond to the part being good and two correspond to the part being bad.



- (f) Each of the experimental outcomes is the intersection of two events, so we can use the multiplication rule to compute the probabilities. For instance,

$$P(A_1, G) = P(A_1 \cap G) = P(A_1)P(G|A_1)$$

- (g) **Probability tree** (Figure 4.11): From left to right through the tree, the probabilities for each branch at step 1 are prior probabilities and the probabilities for each branch at step 2 are conditional probabilities.



- (h) To find the probabilities of each experimental outcome, we simply multiply the probabilities on the branches leading to the outcome. Each of these joint probabilities is shown in Figure 4.11 along with the known probabilities for each branch.

 Question (p209)

Suppose now that the parts from the two suppliers are used in the firm's manufacturing process and that a machine breaks down because it attempts to process a bad part. Given the information that the part is bad, what is the probability that it came from supplier 1 and what is the probability that it came from supplier 2?

sol:

- Letting B denote the event that the part is bad, we are looking for the posterior probabilities $P(A_1|B)$ and $P(A_2|B)$. We know that

$$P(A_1|B) = \frac{P(A_1 \cap B)}{P(B)} \Rightarrow \underline{P(A_1 \cap B) = P(A_1)P(B|A_1)}$$

- Find $P(B)$:

$$P(B) = \underline{P(A_1 \cap B) + P(A_2 \cap B) = P(A_1)P(B|A_1) + P(A_2)P(B|A_2)}$$

- We have

$$\begin{aligned} P(A_1|B) &= \frac{P(A_1)P(B|A_1)}{P(A_1)P(B|A_1) + P(A_2)P(B|A_2)} \\ &= \frac{(0.65)(0.02)}{(0.65)(0.02) + (0.35)(0.05)} = \frac{0.0130}{0.0130 + 0.0175} \\ &= \frac{0.0130}{0.0305} = .4262 \\ P(A_2|B) &= \frac{(0.35)(0.05)}{(0.65)(0.02) + (0.35)(0.05)} \\ &= \frac{0.0175}{0.0130 + 0.0175} = \frac{0.0175}{0.0305} = 0.5738 \end{aligned}$$

- Note that in this application we started with a probability of 0.65 that a part selected at random was from supplier 1. However, given information that the part is bad, the probability that the part is from supplier 1 drops to 0.4262. In fact, if the part is bad, it has better than a 50–50 chance that it came from supplier 2; that is, $P(A_2|B) = 0.5738$.

3. Bayes' theorem is applicable when the events for which we want to compute posterior probabilities are mutually exclusive and their union is the entire sample space. For the case of n mutually exclusive events A_1, A_2, \dots, A_n , whose union is the entire sample space, Bayes' theorem can be used to compute any posterior probability $P(A_i|B)$:

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{P(A_1)P(B|A_1) + P(A_2)P(B|A_2) + \dots + P(A_n)P(B|A_n)}$$

Tabular Approach

1. A tabular approach is helpful in conducting the Bayes' theorem calculations.

(1) Events A_i	(2) Prior Probabilities $P(A_i)$	(3) Conditional Probabilities $P(B A_i)$	(4) Joint Probabilities $P(A_i \cap B)$	(5) Posterior Probabilities $P(A_i B)$
A_1	.65	.02	.0130	.0130/.0305 = .4262
A_2	.35	.05	.0175	.0175/.0305 = .5738
	<u>1.00</u>		$P(B) = .0305$	<u>1.0000</u>

Step 1. Prepare the following three columns:

Column 1 : The mutually exclusive events A_i for which posterior probabilities are desired.

Column 2 : The prior probabilities $P(A_i)$ for the events.

Column 3 : The conditional probabilities $P(B|A_i)$ of the new information B given each event.

Step 2. In column 4, compute the joint probabilities $P(A_i \cap B)$ for each event and the new information B by using the multiplication law: $P(A_i \cap B) = P(A_i)P(B|A_i)$.

Step 3. Sum the joint probabilities in column 4. The sum is the probability of the new information, $P(B)$.

Example Thus we see in Table 4.7 that there is a 0.0130 probability that the part came from supplier 1 and is bad and a 0.0175 probability that the part came from supplier 2 and is bad. Because these are the only two ways in which a bad part can be obtained, the sum $0.0130 + 0.0175$ shows an overall probability of 0.0305 of finding a bad part from the combined shipments of the two suppliers.

Step 4. In column 5, compute the posterior probabilities using the basic relationship of conditional probability.

$$P(A_i|B) = \frac{P(A_i \cap B)}{P(B)}$$

Note that the joint probabilities $P(A_i \cap B)$ are in column 4 and the probability $P(B)$ is the sum of column 4.

☺ EXERCISES

4.1 : 4, 6, 11, 13

4.2 : 16, 17, 19

4.3 : 23, 26, 29

4.4 : 31, 32, 35, 36

4.5 : 40, 42, 45

SUP : 49, 51, 56, 59

“你必須找到你的熱情所在”

“You’ve got to find what you love”

— *Steve Jobs (February 24, 1955 –October 5, 2011)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 5: Discrete Probability Distributions

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

5.1 Random Variables

1. A random variable is a numerical description of the outcome of an experiment.
2. A random variable can be classified as being either discrete or continuous depending on the numerical values it assumes.
3. A random variable, often denoted by X (or some other capital letter), is a function that maps outcomes to numbers on the real line.

Discrete Random Variables

1. A discrete random variable assumes either a finite number of values or an infinite sequence of values such as $0, 1, 2, \dots$.

Random Experiment	Random Variable (x)	Possible Values for the Random Variable
Flip a coin	Face of coin showing	1 if heads; 0 if tails
Roll a die	Number of dots showing on top of die	1, 2, 3, 4, 5, 6
Contact five customers	Number of customers who place an order	0, 1, 2, 3, 4, 5
Operate a health care clinic for one day	Number of patients who arrive	0, 1, 2, 3, ...
Offer a customer the choice of two products	Product chosen by customer	0 if none; 1 if choose product A; 2 if choose product B

Continuous Random Variables

1. A continuous random variable assumes any numerical value in an interval or collection of intervals.
2. Experimental outcomes based on measurement scales such as time, weight, distance, and temperature can be described by continuous random variables.

TABLE 5.2 Examples of Continuous Random Variables

Random Experiment	Random Variable (x)	Possible Values for the Random Variable
Customer visits a web page	Time customer spends on web page in minutes	$x \geq 0$
Fill a soft drink can (max capacity = 360 milliliters)	Number of milliliters	$0 \leq x \leq 360$
Test a new chemical process	Temperature when the desired reaction takes place (min temperature = 65°C; max temperature = 100°C)	$65 \leq x \leq 100$
Invest \$10,000 in the stock market	Value of investment after one year	$x \geq 0$

5.2 Developing Discrete Probability Distributions

1. The **probability distribution** for a random variable describes how probabilities are distributed over the values of the random variable.
2. For a discrete random variable X , a probability function, denoted by $f(x)$ (or $f_X(x)$), provides the probability for each value of the random variable.
3. The classical, subjective, and relative frequency methods of assigning probabilities can be used to develop discrete probability distributions.

4. *The classical method:*

- (a) The classical method of assigning probabilities to values of a random variable is applicable when the experimental outcomes generate values of the random variable that are equally likely.
- (b) **Example** The experiment of rolling a die.
- i. Each of the outcomes (the numbers 1, 2, 3, 4, 5, or 6) is equally likely.
 - ii. Let $X = \text{number obtained on one roll of a die}$.
 - iii. (Table 5.3) the probability distribution function $f(x)$ of X .


Number Obtained x	Probability of x $f(x)$
1	1/6
2	1/6
3	1/6
4	1/6
5	1/6
6	1/6

5. *The subjective method:*

- (a) Each probability is assigned by user's best judgment. Different people can be expected to obtain different probability distributions.

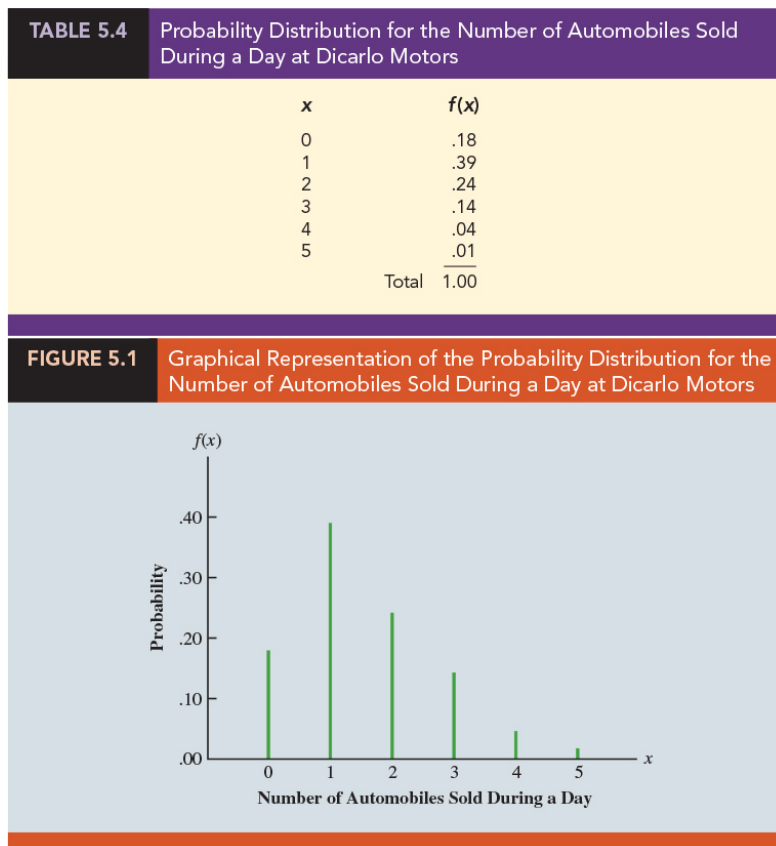
6. *The relative frequency method:*

- (a) We treat the (large amounts of) data as if they were the population and use the relative frequency method to assign probabilities to the experimental outcomes.
- (b) The use of the relative frequency method to develop discrete probability distributions leads to what is called an empirical discrete distribution. (more widely used in practice)

 Question (p229)

Over the past 300 days, DiCarlo DiCarlo Motors in Saratoga, New York, has experienced 54 (117, 72, 42, 12, 3) days with no (1, 2, 3, 4, 5) automobiles sold. Suppose we consider the experiment of observing a day of operations at DiCarlo Motors and define the random variable of interest as X = the number of automobiles sold during a day. (a) Use the relative frequency method to develop a probability distribution for the number of cars sold per day. (b) Graph the probability distribution such that the values of the random variable X for DiCarlo Motors are shown on the horizontal axis and the probability associated with these values is shown on the vertical axis.

sol:



7. A primary advantage of defining a random variable and its probability distribution is that once the probability distribution is known, it is relatively easy to determine the probability of a variety of events that may be of interest to a decision maker.
8. *Required Conditions for a Discrete Probability Function:*

$$\underline{f(x) \geq 0} \quad \text{and} \quad \underline{\sum f(x) = 1}$$

9. **Discrete Uniform Probability Distribution:** $f(x) = 1/n$ where n is the number of values the random variable may assume.

 **Question** (p230)

For the experiment of rolling a fair die, we define the random variable X to be the number of dots on the upward face. Find the probability function for this discrete random variable.

sol:

For this experiment, $n = 6$ values are possible for the random variable; $X = 1, 2, \dots, 6$. Since the probabilities are equally likely, the discrete probability function for this discrete random variable is

$$f_X(x) = 1/6, \quad x = 1, 2, 3, 4, 5, 6.$$

X is called the discrete uniform random variable.

5.3 Expected Value and Variance

Expected Value

1. The expected value, or mean, of a random variable X is a measure of the central location for the discrete random variable:

$$E(X) = \underline{\mu} = \underline{\sum xf(x)} .$$

2. The expected value is a weighted average of the values of the random variable where the weights are the probabilities .

Variance

1. Use variance to summarize the variability in the values of a discrete random variable:

$$Var(X) = \underline{\sigma^2} = \underline{\sum (x - \mu)^2 f(x)}$$

2. The variance is a weighted average of the squared deviations of a random variable from its mean. The weights are the probabilities.

 Question (p234)

Let the random variable X be the number of auto mobiles sold during a day in the DiCarlo Motors automobile sales example from Section 5.2, find the expected value and the variance of this discrete random variable.

sol:

TABLE 5.5 Calculation of the Expected Value for the Number of Automobiles Sold During a Day at Dicarlo Motors		
x	$f(x)$	$xf(x)$
0	.18	0(.18) = .00
1	.39	1(.39) = .39
2	.24	2(.24) = .48
3	.14	3(.14) = .42
4	.04	4(.04) = .16
5	.01	5(.01) = .05
		1.50

$E(x) = \mu = \sum xf(x)$

x	$x - \mu$	$(x - \mu)^2$	$f(x)$	$(x - \mu)^2 f(x)$
0	$0 - 1.50 = -1.50$	2.25	.18	$2.25(.18) = .4050$
1	$1 - 1.50 = -.50$.25	.39	$.25(.39) = .0975$
2	$2 - 1.50 = .50$.25	.24	$.25(.24) = .0600$
3	$3 - 1.50 = 1.50$	2.25	.14	$2.25(.14) = .3150$
4	$4 - 1.50 = 2.50$	6.25	.04	$6.25(.04) = .2500$
5	$5 - 1.50 = 3.50$	12.25	.01	$12.25(.01) = .1225$
				1.2500

$\sigma^2 = \sum(x - \mu)^2 f(x)$

5.4 Bivariate Distributions, Covariance, and Financial Portfolios

1. A probability distribution involving two random variables is called a bivariate probability distribution.
2. **Bivariate experiment:** Each outcome for a bivariate experiment consists of two values, one for each random variable.
3. **Example** Consider the bivariate experiment of rolling a pair of dice. The outcome consists of two values, the number obtained with the first die and the number obtained with the second die.
4. **Example** Consider the experiment of observing the financial markets for a year and recording the percentage gain for a stock fund and a bond fund. Again, the experimental outcome provides a value for two random variables, the percent gain in the stock fund and the percent gain in the bond fund.

A Bivariate Empirical Discrete Probability Distribution

1. **Example** (Section 5.2) DiCarlo Motors automobile dealership in Saratoga, New York.

- (a) (Table 5.7) DiCarlo has another dealership in Geneva, New York. Table 5.7 shows the number of cars sold at each of the dealerships over a 300-day period.
- (b) The numbers in the bottom (total) row are the frequencies we used to develop an empirical probability distribution for daily sales at DiCarlo's Saratoga dealership. The numbers in the rightmost (total) column are the frequencies of daily sales for the Geneva dealership.
- (c) Entries in the body of the table give the number of days the Geneva dealership had a level of sales indicated by the row, when the Saratoga dealership had the level of sales indicated by the column.

TABLE 5.7 Number of Automobiles Sold at DiCarlo's Saratoga and Geneva Dealerships Over 300 Days

Geneva Dealership	Saratoga Dealership						Total
	0	1	2	3	4	5	
0	21	30	24	9	2	0	86
1	21	36	33	18	2	1	111
2	9	42	9	12	3	2	77
3	3	9	6	3	5	0	26
Total	54	117	72	42	12	3	300

- (d) Consider the bivariate experiment of observing a day of operations at DiCarlo Motors and recording the number of cars sold. Define

X = number of cars sold at the Geneva dealership


Y = the number of cars sold at the Saratoga dealership

- (e) (Table 5.8) divide all of the frequencies in Table 5.7 by the number of observations (300) to develop a bivariate empirical discrete probability distribution for automobile sales at the two DiCarlo dealerships.

TABLE 5.8 Bivariate Empirical Discrete Probability Distribution for Daily Sales at DiCarlo Dealerships in Saratoga and Geneva, New York

Geneva Dealership	Saratoga Dealership						Total
	0	1	2	3	4	5	
0	.0700	.1000	.0800	.0300	.0067	.0000	.2867
1	.0700	.1200	.1100	.0600	.0067	.0033	.3700
2	.0300	.1400	.0300	.0400	.0100	.0067	.2567
3	.0100	.0300	.0200	.0100	.0167	.0000	.0867
Total	.18	.39	.24	.14	.04	.01	1.0000

- (f) The probabilities in the lower (right) margin provide the marginal distribution for the DiCarlo Motors Saratoga (Geneva) dealership.
 - (g) The probabilities in the body of the table provide the bivariate probability distribution for sales at both dealerships.
2. Bivariate probabilities are often called joint probabilities.
 3. Note that there is one bivariate probability for each experimental outcome.
 4. With C_1 possible values for X and C_2 possible values for Y , there are $C_1 \times C_2$ experimental outcomes and bivariate probabilities.

 **Question** (p240)

Refer to the DiCarlo Motors automobile dealership data, find the probability distribution for total sales at both DiCarlo dealerships and the expected value and variance of total sales.

sol:

TABLE 5.9		Calculation of the Expected Value and Variance for Total Daily Sales at DiCarlo Motors			
s	$f(s)$	$sf(s)$	$s - E(s)$	$(s - E(s))^2$	$(s - E(s))^2 f(s)$
0	.0700	.0000	-2.6433	6.9872	.4891
1	.1700	.1700	-1.6433	2.7005	.4591
2	.2300	.4600	-.6433	.4139	.0952
3	.2900	.8700	.3567	.1272	.0369
4	.1267	.5067	1.3567	1.8405	.2331
5	.0667	.3333	2.3567	5.5539	.3703
6	.0233	.1400	3.3567	11.2672	.2629
7	.0233	.1633	4.3567	18.9805	.4429
8	.0000	.0000	5.3567	28.6939	.0000
		$E(s) = 2.6433$			$Var(s) = 2.3895$

5. The covariance between two random variables X and Y is:

$$\sigma_{xy} = Cov(X, Y) = \frac{[Var(X + Y) - Var(X) - Var(Y)]}{2} \quad \text{or}$$

$$\sigma_{xy} = Cov(X, Y) = \frac{\sum (x - E(X))(y - E(Y)) f(x, y)}{n}$$

6. The correlation between two random variables X and Y is the covariance divided by the product of the standard deviations for the two random variables:

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} .$$


7. *Expected Value of a Linear Combination of Random Variables X and Y :*

$$E(aX + bY) = \underline{aE(X) + bE(Y)} .$$

8. *Variance of a Linear Combination of Two Random Variables:*

$$\text{Var}(aX + bY) = \underline{a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)}$$

where $\text{Cov}(X, Y) = \sigma_{xy}$ is the covariance of X and Y .

 **Question** (p241)


Refer to the DiCarlo Motors automobile dealership data, compute the covariance and the correlation coefficient between daily sales at the two DiCarlo dealerships.

sol:

TABLE 5.10		Calculation of the Expected Value and Variance of Daily Automobile Sales at DiCarlo Motors' Geneva Dealership			
x	$f(x)$	$xf(x)$	$x - E(x)$	$[(x - E(x))^2]$	$[x - E(x)]^2 f(x)$
0	.2867	.0000	-1.1435	1.3076	.3749
1	.3700	.3700	-.1435	.0206	.0076
2	.2567	.5134	.8565	.7336	.1883
3	.0867	.2601	1.8565	3.447	.2988
		$E(x) = 1.1435$			$\text{Var}(x) = .8696$

Financial Applications

- To see how what we have learned can be useful in constructing financial portfolios that provide a good balance of risk and return.
- The standard deviation of percent return is often used as a measure of risk associated with an investment.

 **Question** (p242)

(Table 5.11) A financial advisor is considering four possible economic scenarios for the coming year and has developed a probability distribution showing the percent return, X , for investing in a largecapstock fund and the percent return, Y , for investing in a long-term government bond fund given each of the scenarios. The bivariate probability distribution for X and Y is shown in Table 5.11.

Economic Scenario	Probability $f(x, y)$	Large-Cap Stock Fund (x)	Long-Term Government Bond Fund (y)
Recession	.10	-40	30
Weak Growth	.25	5	5
Stable Growth	.50	15	4
Strong Growth	.15	30	2

Compute the expected percent return for investing in the stock fund, $E(X)$, and the expected percent return for investing in the bond fund, $E(Y)$. Draw the conclusion.

sol:

$$E(X) = 0.10(-40) + 0.25(5) + 0.5(15) + 0.15(30) = 9.25$$

$$E(Y) = 0.10(30) + 0.25(5) + 0.5(4) + 0.15(2) = 6.55$$

Using this information, we might conclude that investing in the stock fund is a better investment. It has a higher expected return, 9.25%.

 Question (p242)

Refer to the financial advisor example. Compute the standard deviation of the percent returns for the stock and bond fund investments.

sol:

$$\begin{aligned} \text{Var}(X) &= 0.1(-40 - 9.25)^2 + 0.25(5 - 9.25)^2 + 0.50(15 - 9.25)^2 + 0.15(30 - 9.25)^2 \\ &= 328.1875 \end{aligned}$$

$$\begin{aligned} \text{Var}(Y) &= 0.1(30 - 6.55)^2 + 0.25(5 - 6.55)^2 + 0.50(4 - 6.55)^2 + 0.15(2 - 6.55)^2 \\ &= 61.9475 \end{aligned}$$

The standard deviation of the return from an investment in the stock fund is $\sigma_x = \sqrt{328.1875} = 18.1159\%$ and the standard deviation of the return from an investment in the bond fund is $\sigma_y = \sqrt{61.9475} = 7.8707\%$. So, we can conclude that investing in the bond fund is less risky. It has the smaller standard deviation.

3. We have already seen that the stock fund offers a greater expected return, so if we want to choose between investing in either the stock fund or the bond fund it depends on our attitude toward risk and return.
4. An aggressive investor might choose the stock fund because of the higher expected return; a conservative investor might choose the bond fund because of the lower risk. But, there are other options.
5. What about the possibility of investing in a portfolio consisting of both an investment in the stock fund and an investment in the bond fund?

 Question (p242)

Refer to the financial advisor example. Suppose we would like to consider a portfolio by investing equal amounts in the largecap stock fund and in the longterm government bond fund. Evaluate this portfolio by computing its expected return and risk.

sol:

1. Define X as the percent return from an investment in the stock fund and Y as the percent return from an investment in the bond fund so the percent return for our portfolio is $R = 0.5X + 0.5Y$.

$$E(0.5X + 0.5Y) = 0.5E(X) + 0.5E(Y) = 0.5(9.25) + 0.5(6.55) = 7.9$$

The expected return for investing in the portfolio is 7.9%.

2. We have $Var(X) = 328.1875$ and $Var(Y) = 61.9475$. Also, it can be shown that

$$Var(X + Y) = 119.46.$$

The covariance of the random variables X and Y is

$$\begin{aligned}\sigma_{xy} &= [Var(X + Y) - Var(X) - Var(Y)]/2 \\ &= [119.46 - 328.1875 - 61.9475]/2 \\ &= -135.3375.\end{aligned}$$

A negative covariance between X and Y means that when X tends to be above its mean, Y tends to be below its mean and vice versa.

3. Compute the variance of return for our portfolio

$$\begin{aligned}Var(0.5X + 0.5Y) &= 0.5^2(328.1875) + 0.5^2(61.9475) + 2(0.5)(0.5)(-135.3375) \\ &= 29.865\end{aligned}$$

The standard deviation of our portfolio is then given by $\sigma_{0.5x+0.5y} = \sqrt{29.865} = 5.4650\%$. This is our measure of risk for the portfolio consisting of investing 50% in the stock fund and 50% in the bond fund.

 Question (p242)

Refer to the financial advisor example. Compare the three investment alternatives according to their expected returns, variances, and standard deviations: investing solely in the stock fund, investing solely in the bond fund, or creating a portfolio by dividing our investment amount equally between the stock and bond funds. Which of these alternatives would you prefer?

Investment Alternative	Expected Return (%)	Variance of Return	Standard Deviation of Return (%)
100% in Stock Fund	9.25	328.1875	18.1159
100% in Bond Fund	6.55	61.9475	7.8707
Portfolio (50% in stock fund, 50% in bond fund)	7.90	29.865	5.4650

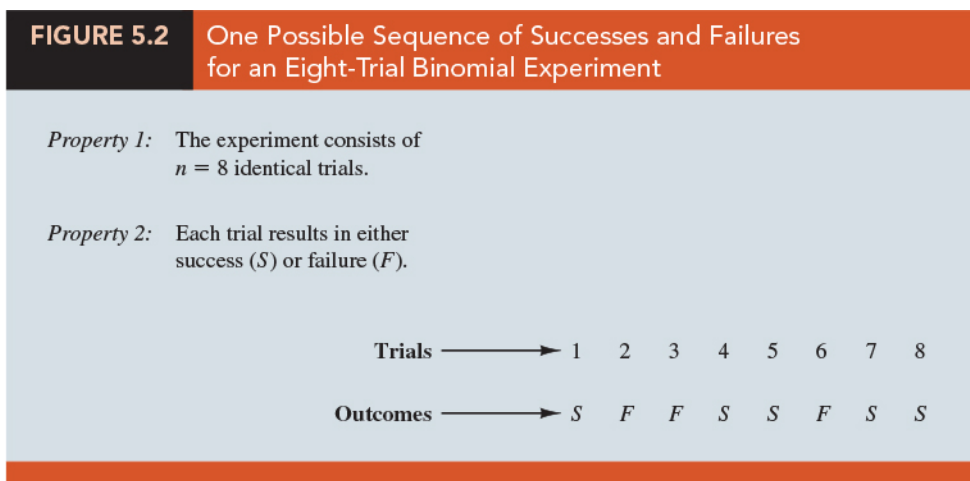
sol:

1. The expected return is highest for investing 100% in the stock fund, but the risk is also highest. The standard deviation is 18.1159%.
2. Investing 100% in the bond fund has a lower expected return, but a significantly smaller risk.
3. Investing 50% in the stock fund and 50% in the bond fund (the portfolio) has an expected return that is halfway between that of the stock fund alone and the bond fund alone. But note that it has less risk than investing 100% in either of the individual funds. It has both a higher return and less risk (smaller standard deviation) than investing solely in the bond fund.
4. Whether you would choose to invest in the stock fund or the portfolio depends on your attitude toward risk.
5. The stock fund has a higher expected return. But the portfolio has significantly less risk and also provides a fairly good return. Many would choose it.
6. It is the negative covariance between the stock and bond funds that has caused the portfolio risk to be so much smaller than the risk of investing solely in either of the individual funds.

5.5 Binomial Probability Distribution


A Binomial Experiment

- A binomial experiment exhibits the following four properties.
 - The experiment consists of n identical trials.
 - Two outcomes are possible on each trial. We refer to one outcome as a success and the other outcome as a failure.
 - The probability of a success, denoted by p , does not change from trial to trial. Consequently, the probability of a failure, denoted by $1 - p$, does not change from trial to trial. (stationarity assumption)
 - The trials are independent.
- If properties (b), (c), and (d) are present, we say the trials are generated by a Bernoulli process.
- (Figure 5.2) depicts one possible sequence of successes and failures for a binomial experiment involving eight trials.



- In a binomial experiment, our interest is in the number of successes occurring in the n trials. If we let X denote the number of successes occurring in the n trials, we see that X can assume the values of $0, 1, 2, 3, \dots, n$.

5. Because the number of values is finite, X is a discrete random variable. The probability distribution associated with this random variable is called the Binomial probability distribution. Denote by $X \sim B(n, p)$.

 Question (p248)


Consider the experiment of tossing a coin five times and on each toss observing whether the coin lands with a head or a tail on its upward face. Suppose we want to count the number of heads appearing over the five tosses. Does this experiment show the properties of a binomial experiment? What is the random variable of interest?

sol:

Four properties of a binomial experiment are satisfied:

1. The experiment consists of five identical trials; each trial involves the tossing of one coin.
2. Two outcomes are possible for each trial: a head or a tail. We can designate head a success and tail a failure.
3. The probability of a head and the probability of a tail are the same for each trial, with $p = 0.5$ and $1 - p = 0.5$.
4. The trials or tosses are independent because the outcome on any one trial is not affected by what happens on other trials or tosses.

The random variable of interest is $X =$ the number of heads appearing in the five trials. In this case, X can assume the values of 0, 1, 2, 3, 4, or 5.

 Question (p249)

Consider an insurance salesperson who visits 10 randomly selected families. The outcome associated with each visit is classified as a success if the family purchases an insurance policy and a failure if the family does not. From past experience, the salesperson knows the probability that a randomly selected family will purchase an insurance policy is 0.10. Checking the properties of a binomial experiment.

sol:

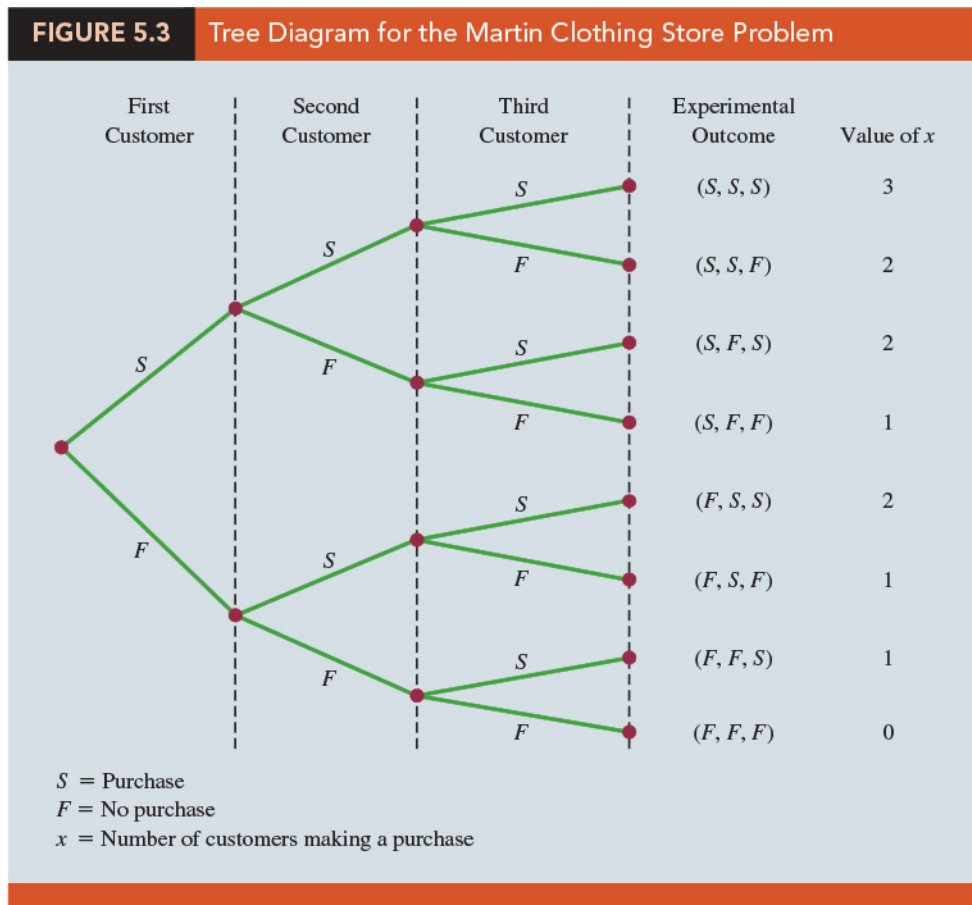
Four assumptions are satisfied, this example is a binomial experiment:

1. The experiment consists of 10 identical trials; each trial involves contacting one family.
2. Two outcomes are possible on each trial: the family purchases a policy (success) or the family does not purchase a policy (failure).
3. The probabilities of a purchase and a nonpurchase are assumed to be the same for each sales call, with $p = 0.10$ and $1-p = 0.90$.
4. The trials are independent because the families are randomly selected.

The random variable of interest is the number of sales obtained in contacting the 10 families. In this case, X can assume the values of $0, 1, \dots, 10$.

Martin Clothing Store Problem

1. Let us consider the purchase decisions of the next three customers who enter the Martin Clothing Store. On the basis of past experience, the store manager estimates the probability that any one customer will make a purchase is 0.30. What is the probability that two of the next three customers will make a purchase?
2. (Figure 5.3) tree diagram shows that the experiment of observing the three customers each making a purchase decision has eight possible outcomes.



3. Using S to denote success (a purchase) and F to denote failure (no purchase), we are interested in experimental outcomes involving two successes in the three trials (purchase decisions).
4. Verify that the experiment involving the sequence of three purchase decisions can be viewed as a binomial experiment:
 - (a) The experiment can be described as a sequence of three identical trials, one trial for each of the three customers who will enter the store.
 - (b) Two outcomes - the customer makes a purchase (success) or the customer does not make a purchase (failure) - are possible for each trial.
 - (c) The probability that the customer will make a purchase (0.30) or will not make a purchase (0.70) is assumed to be the same for all customers.
 - (d) The purchase decision of each customer is independent of the decisions of the other customers.

5. The number of experimental outcomes resulting in exactly x successes in n trials can be computed using:

$$\binom{n}{x} = \frac{n!}{x!(n-x)!},$$

where $n! = n(n-1)(n-2)\cdots(2)(1)$ and $0! = 1$.

6. Because the trials of a binomial experiment are independent, we can simply multiply the probabilities associated with each trial outcome to find the probability of a particular sequence of successes and failures.
7. The probability of purchases by the first two customers and no purchase by the third customer, denoted (S, S, F) , is given by $pp(1-p)$. With a 0.30 probability of a purchase on any one trial, the probability of a purchase on the first two trials and no purchase on the third is given by $(0.30)(0.30)(0.70) = 0.063$.
8. In any binomial experiment, all sequences of trial outcomes yielding x successes in n trials have the same probability of occurrence. The probability of each sequence of trials yielding x successes in n trials is $p^x(1-p)^{n-x}$.

9. The binomial probability function

$$f(x) = \frac{\binom{n}{x} p^x (1-p)^{n-x}}{\quad},$$

where

x = the number of successes, $x = 0, 1, 2, \dots, n$

p = the probability of a success on one trial

n = the number of trials

$f(x)$ = the probability of x successes in n trials

10. For the binomial probability distribution, X is a discrete random variable with the probability function $f(x)$ (or $f_X(x)$) applicable for values of $x = 0, 1, 2, \dots, n$.
($X \sim B(n, p)$)

 Question (p113)

Refer to the Martin Clothing Store example, compute the probability that no customer makes a purchase, exactly one customer makes a purchase, exactly two customers make a purchase, and all three customers make a purchase.

sol: The calculations are summarized in Table 5.13, which gives the probability distribution of the number of customers making a purchase. Figure 5.4 is a graph of this probability distribution.

TABLE 5.13 Probability Distribution for the Number of Customers Making a Purchase	
x	$f(x)$
0	$\frac{3!}{0!3!} (.30)^0(.70)^3 = .343$
1	$\frac{3!}{1!2!} (.30)^1(.70)^2 = .441$
2	$\frac{3!}{2!1!} (.30)^2(.70)^1 = .189$
3	$\frac{3!}{3!0!} (.30)^3(.70)^0 = \frac{.027}{1.000}$

FIGURE 5.4 Graphical Representation of the Probability Distribution for the Number of Customers Making a Purchase



Using Tables of Binomial Probabilities

- (Table 5.14) Table 5 of Appendix B provides a table of binomial probabilities. To use this table, we must specify the values of n , p , and x for the binomial experiment of interest.

TABLE 5.14 Selected Values from the Binomial Probability Table
Example: $n = 10$, $x = 3$, $P = .40$; $f(3) = .2150$

n	x	.05	.10	.15	.20	p .25	.30	.35	.40	.45	.50
9	0	.6302	.3874	.2316	.1342	.0751	.0404	.0207	.0101	.0046	.0020
	1	.2985	.3874	.3679	.3020	.2253	.1556	.1004	.0605	.0339	.0176
	2	.0629	.1722	.2597	.3020	.3003	.2668	.2162	.1612	.1110	.0703
	3	.0077	.0446	.1069	.1762	.2336	.2668	.2716	.2508	.2119	.1641
	4	.0006	.0074	.0283	.0661	.1168	.1715	.2194	.2508	.2600	.2461
	5	.0000	.0008	.0050	.0165	.0389	.0735	.1181	.1672	.2128	.2461
	6	.0000	.0001	.0006	.0028	.0087	.0210	.0424	.0743	.1160	.1641
	7	.0000	.0000	.0000	.0003	.0012	.0039	.0098	.0212	.0407	.0703
	8	.0000	.0000	.0000	.0000	.0001	.0004	.0013	.0035	.0083	.0176
	9	.0000	.0000	.0000	.0000	.0000	.0000	.0001	.0003	.0008	.0020
10	0	.5987	.3487	.1969	.1074	.0563	.0282	.0135	.0060	.0025	.0010
	1	.3151	.3874	.3474	.2684	.1877	.1211	.0725	.0403	.0207	.0098
	2	.0746	.1937	.2759	.3020	.2816	.2335	.1757	.1209	.0763	.0439
	3	.0105	.0574	.1298	.2013	.2503	.2668	.2522	.2150	.1665	.1172
	4	.0010	.0112	.0401	.0881	.1460	.2001	.2377	.2508	.2384	.2051
	5	.0001	.0015	.0085	.0264	.0584	.1029	.1536	.2007	.2340	.2461
	6	.0000	.0001	.0012	.0055	.0162	.0368	.0689	.1115	.1596	.2051
	7	.0000	.0000	.0001	.0008	.0031	.0090	.0212	.0425	.0746	.1172
	8	.0000	.0000	.0000	.0001	.0004	.0014	.0043	.0106	.0229	.0439
	9	.0000	.0000	.0000	.0000	.0000	.0001	.0005	.0016	.0042	.0098
	10	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0001	.0003	.0010

Expected Value and Variance for the Binomial Distribution

- Expected Value and Variance for the Binomial Distribution:

$$E(X) = \mu = \underline{np}$$

$$Var(X) = \sigma^2 = \underline{np(1-p)}$$

 Question (p255)

Refer to the Martin Clothing Store problem with three customers, (a) compute the expected number, the variance and standard deviation of customers who will make a purchase. (b) Suppose that for the next month the Martin Clothing Store forecasts 1000 customers will enter the store. (c) What is the expected number, the variance and standard deviation of customers who will make a purchase?

sol:

5.6 Poisson Probability Distribution

1. The probability distribution of Poisson random variable is called a Poisson distribution. It is a discrete random variable that is often useful in estimating the number of occurrences over a specified interval of time or space.
2. Example The random variable of interest might be the number of arrivals at a car wash in one hour, the number of repairs needed in 10 miles of highway, or the number of leaks in 100 miles of pipeline.
3. Two properties must be satisfied, the number of occurrences is a random variable described by the Poisson probability distribution:
 - (a) The probability of an occurrence is the same for any two intervals of equal length.

- (b) The occurrence or nonoccurrence in any interval is independent of the occurrence or nonoccurrence in any other interval.
4. The **Poisson probability mass function** is defined by

$$f(x) = \frac{\mu^x e^{-\mu}}{x!},$$

where

- x = the number of occurrences in an interval, $x = 0, 1, 2, \dots$
 $f(x)$ = the probability of X occurrences in an interval
 μ = expected value or mean number of occurrences
 e = 2.71828.

5. For the Poisson probability distribution, X is a discrete random variable indicating the number of occurrences in the interval. (Denote $X \sim P(\mu)$ or $X \sim Poi(\mu)$)
6. Since there is no stated upper limit for the number of occurrences, the probability function $f(x)$ is applicable for values $x = 0, 1, 2, \dots$ without limit.
7. In practical applications, X will eventually become large enough so that $f(x)$ is approximately zero and the probability of any larger values of X becomes negligible.

An Example Involving Time Intervals

- Suppose that we are interested in the number of patients who arrive at the emergency room of a large hospital during a 15-minute period on weekday mornings. If we can assume that the probability of a patient arriving is the same for any two time periods of equal length and that the arrival or nonarrival of a patient in any time period is independent of the arrival or nonarrival in any other time period, the Poisson probability function is applicable.
- Suppose these assumptions are satisfied and an analysis of historical data shows that the average number of patients arriving in a 15-minute period of time is 10; in this case, the following probability function applies:

$$f(x) = \frac{10^x e^{-10}}{x!}$$

The random variable here is $X =$ number of patients arriving in any 15-minute period.

3. If management wanted to know the probability of exactly five arrivals in 15 minutes, we would set $X = 5$ and thus obtain

$$\text{Probability of exactly 5 arrivals in 15-minutes} = f(5) = \frac{10^5 e^{-10}}{5!} = 0.0378$$

4. (Table 5.15) Table 7 of Appendix B provides probabilities for specific values of x and μ .
5. In the preceding example, the mean of the Poisson distribution is $\mu = 10$ arrivals per 15-minute period.
6. A property of the Poisson distribution is that the mean of the distribution and the variance of the distribution are equal. Thus, the variance for the number of arrivals during 15-minute periods is $\sigma^2 = 10$. The standard deviation is $\sigma = \sqrt{10} = 3.16$.
7. When computing a Poisson probability for a different time interval, we must first convert the mean arrival rate to the time period of interest and then compute the probability.

 **Question** (p260)

Suppose the number of patients who arrive at the emergency room of a large hospital during a 15-minute period on weekday morning is the Poisson distribution with the average number of patients arriving in a 15-minute period of time is 10, compute the probability of one arrival in a 3-minute period.

sol:

An Example Involving Length or Distance Intervals

 Question (p260)

Suppose we are concerned with the occurrence of major defects in a highway one month after resurfacing. We will assume that the probability of a defect is the same for any two highway intervals of equal length and that the occurrence or nonoccurrence of a defect in any one interval is independent of the occurrence or nonoccurrence of a defect in any other interval. Hence, the Poisson distribution can be applied. Suppose we learn that major defects one month after resurfacing occur at the average rate of two per mile. Find the probability of no major defects in a particular three mile section of the highway.

sol:

5.7 Hypergeometric Probability Distribution

1. The hypergeometric probability distribution is closely related to the binomial distribution.
2. The two probability distributions differ in two key ways. With the hypergeometric distribution, the trials are not independent; and the probability of success changes from trial to trial.
3. Notation for the hypergeometric distribution, r denotes the number of elements in the population of size N labeled success, and $N-r$ denotes the number of elements in the population labeled failure.

4. The hypergeometric probability function is used to compute the probability that in a random selection of n elements, selected without replacement, we obtain x elements labeled success and $n-x$ elements labeled failure.
5. For this outcome to occur, we must obtain x successes from the r successes in the population and $n-x$ failures from the $N-r$ failures.
6. **Hypergeometric Probability Function** (Denote by $X \sim Hyp(N, r, n)$):

$$f(x) = \frac{\binom{r}{x} \binom{N-r}{n-x}}{\binom{N}{n}}$$

where:

x = number of successes, $x = 0, 1, 2, \dots, n$

n = number of trials

$f(x)$ = probability of x successes in n trials

N = number of elements in the population

r = number of elements in the population labeled success, $x \leq r$

$\binom{N}{n}$: the number of ways n elements can be selected from a population of size N

$\binom{r}{x}$: the number of ways that x successes can be selected from a total of r successes in the

$\binom{N-r}{n-x}$: the number of ways that $n-x$ failures can be selected from a total of $N-r$ failures in

7. The mean and variance of a hypergeometric distribution are

$$E(X) = \mu = \frac{n \left(\frac{r}{N}\right)}{n \left(\frac{r}{N}\right) \left(1 - \frac{r}{N}\right) \left(\frac{N-n}{N-1}\right)}$$

 Question (p263)

Electric fuses produced by Ontario Electric are packaged in boxes of 12 units each. Suppose an inspector randomly selects three of the 12 fuses in a box for testing.

(a) If the box contains exactly five defective fuses, what is the probability that the inspector will find exactly one of the three fuses defective? (b) What is the probability of finding at least one defective fuse. (c) Find the mean and variance for the number of defective fuses.

sol:

 **EXERCISES**

5.1 : 1, 3, 5

5.2 : 9, 10

5.3 : 16, 21, 23

5.4 : 26, 28, 29

5.5 : 32, 35, 39, 42

5.6 : 45, 49, 50

5.7 : 55, 57, 58

SUP : 59, 62, 66, 73

 **See Also:** SUMMARY, GLOSSARY, KEY FORMULAS

“我們每天做什麼，就會成為什麼樣的人。因此，卓越不是一種行為，而是一個習慣”

“We are what we repeatedly do. Excellence, therefore, is not an act, but a habit.”

— *Aristotle (384–322 BC)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 6: Continuous Probability Distributions

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

系級: _____ 學號: _____ 姓名: _____

Overview

1. For a discrete random variable, the probability function $f(x)$ provides the probability that the random variable assumes a particular value.
2. With continuous random variables, the counterpart of the probability function is the probability density function (pdf), also denoted by $f(x)$.
3. The area under the graph of $f(x)$ corresponding to a given interval does provide the probability that the continuous random variable X assumes a value in that interval.
4. The normal distribution is of major importance because of its wide applicability and its extensive use in statistical inference.
5. The exponential distribution is useful in applications involving such factors as waiting times and service times.

6.1 Uniform Probability Distribution

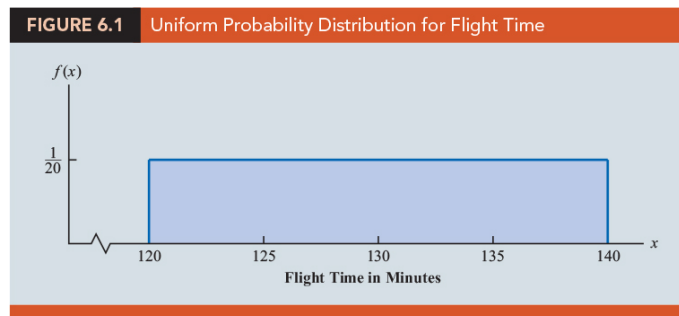
1. **Uniform Probability Density Function:** the uniform probability density function for a random variable X (denoted by $(X \sim U(a, b))$):

$$f(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{elsewhere} \end{cases}$$

2. For a continuous random variable, we consider probability only in terms of the likelihood that a random variable assumes a value within a specified interval.
3. **Example** Flight Time Example:

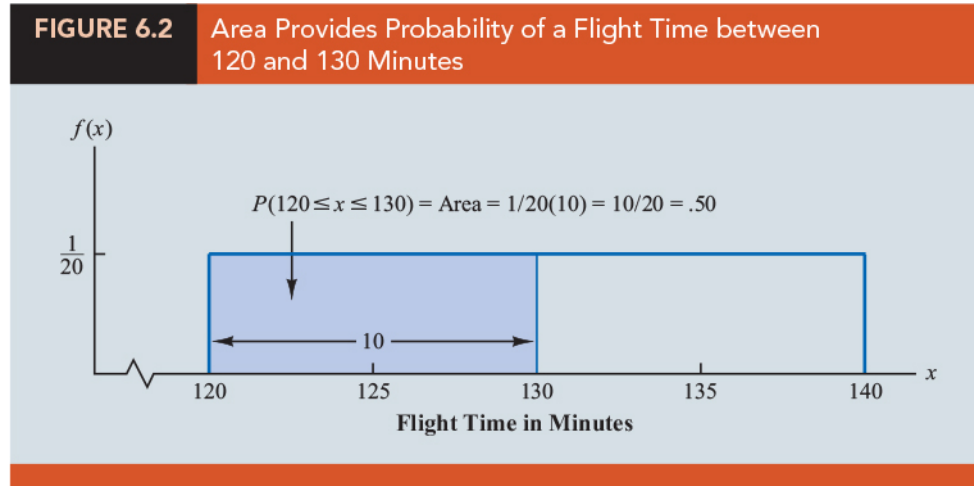
- (a) Consider the random variable X representing the flight time in the interval from 120 minutes to 140 minutes of an airplane traveling from Chicago to New York.
- (b) Because the random variable X can assume any value in that interval, X is a continuous rather than a discrete random variable.
- (c) Let us assume that sufficient actual flight data are available to conclude that the probability of a flight time within any 1-minute interval is the same as the probability of a flight time within any other 1-minute interval contained in the larger interval from 120 to 140 minutes.
- (d) With every 1-minute interval being equally likely, the random variable X is said to have a uniform probability distribution.
- (e) (Figure 6.1) The probability density function, which defines the uniform distribution for the flight-time random variable, is

$$f(x) = \frac{1}{20}, 120 \leq x \leq 140$$



Area as a Measure of Probability

1. (Figure 6.2) The area under the graph of $f(x)$ in the interval from 120 to 130 provides probability of a flight time between 120 and 130 minutes: $10(1/20) = 10/20 = 0.50$.



2. Once a probability density function $f(x)$ is identified, the probability that X takes a value between some lower value a and some higher value b can be found by computing the area under the graph of $f(x)$ over the interval from a to b .
3. Note that $P(a \leq X \leq b) = (b - a) \times (1/(b - a)) = 1$; that is, the total area under the graph of $f(x)$ is equal to 1. This property holds for all continuous probability distributions.
4. For a continuous probability density function, we must also require that $f(x) \geq 0$ for all values of X .
5. Given the uniform distribution for flight time, what is the probability of a flight time between 128 and 136 minutes?

$$\underline{P(128 \leq X \leq 136) = 8(1/20) = 0.40}$$

sol:

6. Two major differences stand out between the treatment of continuous random variables and the treatment of their discrete counterparts.

- (a) We no longer talk about the probability of the random variable assuming a particular value. Instead, we talk about the probability of the random variable assuming a value within some given interval.
- (b) The probability of a continuous random variable assuming a value within some given interval from a to b is defined to be the area under the graph of the probability density function between a and b .
7. Because a single point is an interval of zero width, this implies that the probability of a continuous random variable assuming any particular value exactly is zero.
8. The expected value and variance for a continuous uniform random variable is

$$E(X) = \frac{a + b}{2}, \quad Var(X) = \frac{(b - a)^2}{12}.$$

9. **Example** Applying $E(X)$ and $Var(X)$ to the uniform distribution for flight times from Chicago to New York, we obtain

$$E(X) = \frac{120 + 140}{2} = 130, \quad Var(X) = \frac{(140 - 120)^2}{12} = 33.33$$

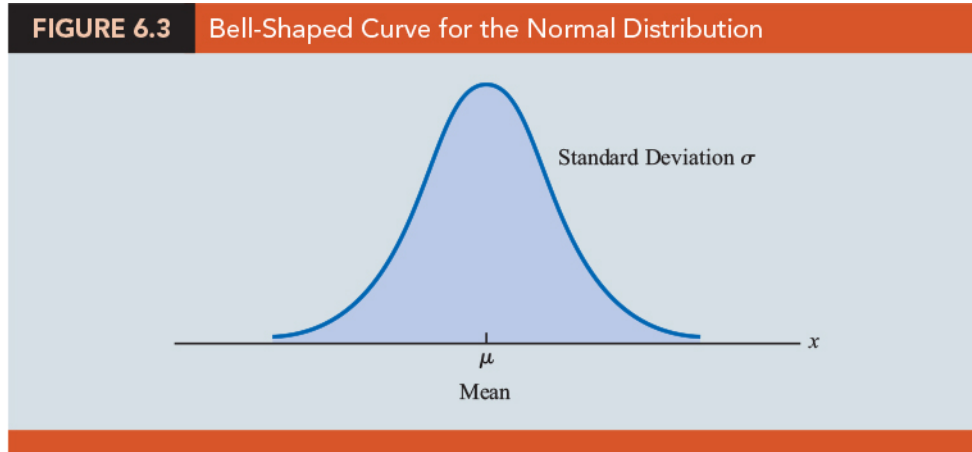
The standard deviation of flight times can be found by taking the square root of the variance. Thus, $\sigma = 5.77$ minutes.

6.2 Normal Probability Distribution

- The normal distribution (Gaussian distribution) has been used in a wide variety of practical applications in which the random variables are heights and weights of people, test scores, scientific measurements, amounts of rainfall, and other similar values.
- It is also widely used in statistical inference.

Normal Curve

- (Figure 6.3) The form, or shape, of the normal distribution is illustrated by the bellshaped normal curve.



- Normal Probability Density Function:**

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

where

μ = mean

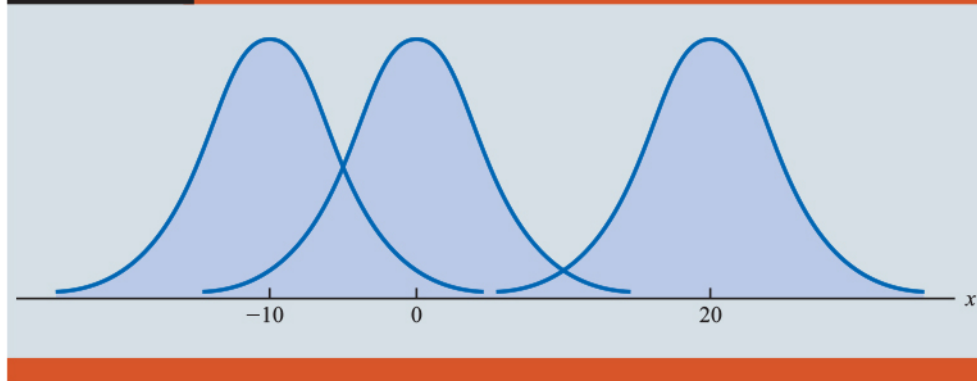
σ = standard deviation

π = 3.14159

e = 2.71828

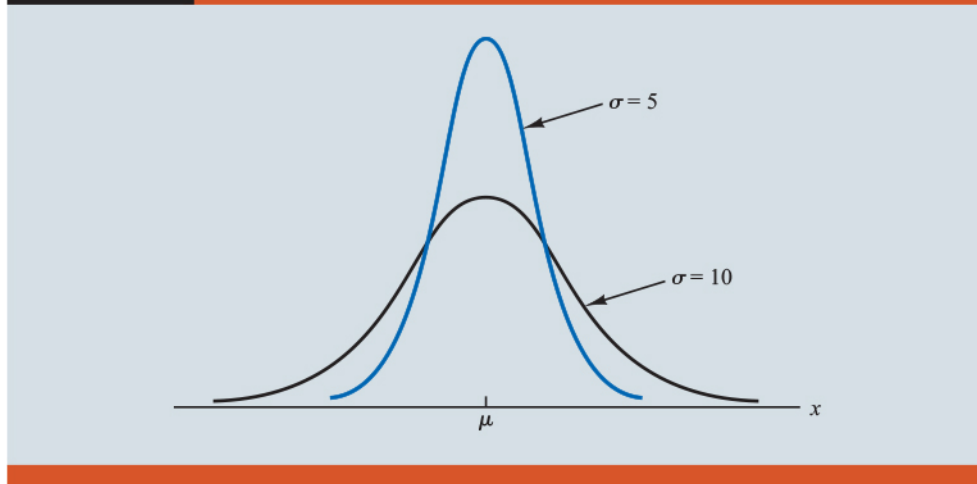
- Several observations about the characteristics of the normal distribution.
 - The entire family of normal distributions is differentiated by two parameters: the mean μ and the standard deviation σ .
 - The highest point on the normal curve is at the mean, which is also the median and mode of the distribution.
 - (Figure 6.4) The mean of the distribution can be any numerical value: negative, zero, or positive.

FIGURE 6.4 Normal Distributions with Same Standard Deviation and Different Means



- (d) The normal distribution is symmetric, with the shape of the normal curve to the left of the mean a mirror image of the shape of the normal curve to the right of the mean.
- (e) The tails of the normal curve extend to infinity in both directions and theoretically never touch the horizontal axis. Because it is symmetric, the normal distribution is not skewed; its skewness measure is zero.
- (f) (Figure 6.5) The standard deviation determines how flat and wide the normal curve is. Larger values of the standard deviation result in wider, flatter curves, showing more variability in the data.

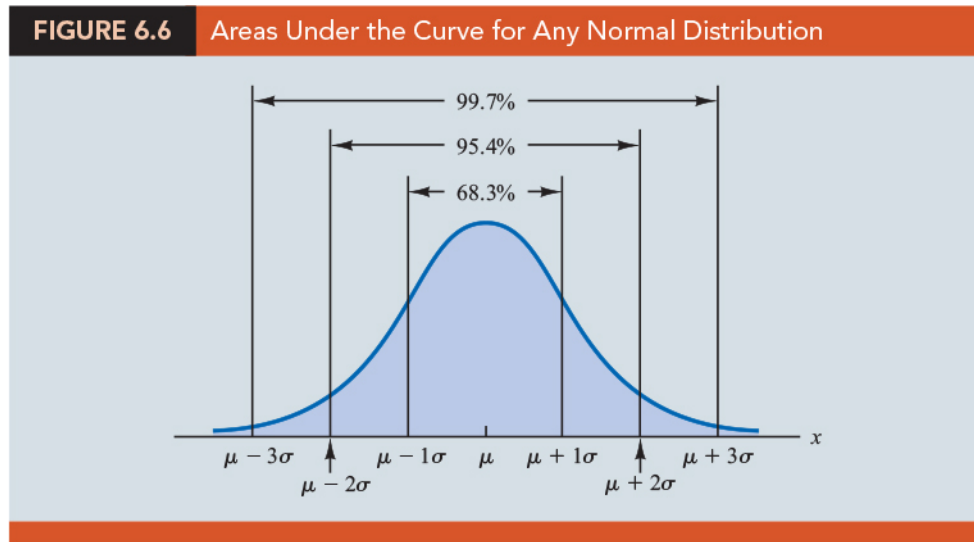
FIGURE 6.5 Normal Distributions with Same Mean and Different Standard Deviations



- (g) Probabilities for the normal random variable are given by areas under the normal curve. The total area under the curve for the normal distribution is

1. Because the distribution is symmetric, the area under the curve to the left (right) of the mean is 0.50 (0.50).

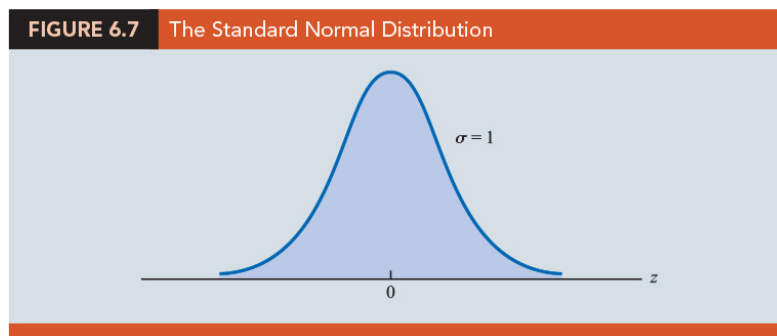
- (h) (Figure 6.6) The percentage of values in some commonly used intervals: 68.3% (95.4%, 99.7%) of the values of a normal random variable are within plus or minus one (two, three) standard deviation of its mean.



Standard Normal Probability Distribution

1. (Figure 6.7) A random variable, Z , that has a normal distribution with a mean of zero and a standard deviation of one is said to have a standard normal probability distribution.

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

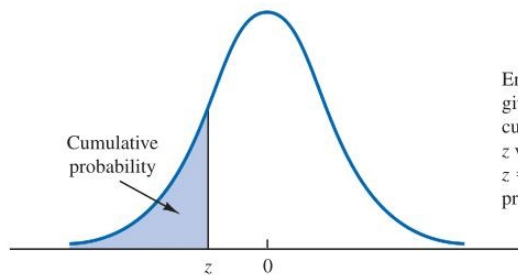


2. The three types of probabilities we need to compute:

- (a) the probability that the standard normal random variable Z will be less than or equal to a given value z ; $P(Z \leq z)$
- (b) the probability that Z will be between two given values z_1, z_2 ; $P(z_1 \leq Z \leq z_2)$
- (c) the probability that Z will be greater than or equal to a given value z ; $P(Z \geq z)$

3. (Appendix B, Table 1) Cumulative probabilities for the standard normal probability:

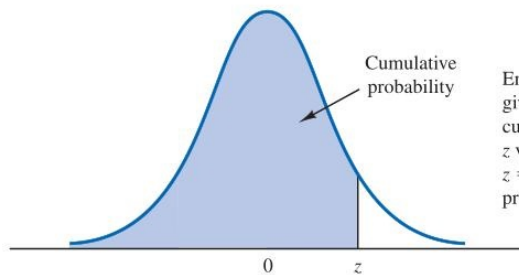
TABLE 1 Cumulative Probabilities for the standard Normal Distribution



Entries in the table give the area under the curve to the left of the z value. For example, for $z = -.85$, the cumulative probability is .1977.

z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
-3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0011	.0010	.0010
-2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
-2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
-2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
-2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
-2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
-2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064


TABLE 1 Cumulative Probabilities for the standard Normal Distribution (Continued)



Entries in the table give the area under the curve to the left of the z value. For example, for $z = 1.25$, the cumulative probability is .8944.

z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
.2	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141
.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
.4	.6554	.6591	.6628	.6665	.6702	.6739	.6776	.6812	.6849	.6886
.5	.6923	.6959	.6995	.7031	.7068	.7104	.7140	.7177	.7213	.7250
.6	.7286	.7322	.7358	.7394	.7429	.7465	.7500	.7535	.7570	.7605
.7	.7640	.7675	.7710	.7744	.7779	.7814	.7849	.7883	.7918	.7952
.8	.7987	.8021	.8055	.8089	.8123	.8157	.8191	.8225	.8259	.8293
.9	.8327	.8360	.8394	.8428	.8461	.8495	.8528	.8561	.8594	.8627
1.0	.8659	.8691	.8724	.8756	.8788	.8820	.8851	.8882	.8913	.8944
1.1	.8974	.9005	.9035	.9065	.9095	.9125	.9154	.9183	.9212	.9241
1.2	.9270	.9298	.9327	.9354	.9381	.9408	.9434	.9461	.9487	.9513
1.3	.9539	.9564	.9589	.9613	.9638	.9661	.9685	.9708	.9731	.9754
1.4	.9776	.9798	.9819	.9841	.9861	.9881	.9900	.9919	.9937	.9955

4. The table of cumulative probabilities for the standard normal probability distribution can be used to find probabilities associated with values of the standard normal random variable Z .
5. Two types of questions can be asked.
 - (a) The first type of question specifies a value, or values, for Z and asks us to use the table to determine the corresponding areas or probabilities.
 - (b) The second type of question provides an area, or probability, and asks us to use the table to determine the corresponding Z value.

 **Question** (p290)

Using the standard normal probability Table to compute

1. (Figure 6.8) the probability that the standard normal random variable Z is less than or equal to 1.00;
2. (Figure 6.9) the probability that Z is in the interval between -0.50 and 1.25 ;
3. (Figure 6.10) the probability that Z is within one standard deviation of the mean;
4. (Figure 6.11) the probability of obtaining a Z value of at least 1.58.

sol:

1. $P(Z \leq 1.00) =$
2. $P(-0.50 \leq Z \leq 1.25) =$
3. $P(-1.00 \leq Z \leq 1.00) =$
4. $P(Z \geq 1.58) =$

FIGURE 6.8 Cumulative Probability for Normal Distribution Corresponding to $P(z \leq 1.00)$

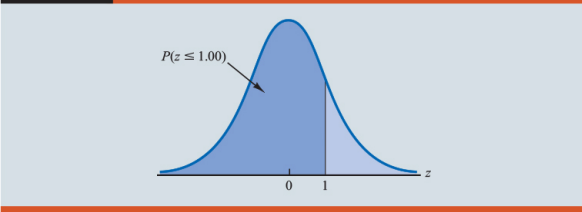


FIGURE 6.9 Cumulative Probability for Normal Distribution Corresponding to $P(-.50 \leq z \leq 1.25)$

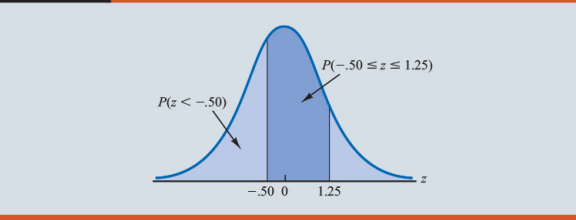


FIGURE 6.10 Cumulative Probability for Normal Distribution Corresponding to $P(-1.00 \leq z \leq 1.00)$

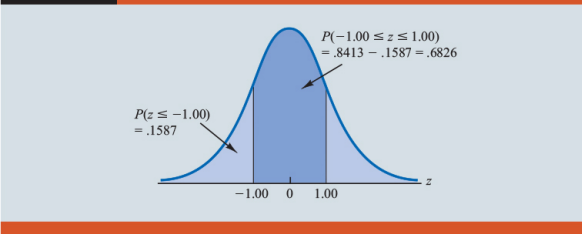
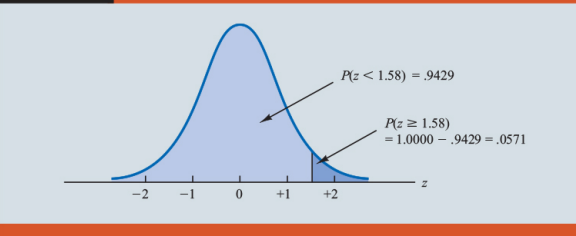



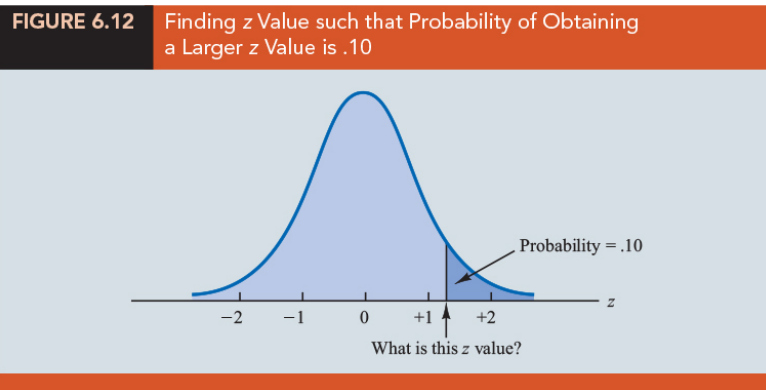
FIGURE 6.11 Cumulative Probability for Normal Distribution Corresponding to $P(z \geq 1.58)$



 **Question** (p293)

(Figure 6.12) Find a Z value such that the probability of obtaining a larger Z value is 0.10.

sol:



Computing Probabilities for Any Normal Probability Distribution

1. Probabilities for all normal distributions can be computed using the standard normal distribution.
2. Convert any normal random variable X with mean μ and standard deviation σ to the standard normal random variable Z :

$$Z = \frac{X - \mu}{\sigma}$$

3. We can interpret Z as the number of standard deviations that the normal random variable X is from its mean μ .

 **Question** (p294)

Suppose we have a normal random variable X with $\mu = 10$ and $\sigma = 2$. What is the probability that the random variable X is between 10 and 14?

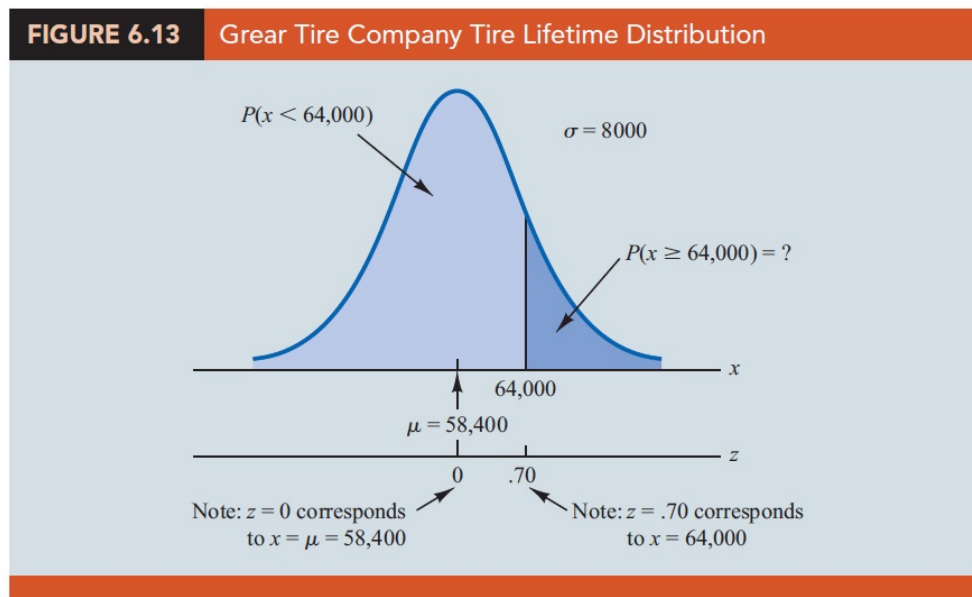
sol:


Grear Tire Company Problem

 Question (p294)

(Grear Tire Company Problem) (Figure 6.13) Suppose the Grear Tire Company developed a new steelbelted radial tire to be sold through a national chain of discount stores. Because the tire is a new product, Grear's managers believe that the mileage guarantee offered with the tire will be an important factor in the acceptance of the product. Before finalizing the tire mileage guarantee policy, Grear's managers want probability information about X = number of miles the tires will last. From actual road tests with the tires, Grear's engineering group estimated that the mean tire mileage is $\mu = 36,500$ miles and that the standard deviation is $\sigma = 5000$. In addition, the data collected indicate that a normal distribution is a reasonable assumption. What percentage of the tires can be expected to last more than 40,000 miles? In other words, what is the probability that the tire mileage, X , will exceed 40,000?

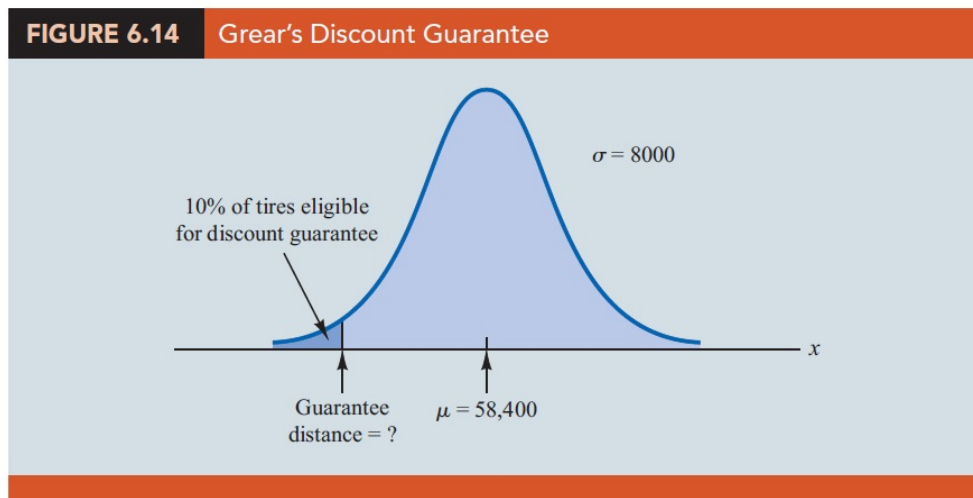
sol:



 Question (p294)

(Gear Tire Company Problem) (Figure 6.14) Let us now assume that Gear is considering a guarantee that will provide a discount on replacement tires if the original tires do not provide the guaranteed mileage. What should the guarantee mileage be if Gear wants no more than 10% of the tires to be eligible for the discount guarantee?

sol:



1. **The important role of the probability distributions:** providing decision - making information. Namely, once a probability distribution is established for a particular application, it can be used to obtain probability information about the problem.
2. Probability does not make a decision recommendation directly, but it provides information that helps the decision maker better understand the risks and uncertainties associated with the problem.

6.3 Normal Approximation of Binomial Probabilities

1. Recall that a binomial experiment:
 - (a) consists of a sequence of n identical independent trials with each trial having two possible outcomes, a success or a failure.
 - (b) The probability of a success on a trial is the same for all trials and is denoted by p .
 - (c) The binomial random variable is the number of successes in the n trials, and probability questions pertain to the probability of x successes in the n trials.
2. In cases where $np \geq 5$, and $n(1-p) \geq 5$, the normal distribution provides an approximation of binomial probabilities.
3. When using the normal approximation to the binomial, we set $\mu = np$ and $\sigma = \sqrt{np(1-p)}$ in the definition of the normal curve.
4. **Continuity correction:** Recall that, with a continuous probability distribution, probabilities are computed as areas under the probability density function. As a result, the probability of any single value for the random variable is zero. Thus to approximate the binomial probability of x successes, we compute the area under the corresponding normal curve between $x - 0.5$ and $x + 0.5$. The 0.5 that we add and subtract from x is called a continuity correction factor.
5. **補充:** 理論上以常態機率逼近二項式機率，可以使用的情境是：「 n 要很大且 p 不要近於 0 或 1」。但即使「 n 有點小且 p 近於 1/2」，逼近效果也是可在接受範圍。依據上面兩點，不同教科書裡有一些不盡相同的經驗法則 (rules of thumb):
 - (a) $np \geq 5$ 且 $n(1-p) \geq 5$ 。[我們的教科書]
 - (b) $n > 5$, 且 $|skewness| < 1/3$ 。
 - (c) $(np) \pm \sqrt{np(1-p)} \in (0, n)$ 。(等同 $np > 9(1-p)$ 且 $n(1-p) > 9p$)。
 - (d) $np(1-p) \geq 10$
 - (e) n 很大且 $np \geq 10$ 且 $n(1-p) \geq 10$

在下面篇 1945 年的論文裡，證明中的確有提到「 $\sqrt{np(1-p)} > 5$ 」的條件。W. Feller, (1945) On the Normal Approximation to the Binomial Distribution, Ann. Math. Statist. 16(4): 319-329 (December, 1945). DOI: 10.1214/aoms/1177731058。

6. **補充**: 使用「連續修正項」需注意「有沒有等號」:

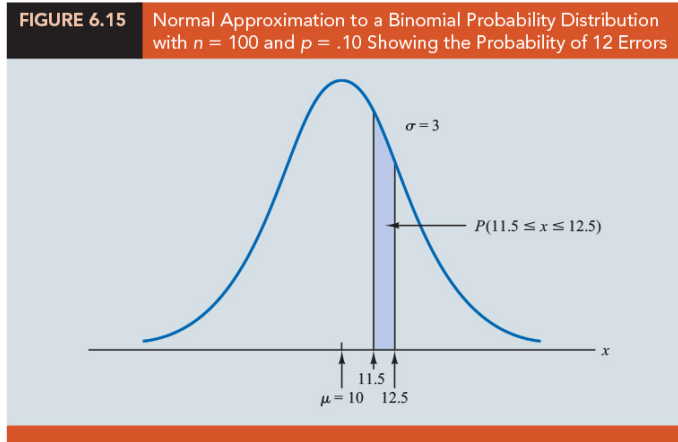
Binomial	Normal
If $P(X = n)$	use $P(n - 0.5 < X < n + 0.5)$
If $P(X > n)$	use $P(X > n + 0.5)$
If $P(X \leq n)$	use $P(X < n + 0.5)$
If $P(X < n)$	use $P(X < n - 0.5)$
If $P(X \geq n)$	use $P(X > n - 0.5)$

7. **補充**: 「Normal Approximation + continuity correction」R 程式範例：
https://www.macmillanlearning.com/studentresources/college/statistics/introstats3e/optional_sections/00_KokoskaIntroStat3e_04962_ch06.5_online_001_010_4PP.pdf

 **Question** (p299)

(Figure 6.15) Suppose that a particular company has a history of making errors in 10% of its invoices. A sample of 100 invoices has been taken, and we want to compute the probability that 12 invoices contain errors. That is, find the binomial probability of 12 successes in 100 trials by applying the normal approximation.

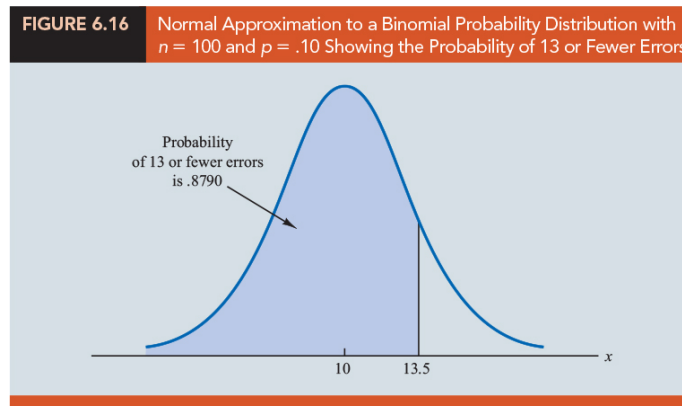
sol:



Question (p300)

(Figure 6.16) Following the previous question, compute the probability of 13 or fewer errors in the sample of 100 invoices.

sol:



6.4 Exponential Probability Distribution

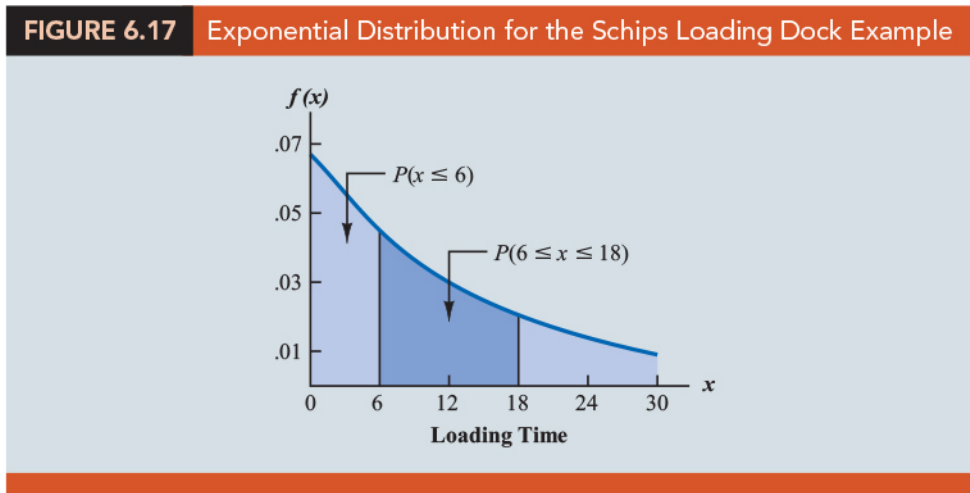
- The exponential probability distribution may be used for random variables such as the time between arrivals at a hospital emergency room, the time required to load a truck, the distance between major defects in a highway, and so on.
- The exponential probability density function (denoted by $X \sim \text{Exp}(\mu)$):

$$f(x) = \frac{1}{\mu} e^{-\frac{x}{\mu}}, \quad x \geq 0,$$

where $\mu =$ expected value or mean.

- (Figure 6.17) (Schips loading dock example) Suppose that X represents the loading time for a truck at the Schips loading dock and follows a exponential distribution. If the mean, or average, loading time is 15 minutes ($\mu = 15$), the appropriate probability density function for X is

$$f(x) = \frac{1}{15} e^{-\frac{x}{15}}$$



- The mean and the standard deviation of the exponential distribution are equal.
- (Figure 6.17) The probability that loading a truck will take 6 minutes or less, $P(X \leq 6)$, is defined to be the area under the curve in Figure 6.17 from $x = 0$ to $x = 6$. The probability that the loading time will be between 6 minutes and 18 minutes, $P(6 \leq x \leq 18)$, is given by the area under the curve from $x = 6$ to $x = 18$.

Computing Probabilities for the Exponential Distribution

- The cumulative probability of obtaining a value for the exponential random variable of less than or equal to some specific value denoted by x :

$$F(x) = P(X \leq x) = \underline{1 - e^{-\frac{x}{\mu}}}$$

 **Question** (p303)

In the Schips loading dock example, Find (a) the probability that loading a truck will take 6 minutes or less; (b) the probability that the loading time will be 18 minutes or less; (c) the probability that the loading time will be between 6 minutes and 18 minutes.

sol:

Relationship Between the Poisson and Exponential Distributions

- (Section 5.6) The Poisson distribution is a discrete probability distribution that is often useful in examining the number of occurrences of an event over a specified interval of time or space.

$$f(x) = \underline{\frac{\mu^x e^{-\mu}}{x!}},$$

where μ = expected value or mean number of occurrences over a specified interval.

2. If the Poisson distribution provides an appropriate description of the number of occurrences per interval, the exponential distribution provides a description of the length of the interval between occurrences.

3. **Example** Suppose the number of patients who arrive at a hospital emergency room during one hour is described by a Poisson probability distribution with a mean of 10 patients per hour. The Poisson probability function that gives the probability of x arrivals per hour is

$$f(x) = \frac{10^x e^{-10}}{x!}$$

Because the average number of arrivals is 10 patients per hour, the average time between patients arriving is

$$\frac{1 \text{ hour}}{10 \text{ patient}} = 0.1 \text{ hour/patient.}$$

Thus, the corresponding exponential distribution that describes the time between the arrivals has a mean of $\mu = 0.1$ hour per patient; as a result, the appropriate exponential probability density function is

$$f(x) = \frac{1}{0.1} e^{-x/0.1} = 10e^{-10x}.$$

😊 R functions:

1. `dnorm`, `pnorm`, `qnorm`, `rnorm`:

<https://stat.ethz.ch/R-manual/R-devel/library/stats/html/Normal.html>

2. `dexp`, `pexp`, `qexp`, `rexp`:

<https://stat.ethz.ch/R-manual/R-devel/library/stats/html/Exponential.html>

☺ **EXERCISES**

6.1 : 2, 6

6.2 : 11, 14, 19, 23

6.3 : 27, 30

6.4 : 33, 37, 38

SUP : 39, 42, 49, 52, 53

“樹林裡有兩條岔路，而我走了人跡較少的那一條，因此有了完全不同的人生”

“Try not to become a man of success, but rather try to become a man of value”

— *Robert Frost (March 26, 1874 – January 29, 1963)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 7: Sampling and Sampling Distributions

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

系級: _____ 學號: _____ 姓名: _____

Overview

- Definitions:
 - An element is the entity on which data are collected.
 - A population is the collection of all the elements of interest.
 - A sample is a subset of the population.
- Why select a sample? Answer: collect data to make an inference and answer research questions about a population.
- The sample contains only a portion of the population. Sample results provide only estimates of the values of the corresponding population characteristics.
- Some sampling error is to be expected. With proper sampling methods, the sample results will provide "good" estimates of the population parameters.
- Numerical characteristics of a population are called parameters.
- This chapter:

- (a) show how simple random sampling can be used to select a sample from a finite population.
- (b) describe how a random sample can be taken from an infinite population that is generated by an ongoing process.
- (c) how to compute estimates of a population mean, a population standard deviation, and a population proportion.
- (d) introduce the important concept of a sampling distribution. It enables us to make statements about how close the sample estimates are to the corresponding population parameters.
- (e) discuss some alternatives to simple random sampling that are often employed in practice and the ramifications of large samples on sampling distributions.

7.1 The Electronics Associates Sampling Problem

1. *Background:* The director of personnel for Electronics Associates, Inc. (EAI), has been assigned the task of developing a profile of the company's 2500 managers (population, $N = 2500$).
2. *Interested:* The characteristics to be identified include the mean annual salary for the managers and the proportion of managers having completed the company's management training program.
3. *Population Data:* The population contains the annual salary and the training program status for all 2500 managers (referring to the firm's personnel records) (Dataset: EAI).
4. *Population Parameters:* The population mean annual salary ($\mu = \$71,800$), the population standard deviation of annual salary ($\sigma = \$4000$), and the population proportion that completed the training program ($p = 1500/2500 = 0.60$) are parameters of the population of EAI managers.

5. *Question:* Suppose that the necessary information on all the EAI managers was not readily available in the company's database. How the firm's director of personnel can obtain estimates of the population parameters by using a sample of managers, say a sample of 30 managers ($n = 30$).
6. *Sample Data:* If the personnel director could be assured that a sample of 30 managers would provide adequate information about the population of 2500 managers, the time and the cost of developing a profile would be substantially less.

7.2 Selecting a Sample

Sampling from a Finite Population

1. **A simple random sample:** A simple random sample of size n from a finite population of size N is a sample selected such that each possible sample of size n has the same probability of being selected.
2. One procedure for selecting a simple random sample from a finite population is to use a table of random numbers to choose the elements for the sample one at a time in such a way that, at each step, each of the elements remaining in the population has the same probability of being selected.
3. **Example** EAI managers
 - (a) First assign the managers the numbers 1 to 2500 in the order that their names appear in the EAI personnel file.
 - (b) (Table 7.1) Select random numbers from anywhere in the table and move systematically in a direction of our choice in sets or groups of four digits (largest number in the population list).
 - (c) E.g., Use the first row of Table 7.1 and move from left to right. The first 7 four-digit random numbers are

6327 1599 8671 7445 1102 1514 1807

63271	59986	71744	51102	15141	80714	58683	93108	13554	79945
88547	09896	95436	79115	08303	01041	20030	63754	08459	28364
55957	57243	83865	09911	19761	66535	40102	26646	60147	15702
46276	87453	44790	67122	45573	84358	21625	16999	13385	22782
55363	07449	34835	15290	76616	67191	12777	21861	68689	03263
69393	92785	49902	58447	42048	30378	87618	26933	40640	16281
13186	29431	88190	04588	38733	81290	89541	70290	40113	08243
17726	28652	56836	78351	47327	18518	92222	55201	27340	10493
36520	64465	05550	30157	82242	29520	69753	72602	23756	54935
81628	36100	39254	56835	37636	02421	98063	89641	64953	99337
84649	48968	75215	75498	49539	74240	03466	49292	36401	45525
63291	11618	12613	75055	43915	26488	41116	64531	56827	30825
70502	53225	03655	05915	37140	57051	48393	91322	25653	06543
06426	24771	59935	49801	11082	66762	94477	02494	88215	27191
20711	55609	29430	70165	45406	78484	31639	52009	18873	96927
41990	70538	77191	25860	55204	73417	83920	69468	74972	38712
72452	36618	76298	26678	89334	33938	95567	29380	75906	91807
37042	40318	57099	10528	09925	89773	41335	96244	29002	46453
53766	52875	15987	46962	67342	77592	57651	95508	80033	69828
90585	58955	53122	16025	84299	53310	67380	84249	25348	04332
32001	96293	37203	64516	51530	37069	40261	61374	05815	06714
62606	64324	46354	72157	67248	20135	49804	09226	64419	29457
10078	28073	85389	50324	14500	15562	64165	06125	71353	77669
91561	46145	24177	15294	10061	98124	75732	00815	83452	97355
13091	98112	53959	79607	52244	63303	10413	63839	74762	50289

- (d) The first number, 6327, is greater than 2500: discarded. The first manager selected for the random sample is the 2nd number 1599 on the list of EAI managers.
- (e) This process continues until the simple random sample of 30 EAI managers has been obtained.
- (f) Any previously used random numbers are ignored because the corresponding manager is already included in the sample.
4. If we selected a sample such that previously used random numbers are acceptable and specific managers could be included in the sample two or more times, we would be sampling with replacement. Otherwise, it is referred to as sampling without replacement.
5. Sampling without replacement is the sampling procedure used most often in practice. When we refer to simple random sampling, we will assume the sampling is

without replacement.

Sampling from an Infinite Population

1. *The infinite population case:* The population is infinitely large or the elements of the population are being generated by an ongoing process for which there is no limit on the number of elements that can be generated. Thus, it is not possible to develop a list of all the elements in the population.

2. A random sample of size n from an infinite population is a sample selected such that two conditions are satisfied.

(a) Each element selected comes from the same population.

(b) Each element is selected independently.

3. The purpose of the second condition of the random sample selection procedure (each element is selected independently) is to prevent selection bias.

4. **Example 1**: Quality Control Inspection

(a) Consider a production line designed to fill boxes of a breakfast cereal with a mean weight of 24 ounces of breakfast cereal per box.

(b) Samples of 12 boxes filled by this process are periodically selected by a quality control inspector to determine if the process is operating properly or if, perhaps, a machine malfunction has caused the process to begin underfilling or overfilling the boxes.

(c) *Condition 1:* The boxes must be selected at approximately the same point in time. This way the inspector avoids the possibility of selecting some boxes when the process is operating properly and other boxes when the process is not operating properly and is underfilling or overfilling the boxes.

(d) *Condition 2:* designing the production process so that each box of cereal is filled independently.

5. **Example 2**: Customers Arrival

(a) Suppose an employee is asked to select and interview a sample of customers in order to develop a profile of customers who visit the fast-food restaurant.

- (b) The customer arrival process is ongoing and there is no way to obtain a list of all customers in the population. So, for practical purposes, the population for this ongoing process is considered infinite.
 - (c) *Condition 1:* The employee collects the sample from people who come into the restaurant and make a purchase to ensure that the same population condition is satisfied. the person came into the restaurant to use the restroom would not be a customer.
 - (d) *Condition 2:* Selection bias (a particular age group, a group of customers would be likely to exhibit similar characteristics) can be avoided by ensuring that the selection of a particular customer does not influence the selection of any other customer.
6. Situations involving sampling from an infinite population are usually associated with a process that operates over time. Examples include parts being manufactured on a production line, repeated experimental trials in a laboratory, transactions occurring at a bank, telephone calls arriving at a technical support center, and customers entering a retail store.
7. As long as the sampled elements are selected from the same population and are selected independently, the sample is considered a random sample from an infinite population.

7.3 Point Estimation

1. Glossary:

- (a) **Parameter:** A numerical characteristic of a population, such as a population mean μ , a population standard deviation σ , a population proportion p , and so on.

- (b) **Random sample:** A random sample from an infinite population is a sample selected such that each element selected comes from the same population and each element is selected independently.
- (c) **Sampling without replacement:** Once an element has been included in the sample, it is removed from the population and cannot be selected a second time.
- (d) **Sampling with replacement:** Once an element has been included in the sample, it is returned to the population. A previously selected element can be selected again and therefore may appear in the sample more than once.
- (e) **Simple random sample:** A simple random sample of size n from a finite population of size N is a sample selected such that each possible sample of size n has the same probability of being selected.
- (f) **Sample statistic:** A sample characteristic, such as a sample mean \bar{x} , a sample standard deviation s , a sample proportion \bar{p} , and so on. The value of the sample statistic is used to estimate the value of the corresponding population parameter.
2. **Example** The EAI problem.
- (a) To estimate the value of a population parameter (e.g., the population mean μ and the population standard deviation σ for the annual salary of EAI managers), we compute a corresponding characteristic of the sample, referred to as a sample statistic.
- (b) (Table 7.2) A simple random sample of 30 managers and the corresponding data on annual salary $(x_1, x_2, \dots, x_{30})$ and management training program participation (Yes, No).

TABLE 7.2 Annual Salary and Training Program Status for a Simple Random Sample of 30 EAI Managers

Annual Salary (\$)	Management Training Program	Annual Salary (\$)	Management Training Program
$x_1 = 69,094.30$	Yes	$x_{16} = 71,766.00$	Yes
$x_2 = 73,263.90$	Yes	$x_{17} = 72,541.30$	No
$x_3 = 69,643.50$	Yes	$x_{18} = 64,980.00$	Yes
$x_4 = 69,894.90$	Yes	$x_{19} = 71,932.60$	Yes
$x_5 = 67,621.60$	No	$x_{20} = 72,973.00$	Yes
$x_6 = 75,924.00$	Yes	$x_{21} = 65,120.90$	Yes
$x_7 = 69,092.30$	Yes	$x_{22} = 71,753.00$	Yes
$x_8 = 71,404.40$	Yes	$x_{23} = 74,391.80$	No
$x_9 = 70,957.70$	Yes	$x_{24} = 70,164.20$	No
$x_{10} = 75,109.70$	Yes	$x_{25} = 72,973.60$	No
$x_{11} = 65,922.60$	Yes	$x_{26} = 70,241.30$	No
$x_{12} = 77,268.40$	No	$x_{27} = 72,793.90$	No
$x_{13} = 75,688.80$	Yes	$x_{28} = 70,979.40$	Yes
$x_{14} = 71,564.70$	No	$x_{29} = 75,860.90$	Yes
$x_{15} = 76,188.20$	No	$x_{30} = 77,309.10$	No

- (c) We use the data in Table 7.2 to calculate the corresponding sample statistics: the sample mean \bar{x} and the sample standard deviation s :

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{2154420}{30} = \$71,814 \quad ,$$

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{325009260}{29}} = \$3348 \quad ,$$

- (d) The corresponding sample proportion \bar{p} , who completed the management training program:

$$\bar{p} = \frac{19}{30} = 0.63 \quad .$$

- Point estimation:** the preceding computations, we perform the statistical procedure called point estimation.
- We refer to the sample mean \bar{x} as the point estimator of the population mean μ , the sample standard deviation s as the point estimator of the population standard deviation σ , and the sample proportion \bar{p} as the point estimator of the population proportion p .

5. **Point estimator:** The sample statistic, such as \bar{x} , s , or \bar{p} , that provides the point estimate of the population parameter.
6. **Point estimate:** The numerical value obtained for \bar{x} , s , or \bar{p} is called the point estimate.
7. (Table 7.3) The point estimates differ somewhat from the corresponding population parameters. This difference is to be expected because a sample, and not a census of the entire population, is being used to develop the point estimates.

TABLE 7.3 Summary of Point Estimates Obtained from a Simple Random Sample of 30 EAI Managers

Population Parameter	Parameter Value	Point Estimator	Point Estimate
μ = Population mean annual salary	\$71,800	\bar{x} = Sample mean annual salary	\$71,814
σ = Population standard deviation for annual salary	\$4000	s = Sample standard deviation for annual salary	\$3348
p = Population proportion having completed the management training program	.60	\bar{p} = Sample proportion having completed the management training program	.63

Practical Advice

1. Point estimation is a form of statistical inference. We use a sample statistic to make an inference about a population parameter.
2. When making inferences about a population based on a sample, it is important to have a close correspondence between the sampled population and the target population.
3. In some cases, it is not as easy to obtain a close correspondence between the sampled and target populations.
4. Example Consider the case of an amusement park selecting a sample of its customers to learn about characteristics such as age and time spent at the park.

- (a) Suppose all the sample elements were selected on a day when park attendance was restricted to employees of a single company.
 - (b) The sampled population would be composed of employees of that company and members of their families.
 - (c) If the target population we wanted to make inferences about were typical park customers over a typical summer, then we might encounter a significant difference between the sampled population and the target population.
 - (d) In such a case, we would question the validity of the point estimates being made.
 - (e) Park management would be in the best position to know whether a sample taken on a particular day was likely to be representative of the target population.
5. In summary, whenever a sample is used to make inferences about a population, we should make sure that the study is designed so that the sampled population and the target population are in close agreement.
6. Good judgment is a necessary ingredient of sound statistical practice.

7.4 Introduction to Sampling Distributions

1. (Table 7.4) For the simple random sample of 30 EAI managers, the point estimate of μ is $\bar{x} = \$71,814$ and the point estimate of p is $\bar{p} = 0.63$. Suppose we repeat the process of selecting a simple random sample of 30 EAI managers over and over again (e.g., 500 simple random samples), each time computing the values of \bar{x} and \bar{p} .

TABLE 7.4 Values of \bar{x} and \bar{p} from 500 Simple Random Samples of 30 EAI Managers

Sample Number	Sample Mean (\bar{x})	Sample Proportion (\bar{p})
1	71,814	.63
2	72,670	.70
3	71,780	.67
4	71,588	.53
.	.	.
.	.	.
.	.	.
500	71,752	.50

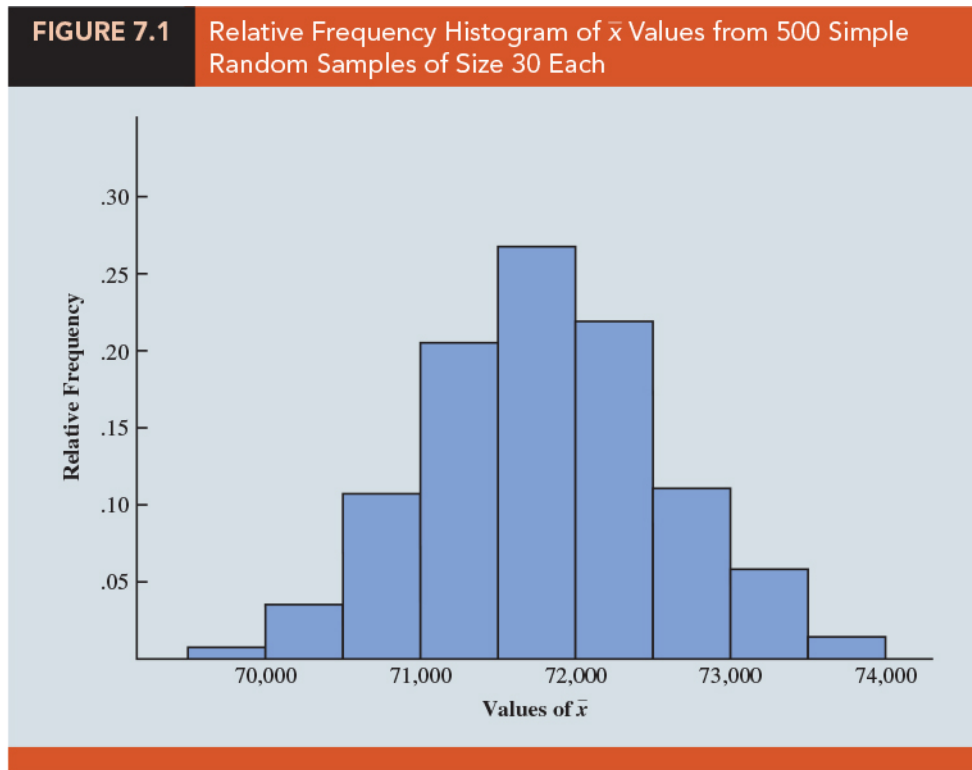
2. (Table 7.5) The frequency and relative frequency distributions for the 500 \bar{x} values.

TABLE 7.5 Frequency and Relative Frequency Distributions of \bar{x} from 500 Simple Random Samples of 30 EAI Managers

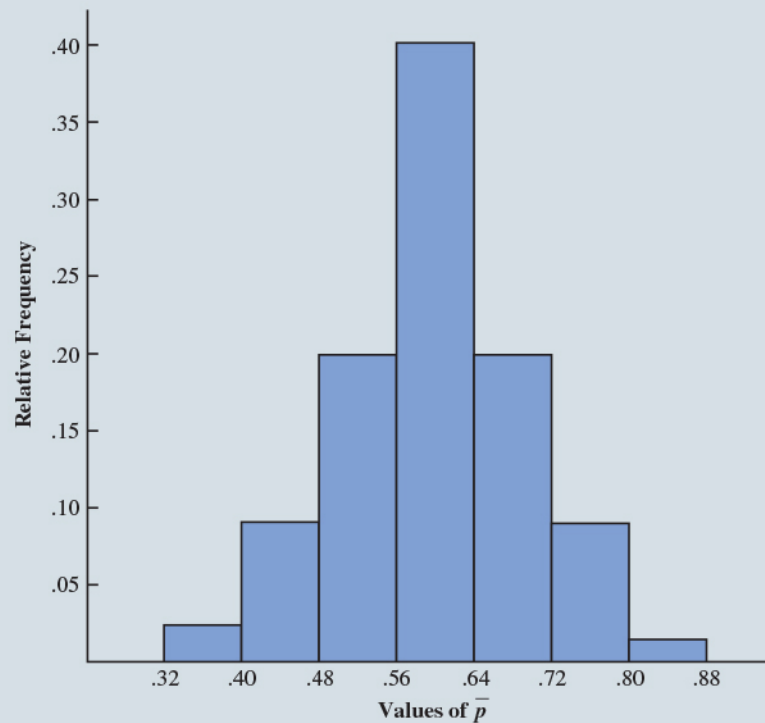
Mean Annual Salary (\$)	Frequency	Relative Frequency
69,500.00–69,999.99	2	.004
70,000.00–70,499.99	16	.032
70,500.00–70,999.99	52	.104
71,000.00–71,499.99	101	.202
71,500.00–71,999.99	133	.266
72,000.00–72,499.99	110	.220
72,500.00–72,999.99	54	.108
73,000.00–73,499.99	26	.052
73,500.00–73,999.99	6	.012
Totals	500	1.000

3. Recall: a random variable is a numerical description of the outcome of an experiment. If we consider the process of selecting a simple random sample as an experiment, the sample mean \bar{x} is a random variable. \bar{x} has a mean or expected value, a standard deviation, and a probability distribution.
4. **The Sampling Distribution:** The various possible values of \bar{x} are the result of different simple random samples, the probability distribution of \bar{x} is called the sampling distribution of \bar{x} .
5. (Figure 7.1) The histogram of 500 \bar{x} values gives an approximation of this sampling distribution. From the approximation we observe the bell-shaped appearance

of the distribution. The largest concentration of the \bar{x} values and the mean of the 500 \bar{x} values are near the population mean $\mu = \$71,800$.



6. (Figure 7.2) The 500 values of the sample proportion \bar{p} are summarized by the relative frequency histogram. The relative frequency histogram of the 500 sample values provides a general idea of the appearance of the sampling distribution of \bar{p} .

FIGURE 7.2 Relative Frequency Histogram of \bar{p} Values from 500 Simple Random Samples of Size 30 Each

7. In practice, we select only one simple random sample from the population. We repeated the sampling process many times and that the different samples generate a variety of values for the sample statistics \bar{x} and \bar{p} .
8. The probability distribution of any particular sample statistic is called the sampling distribution of the statistic.

7.5 Sampling Distribution of \bar{X}

1. **Sampling Distribution of \bar{x} :** The sampling distribution of \bar{x} is the probability distribution of all possible values of the sample mean \bar{x} .

- The sampling distribution of \bar{x} has an expected value or mean, a standard deviation, and a characteristic shape or form.


Expected Value of \bar{x}

- Different simple random samples result in a variety of values for the sample mean \bar{x} .
- Because many different values of the random variable \bar{x} are possible, we are often interested in the mean of all possible values of \bar{x} that can be generated by the various simple random samples.
- The mean of the \bar{x} random variable is the expected value of \bar{x} .
- Let $E(\bar{x})$ represent the expected value of \bar{x} and μ represent the mean of the population from which we are selecting a simple random sample. It can be shown that with simple random sampling, the expected value of \bar{x} and the population mean are equal. ($E(\bar{x}) = \mu$).
- When the expected value of a point estimator equals the population parameter, we say the point estimator is unbiased. The sample mean \bar{x} is an unbiased estimator of the population mean μ .

Standard Deviation of \bar{x}

- Notations:
 - $\sigma_{\bar{x}}$: the standard deviation of \bar{x} . σ : the standard deviation of the population.
 - n : the sample size. N : the population size.
- The two formulas for the standard deviation of \bar{x} :
 - Finite Population: $\sigma_{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \left(\frac{\sigma}{\sqrt{n}} \right)$.
 - Infinite Population: $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$.
- The finite population correction factor:** $\sqrt{(N-n)/(N-1)}$.

4. In many practical sampling situations, we find that the population involved, although finite, is "large," whereas the sample size is relatively "small." In such cases the finite population correction factor $\sqrt{(N-n)/(N-1)}$ is close to 1. As a result, the difference between the values of the standard deviation of \bar{x} for the finite and infinite population cases becomes negligible. Then, $\sigma_{\bar{x}} = \sigma/\sqrt{n}$ becomes a good approximation to the standard deviation of \bar{x} even though the population is finite.
5. *General guideline*, or rule of thumb, for computing the standard deviation of \bar{x} : Use $\sigma_{\bar{x}} = \sigma/\sqrt{n}$ whenever
- The population is infinite; or
 - The population is finite and the sample size is less than or equal to 5% of the population size; that is, $n/N \leq 0.05$.
6. In cases where $n/N > 0.05$, the finite population version of formula should be used in the computation of $\sigma_{\bar{x}}$. Unless otherwise noted, throughout the text we will assume that the population size is "large," $n/N < 0.05$, and expression σ/\sqrt{n} can be used to compute $\sigma_{\bar{x}}$.
7. We refer to the standard deviation of \bar{x} , $\sigma_{\bar{x}}$, as the standard error of the mean.
8. In general, the term standard error refers to the standard deviation of a point estimator.

 **Question** (p335)

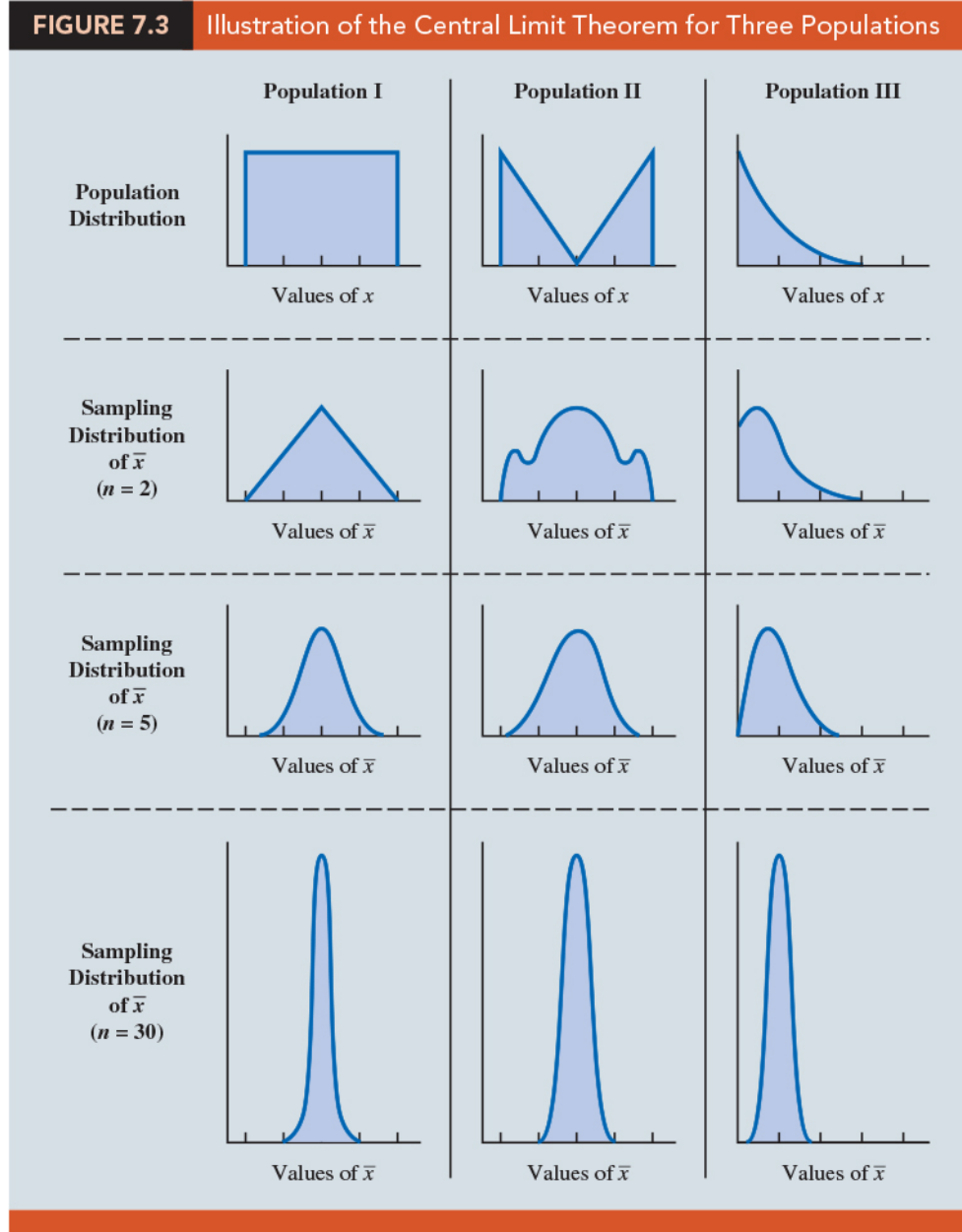
(EAI managers Problem) In Section 7.1 we saw that the standard deviation of annual salary for the population of 2500 EAI managers is $\sigma = 4000$. Compute the standard error.

sol: The population is finite, with $N = 2500$. However, with a sample size of 30, we have $n/N = 30/2500 = 0.012$. Because the sample size is less than 5% of the population size, we can ignore the finite population correction factor and use equation (7.3) to compute the standard error:

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{4000}{\sqrt{30}} = 730.3$$

Form of the Sampling Distribution of \bar{X}

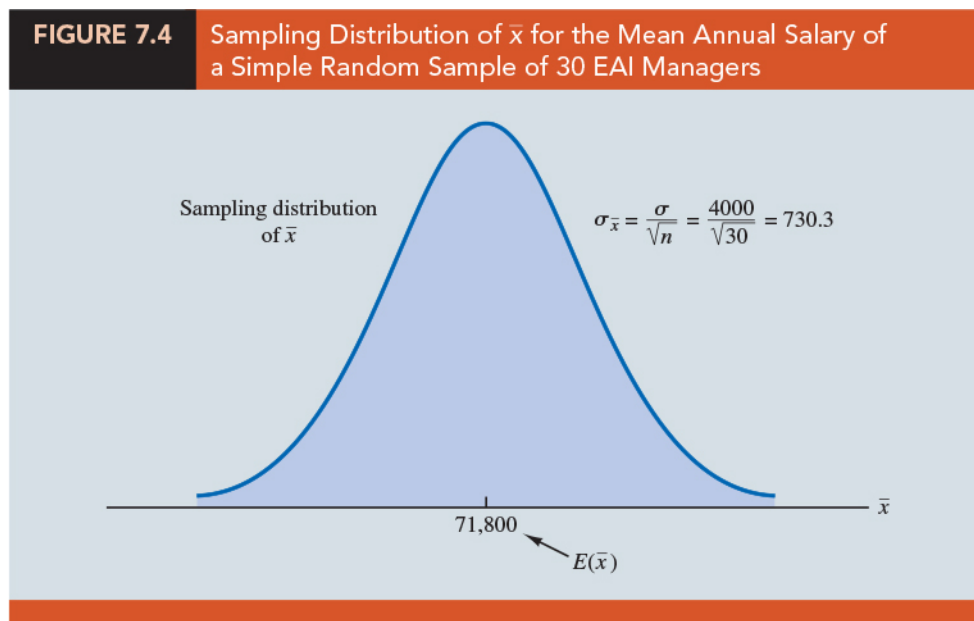
1. Consider two cases to determine the form or shape of \bar{x} : (1) The population has a normal distribution; and (2) the population does not have a normal distribution.
2. *Population has a Normal Distribution:* When the population has a normal distribution, the sampling distribution of \bar{x} is normally distributed for any sample size.
3. *Population does not have a Normal Distribution:* When the population from which we are selecting a random sample does not have a normal distribution, the central limit theorem is helpful in identifying the shape of the sampling distribution of \bar{x} .
4. **Central Limit Theorem (CLT):** In selecting random samples of size n from a population, the sampling distribution of the sample mean \bar{x} can be approximated by a normal distribution as the sample size becomes large.
5. (Figure 7.3) Illustration of the Central Limit Theorem for Three Populations
 - (a) Population I follows a uniform distribution. Population II is called the rabbit-eared distribution. It is symmetric, but the more likely values fall in the tails of the distribution. Population III is shaped like the exponential distribution; it is skewed to the right.
 - (b) When the sample size is $n = 2$, the shape of each sampling distribution is different from the shape of the corresponding population distribution.
 - (c) For samples of size $n = 5$, the shapes of the sampling distributions for populations I and II begin to look similar to the shape of a normal distribution. Even though the shape of the sampling distribution for population III begins to look similar to the shape of a normal distribution, some skewness to the right is still present.
 - (d) For samples of size 30, the shapes of each of the three sampling distributions are approximately normal.



6. General statistical practice is to assume that, for most applications, the sampling distribution of \bar{x} can be approximated by a normal distribution whenever the sample size is size 30 or more.
7. In cases where the population is highly skewed or outliers are present, samples of size 50 may be needed. Finally, if the population is discrete, the sample size needed for a normal approximation often depends on the population proportion.

Sampling Distribution of \bar{x} for the EAI Problem

1. The expected value and the standard deviation of the annual salary of EAI managers are $E(x) = \$71,800$ and $\sigma_{\bar{x}} = 730.3$. We do not have any information about the population distribution
2. If the population has a normal distribution, the sampling distribution of \bar{x} is normally distributed. If the population does not have a normal distribution, the simple random sample of 30 managers and the central limit theorem enable us to conclude that the sampling distribution of \bar{x} can be approximated by a normal distribution.
3. (Figure 7.4) In either case, we are comfortable proceeding with the conclusion that the sampling distribution of \bar{x} can be described by the normal distribution.




Practical Value of the Sampling Distribution of \bar{x}

1. Whenever a simple random sample is selected, we cannot expect the sample mean to exactly equal the population mean.
2. The practical reason we are interested in the sampling distribution of \bar{x} is that it can be used to provide probability information about the difference between the sample mean and the population mean.

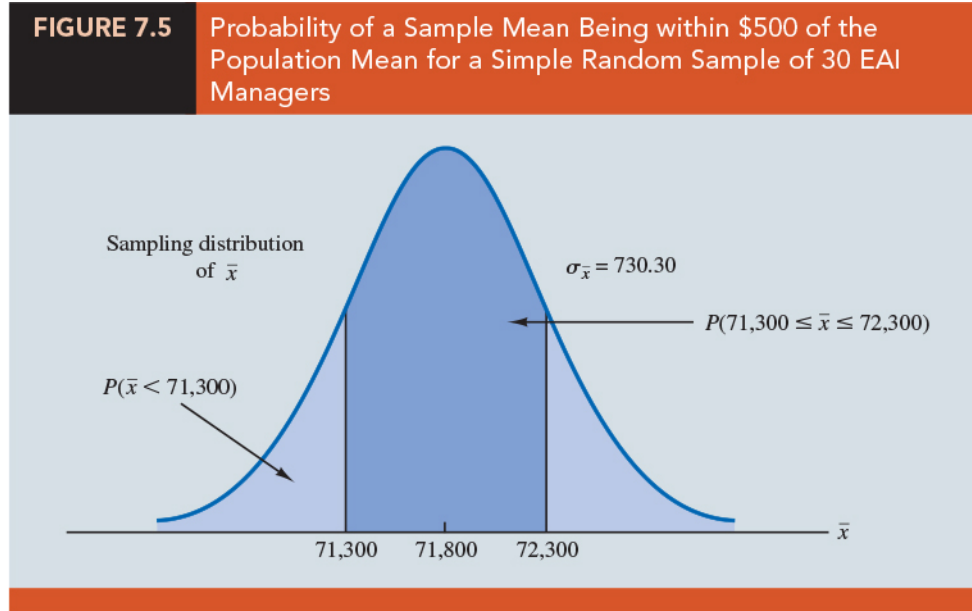
3. Example EAI problem:

- (a) Suppose the personnel director believes the sample mean will be an acceptable estimate of the population mean if the sample mean is within \$500 of the population mean. However, it is not possible to guarantee that the sample mean will be within \$500 of the population mean.
- (b) Indeed, Table 7.5 and Figure 7.1 show that some of the 500 sample means differed by more than \$2000 from the population mean. So we must think of the personnel director's request in probability terms.
- (c) That is, the personnel director is concerned with the following question: What is the probability that the sample mean computed using a simple random sample of 30 EAI managers will be within \$500 of the population mean?

(d)  **Question** (p338)

Consider the EAI problem, the mean and the standard deviation of annual salary for the population of $N = 2500$ EAI managers is $\mu = \$71,800$ and $\sigma = 4000$. With a sample size of $n = 30$, the personnel director wants to know the probability that \bar{x} is between \$71,300 and \$72,300.

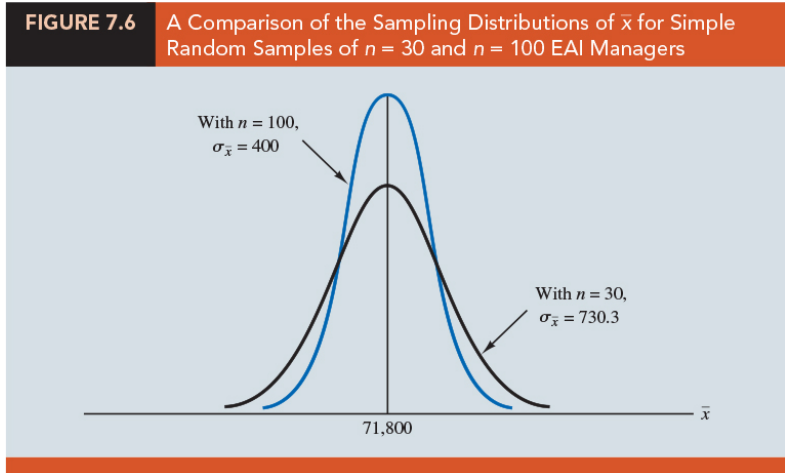
sol: (Figure 7.5) Because the sampling distribution is normally distributed, with mean \$71,800 and standard error of the mean 730.3, we can use the standard normal probability table to find the area or probability.



- (e) The preceding computations show that a simple random sample of 30 EAI managers has a 0.5034 probability of providing a sample mean \bar{x} that is within \$500 of the population mean.
- (f) Thus, there is a $1 - 0.5034 = 0.4966$ probability that the difference between \bar{x} and $\mu = \$71,800$ will be more than \$500. In other words, a simple random sample of 30 EAI managers has roughly a 50–50 chance of providing a sample mean within the allowable \$500. Perhaps a larger sample size should be considered.

Relationship Between the Sample Size and the Sampling Distribution of \bar{X}

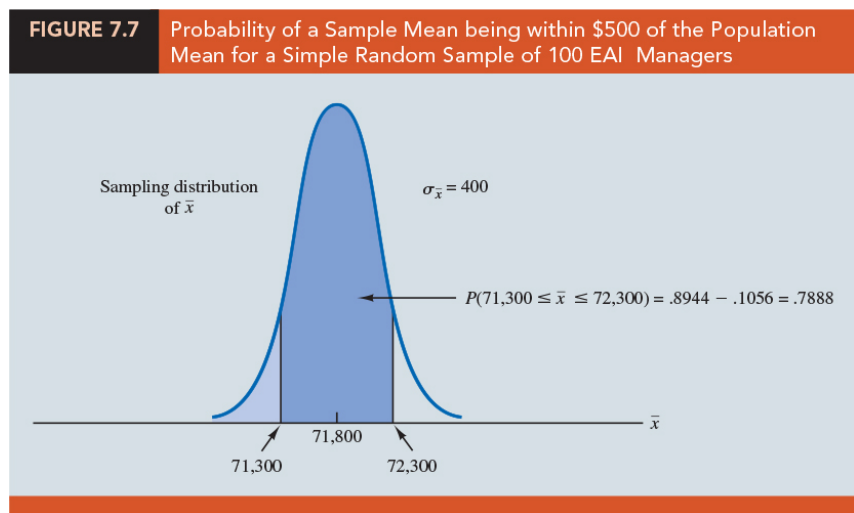
- (Figure 7.6) The sampling distributions of \bar{x} with $n = 30$ and $n = 100$. Because the sampling distribution with $n = 100$ has a smaller standard error, the values of \bar{x} have less variation and tend to be closer to the population mean than the values of \bar{x} with $n = 30$.



2. Question (p339)

Consider the EAI problem, With a sample size of $n = 100$ EAI managers, compute the probability that \bar{x} is within \$500 of the population mean $\mu = \$71,800$.

sol:



- Thus, by increasing the sample size from 30 to 100 EAI managers, we increase the probability of obtaining a sample mean within \$500 of the population mean from 0.5034 to 0.7888. The important point in this discussion is that as the sample size is increased, the standard error of the mean decreases.
- As a result, the larger sample size provides a higher probability that the sample mean is within a specified distance of the population mean.

7.6 Sampling Distribution of \bar{p}

- The sample proportion \bar{p} is the point estimator of the population proportion p :

$$\bar{p} = \frac{x}{n}$$

where x is the number of elements in the sample that possess the characteristic of interest and n is the sample size.

- The sample proportion \bar{p} is a random variable and its probability distribution is called the sampling distribution of \bar{p} .
- Sampling Distribution of \bar{p}** The sampling distribution of \bar{p} is the probability distribution of all possible values of the sample proportion \bar{p} .

Expected Value of \bar{p}

- The expected value of \bar{p} , the mean of all possible values of \bar{p} , is equal to the population proportion p :

$$E(\bar{p}) = p$$

- Because $E(\bar{p}) = p$, \bar{p} is an unbiased estimator of p .
- Example** (Section 7.1) we noted that $p = 0.60$ for the EAI population, where p is the proportion of the population of managers who participated in the company's management training program. Thus, the expected value of \bar{p} for the EAI sampling problem is 0.60.

Standard Deviation of \bar{p}

1. The standard deviation of \bar{p} depends on whether the population is finite or infinite:

(a) Finite Population:
$$\sigma_{\bar{p}} = \sqrt{\frac{N-n}{N-1}} \sqrt{\frac{p(1-p)}{n}} .$$

(b) Infinite Population:
$$\sigma_{\bar{p}} = \sqrt{\frac{p(1-p)}{n}} .$$

2. The difference between the expressions for the finite population and the infinite population becomes negligible if the size of the finite population is large in comparison to the sample size.

3. Rule of thumb:

(a) if the population is finite with $n/N \leq 0.05$, use $\sigma_{\bar{p}} = \sqrt{p(1-p)/n}$.

(b) if the population is finite with $n/N > 0.05$, the finite population correction factor should be used.

(c) Unless specifically noted, throughout the text we will assume that the population size is large in relation to the sample size and thus the finite population correction factor is unnecessary.

4. We stated that in general the term standard error refers to the standard deviation of a point estimator. Thus, for proportions we use standard error of the proportion to refer to the standard deviation of \bar{p} .

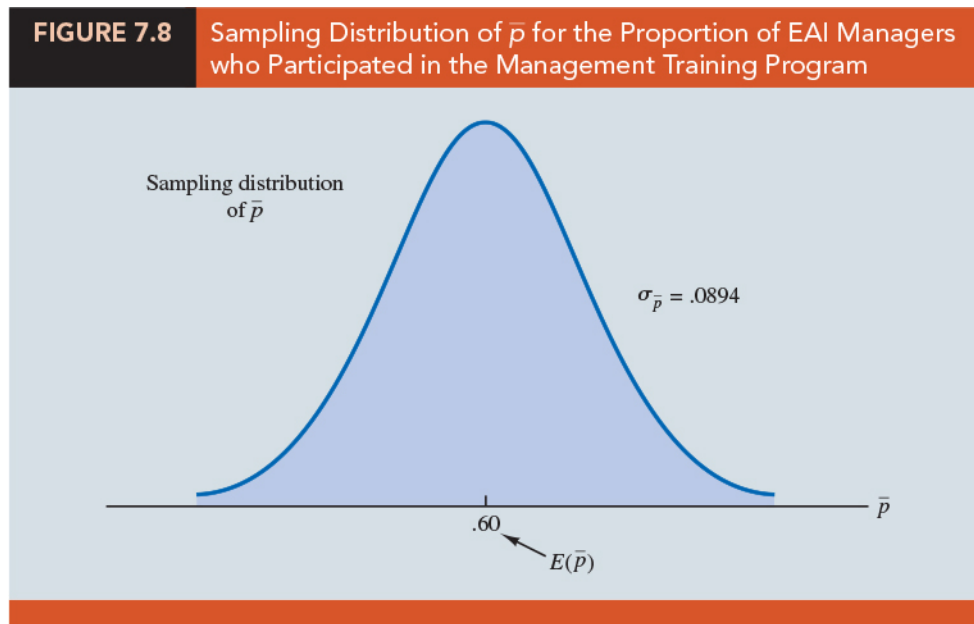
5. **Example** (The EAI example) With $n/N = 30/2500 = 0.012$, we can ignore the finite population correction factor when we compute the standard error of the proportion. For the simple random sample of 30 managers, the standard error of the proportion associated with simple random samples of 30 EAI managers:

$$\sigma_{\bar{p}} = \sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{0.6(1-0.6)}{30}} = 0.0894 .$$

Form of the Sampling Distribution of \bar{p}


1. For a simple random sample from a large population, the value of x is a binomial random variable indicating the number of elements in the sample with the characteristic of interest.

2. Because n is a constant, the probability of x/n is the same as the binomial probability of x , which means that the sampling distribution of \bar{p} is also a discrete probability distribution and that the probability for each value of x/n is the same as the probability of x .
3. (Chapter 6) we showed that a binomial distribution can be approximated by a normal distribution whenever the sample size is large enough to satisfy two conditions:
(1) $np \geq 5$ and $n(1-p) \geq 5$.
4. Assuming these two conditions are satisfied, the probability distribution of x in the sample proportion, $\bar{p} = x/n$, can be approximated by a normal distribution. And because n is a constant, the sampling distribution of \bar{p} can also be approximated by a normal distribution.
5. The sampling distribution of \bar{p} can be approximated by a normal distribution whenever $np \geq 5$ and $n(1-p) \geq 5$.
6. In practical applications, when an estimate of a population proportion is desired, we find that sample sizes are almost always large enough to permit the use of a normal approximation for the sampling distribution of \bar{p} .
7. **Example** (Figure 7.8) (The EAI example) Recall that we know that the population proportion of managers who participated in the training program is $p = 0.60$. With a simple random sample of size 30, we have $np = 30(0.60) = 18$ and $n(1-p) = 30(.40) = 12$. Thus, the sampling distribution of \bar{p} can be approximated by a normal distribution.



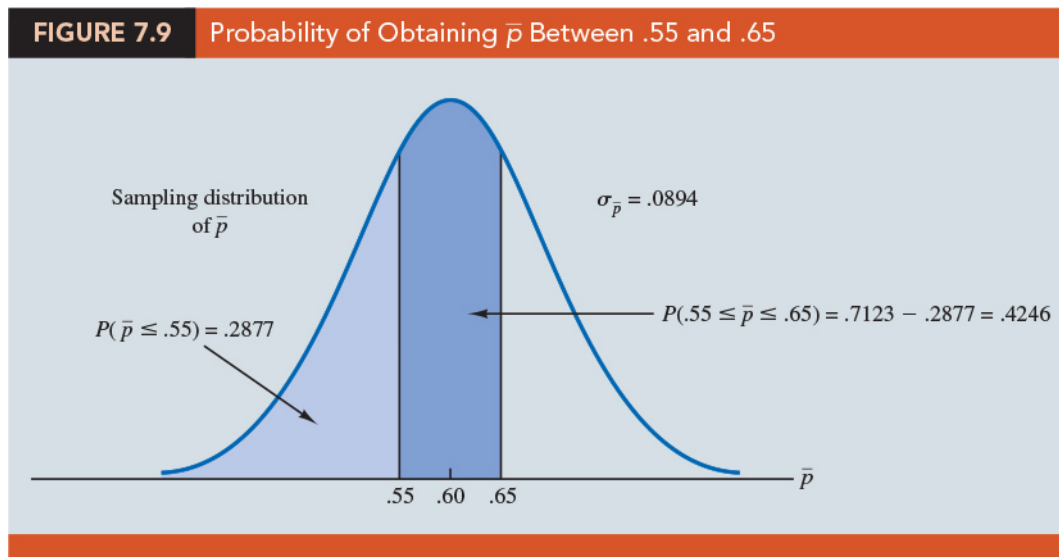
Practical Value of the Sampling Distribution of \bar{p}

1. The practical value of the sampling distribution of \bar{p} is that it can be used to provide probability information about the difference between the sample proportion and the population proportion.

 **Question** (p346)

Consider the EAI problem, suppose that the personnel director wants to know the probability of obtaining a value of \bar{p} that is within 0.05 of the population proportion of EAI managers who participated in the training program. That is, what is the probability of obtaining a sample with a sample proportion \bar{p} between 0.55 and 0.65? Consider two cases: $n = 30$ and $n = 100$ using a normal approximation.

sol:

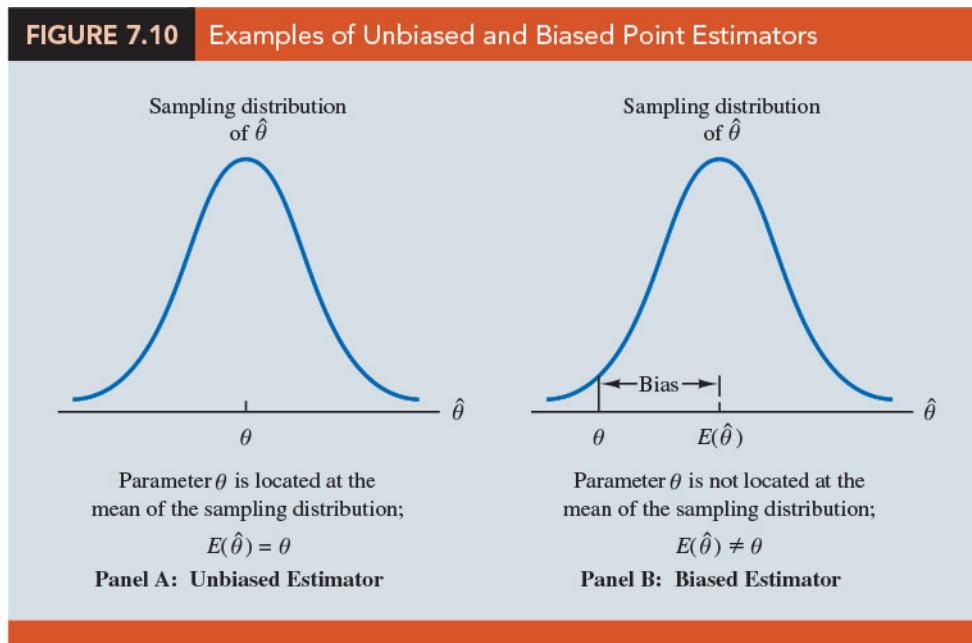


7.7 Properties of Point Estimators

1. Before using a sample statistic as a point estimator, statisticians check to see whether the sample statistic demonstrates certain properties associated with good point estimators.
2. We discuss three properties of good point estimators: unbiased, efficiency, and consistency.
3. General notation:
 - (a) θ : the population parameter of interest,
 - (b) $\hat{\theta}$: the sample statistic or point estimator of θ .

Unbiased

1. If the expected value of the sample statistic is equal to the population parameter being estimated, the sample statistic is said to be an unbiased estimator of the population parameter.
2. **Unbiased estimator:** The sample statistic $\hat{\theta}$ is an unbiased estimator of the population parameter θ if
$$E(\hat{\theta}) = \theta,$$
where $E(\hat{\theta})$ is the expected value of the sample statistic $\hat{\theta}$.
3. (Figure 7.10) the cases of unbiased and biased point estimators.
 - (a) The unbiased estimator: the mean of the sampling distribution is equal to the value of the population parameter.
 - (b) A biased estimator: the mean of the sampling distribution is less than or greater than the value of the population parameter. The amount of the bias is shown.

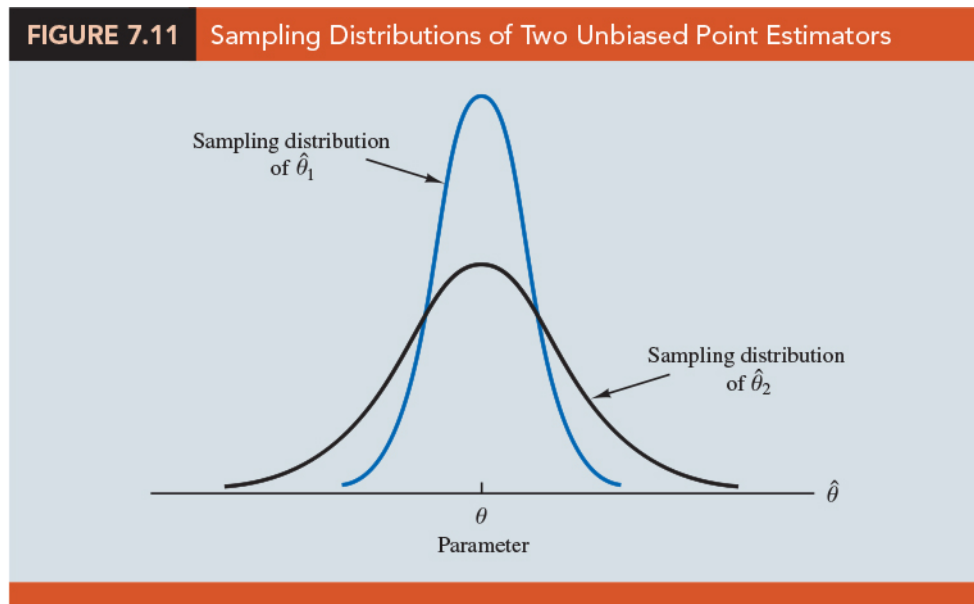


- We stated that $E(\bar{x}) = \mu$ and $E(\bar{p}) = p$. Both \bar{x} and \bar{p} are unbiased estimators of their corresponding population parameters μ and p .
- It can be shown that $E(s^2) = \sigma^2$, the sample variance s^2 is an unbiased estimator of the population variance σ^2 .
- The reason for using $n-1$ rather than n in sample variance (and the sample standard deviation) is to make the sample variance an **unbiased** estimator of the population variance.

Efficiency

- Assume that a simple random sample of n elements can be used to provide **two unbiased point estimators** of the same population parameter. In this situation, we would prefer to use the point estimator with the **smaller standard error**, because it tends to provide estimates closer to the population parameter.
- The point estimator with the smaller standard error is said to have greater **relative efficiency** than the other.
- (Figure 7.11) the sampling distributions of two unbiased point estimators, $\hat{\theta}_1$ and $\hat{\theta}_2$.

- (a) The standard error of $\hat{\theta}_1$ is less than the standard error of $\hat{\theta}_2$;
- (b) Values of $\hat{\theta}_1$ have a greater chance of being close to the parameter θ than do values of $\hat{\theta}_2$.
- (c) $\hat{\theta}_1$ is relatively more efficient than $\hat{\theta}_2$ and is the preferred point estimator.

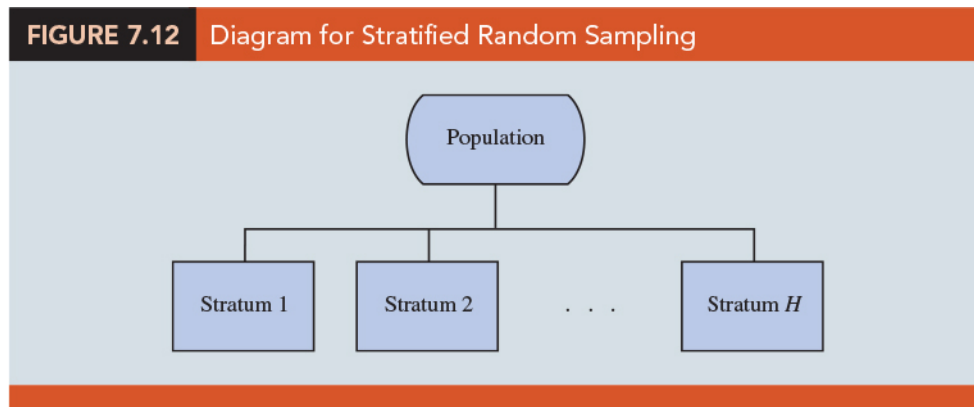


Consistency

1. A point estimator is consistent if the values of the point estimator tend to become closer to the population parameter as the sample size becomes larger.
2. A large sample size tends to provide a better point estimate than a small sample size.
3. The standard error of \bar{x} , $\sigma_{\bar{x}} = \sigma/\sqrt{n}$, is related to the sample size such that larger sample sizes provide smaller values for $\sigma_{\bar{x}}$, we conclude the sample mean \bar{x} is a consistent estimator of the population mean μ .
4. The sample proportion \bar{p} is a consistent estimator of the population proportion p .

7.8 Other Sampling Methods¹

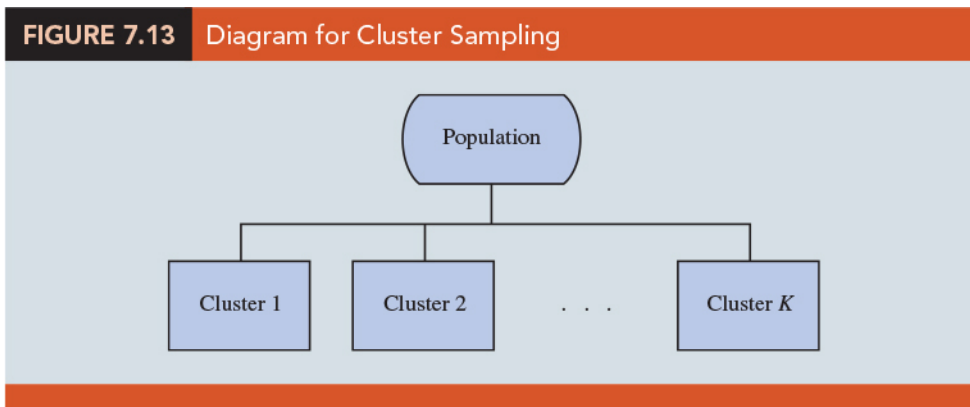
Stratified Random Sampling



1. (Figure 7.12) In stratified random sampling, the elements in the population are first divided into groups called strata, such that each element in the population belongs to one and only one stratum. A simple random sample is taken from each stratum.
2. The basis for forming the strata, such as department, location, age, industry type, and so on, is at the discretion of the designer of the sample.
3. However, the best results are obtained when the elements within each stratum are as much alike as possible.
4. Formulas are available for combining the results for the individual stratum samples into one estimate of the population parameter of interest.

¹A more in-depth treatment is provided in Chapter 22.

Cluster Sampling



1. (Figure 7.13) In cluster sampling, the elements in the population are first divided into separate groups called clusters.
2. Each element of the population belongs to one and only one cluster. A simple random sample of the clusters is then taken. All elements within each sampled cluster form the sample.
3. Cluster sampling tends to provide the best results when the elements within the clusters are not alike.
4. In the ideal case, each cluster is a representative small-scale version of the entire population. The value of cluster sampling depends on how representative each cluster is of the entire population.
5. If all clusters are alike in this regard, sampling a small number of clusters will provide good estimates of the population parameters.
6. One of the primary applications of cluster sampling is area sampling, where clusters are city blocks or other well-defined areas.

Systematic Sampling

1. Systematic sampling is a type of probability sampling method in which sample members from a larger population are selected according to a random starting point but with a fixed, periodic interval.

2. **Example** If a sample size of 50 is desired from a population containing 5000 elements, we will sample one element for every $5000/50 = 100$ elements in the population.
3. A systematic sample for this case involves selecting randomly one of the first 100 elements from the population list.
4. Other sample elements are identified by starting with the first sampled element and then selecting every 100th element that follows in the population list.
5. In effect, the sample of 50 is identified by moving systematically through the population and identifying every 100th element after the first randomly selected element.
6. Because the first element selected is a random choice, a systematic sample is usually assumed to have the properties of a simple random sample.
7. This assumption is especially applicable when the list of elements in the population is a random ordering of the elements.

Convenience Sampling

1. **Probability sampling techniques:** elements selected from the population have a known probability of being included in the sample.
2. Convenience sampling is a nonprobability sampling technique. The sample is identified primarily by convenience. Elements are included in the sample without prespecified or known probabilities of being selected.
3. **Example** A professor conducting research at a university may use student volunteers to constitute a sample simply because they are readily available and will participate as subjects for little or no cost.
4. Convenience samples have the advantage of relatively easy sample selection and data collection; however, it is impossible to evaluate the "goodness" of the sample in terms of its representativeness of the population.

5. A convenience sample may provide good results or it may not; no statistically justified procedure allows a probability analysis and inference about the quality of the sample results.

Judgment Sampling

1. Judgment sampling is a nonprobability sampling technique.
2. In this approach, the person most knowledgeable on the subject of the study selects elements of the population that he or she feels are most representative of the population.
3. The quality of the sample results depends on the judgment of the person selecting the sample.

7.9 Big Data and Standard Errors of Sampling Distributions

1. The purpose of statistical inference is to use sample data to quickly and inexpensively gain insight into some characteristic of a population.
2. There are two general reasons a sample may fail to be representative of the population of interest: sampling error and nonsampling error.

Sampling Error

1. One reason a sample may fail to represent the population from which it has been taken is sampling error, or deviation of the sample from the population that results from random sampling.
2. When collecting a sample randomly, the data in the sample cannot be expected to be perfectly representative of the population from which it has been taken.

3. Sampling error is unavoidable when collecting a random sample; this is a risk we must accept when we chose to collect a random sample rather than incur the costs associated with taking a census of the population.
4. The standard errors of the sampling distributions of the sample mean \bar{x} and the sample proportion of \bar{p} reflect the potential for sampling error when using sample data to estimate the population mean μ and the population proportion p , respectively.
5. For an extremely large sample there may be little potential for sampling error.

Nonsampling Error

1. We can not conclude that an extremely large sample will always provide reliable information about the population of interest.
2. Deviations of the sample from the population that occur for reasons other than random sampling are referred to as nonsampling error.
3. Nonsampling error can occur for a variety of reasons.
 - (a) **Coverage error:** Nonsampling error that results when the research objective and the population from which the sample is to be drawn are not aligned.
 - (b) **Nonresponse error:** Nonsampling error that results when some segments of the population are either more or less likely to respond to the survey mechanism.
 - (c) **Measurement error:** Nonsampling error that results from the incorrect measurement of the population characteristic of interest.
4. Understanding these limitations of sampling will enable us to be more realistic and pragmatic when interpreting sample data and using sample data to draw conclusions about the target population.

Understanding What Big Data Is

1. The processes that generate big data can be described by four attributes or dimensions that are referred to as the four V's:

- (a) **Volume**: the amount of data generated
 - (b) **Variety**: the diversity in types and structures of data generated
 - (c) **Veracity**: the reliability of the data generated
 - (d) **Velocity**: the speed at which the data are generated
2. **Tall data**: A data set that has so many observations that traditional statistical inference has little meaning.
 3. **Wide data**: A data set that has so many variables that simultaneous consideration of all variables is infeasible.
 4. Statistics are useful tools for understanding the information embedded in a big data set, but we must be careful when using statistics to analyze big data.
 5. It is important that we understand the limitations of statistics when applied to big data and we temper our interpretations accordingly.

Implications of Big Data for Sampling Error

1. (Table 7.7) the standard error of the sampling distribution of the sample mean time spent by individual customers when they visit PDT's website decreases as the sample size increases.
2. (Table 7.8) the standard error of the sampling distribution of the proportion of the sample that clicked on any of the ads featured on PenningtonDailyTimes.com decreases as the sample size increases.

Sample Size n	Standard Error $s_x = s/\sqrt{n}$
10	6.32456
100	2.00000
1,000	.63246
10,000	.20000
100,000	.06325
1,000,000	.02000
10,000,000	.00632
100,000,000	.00200
1,000,000,000	.00063

Sample Size n	Standard Error $\sigma_{\bar{p}} = \sqrt{\bar{p}(1 - \bar{p})/n}$
10	.15808
100	.04999
1,000	.01581
10,000	.00500
100,000	.00158
1,000,000	.00050
10,000,000	.00016
100,000,000	.00005
1,000,000,000	.00002

☺ EXERCISES

7.2 : 2, 10

7.3 : 15, 17

7.5 : 18, 24, 25, 30

7.6 : 34, 35, 39

7.9 : 42, 44, 47, 49

SUP : 52, 57, 62, 65

“我哭泣是因為我沒有鞋子，直到我遇到了一個人 — 他沒有腳”

“I cried because I had no shoes until I met a man who had no feet”

— *Helen Keller (June 27, 1880 – June 1, 1968)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 8: Interval Estimation

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

系級: _____ 學號: _____ 姓名: _____

Overview

1. A point estimator is a sample statistic used to estimate a population parameter.
 - (a) The sample mean \bar{x} is a point estimator of the population mean μ .
 - (b) The sample proportion \bar{p} is a point estimator of the population proportion p .
2. A point estimator cannot be expected to provide the exact value of the population parameter, an interval estimate is often computed by adding and subtracting a value, called the margin of error, to the point estimate.
3. The general form of an interval estimate:
Point estimate \pm Margin of error
4. The purpose of an interval estimate is to provide information about how close the point estimate, provided by the sample, is to the value of the population parameter.

5. This chapter shows how to compute interval estimates of a population mean μ and a population proportion p .
6. The sampling distributions of \bar{x} and \bar{p} play key roles in computing these interval estimates.

8.1 Population Mean: σ Known

1. To develop an interval estimate of a population mean, either the population standard deviation σ or the sample standard deviation s must be used to compute the margin of error.
2. The cases that the σ is known:
 - (a) Large amounts of relevant historical data are available and can be used to estimate the population standard deviation prior to sampling.
 - (b) In quality control applications where a process is assumed to be operating correctly, or "in control," it is appropriate to treat the population standard deviation as known.
3. Example Lloyd's example
 - (a) Each week Lloyd's Department Store selects a simple random sample of 100 customers in order to learn about the amount spent per shopping trip.
 - (b) With x representing the amount spent per shopping trip, the sample mean \bar{x} provides a point estimate of μ , the mean amount spent per shopping trip for the population of all Lloyd's customers.
 - (c) Based on the historical data, Lloyd's now assumes a known value of $\sigma = \$20$ for the population standard deviation and the population follows a normal distribution.

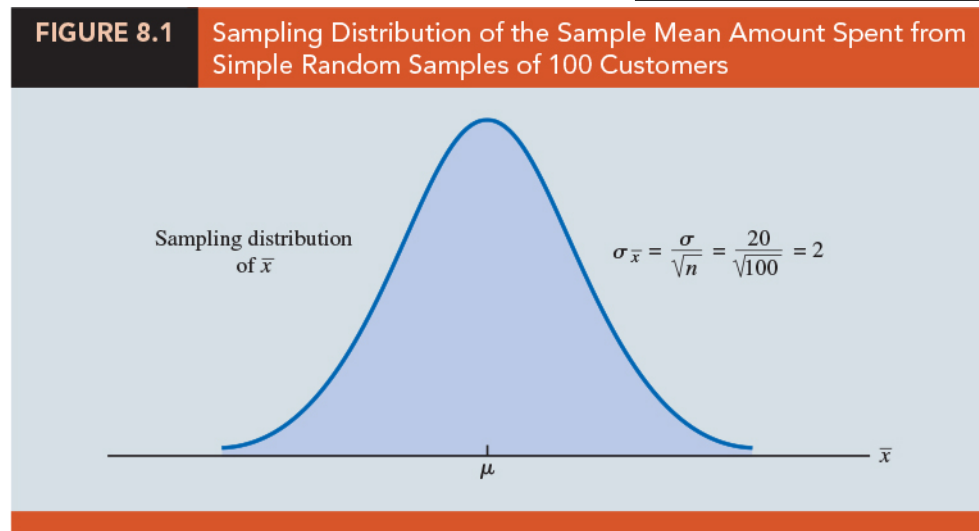
- (d) During the most recent week, Lloyd's surveyed 100 customers ($n = 100$) and obtained a sample mean of $\bar{x} = \$82$.
4. How to compute the margin of error for the sample mean and develop an interval estimate of the population mean.

Margin of Error and the Interval Estimate

1. (Chapter 7) The sampling distribution of \bar{x} can be used to compute the probability that \bar{x} will be within a given distance of μ .

2. **Example** Lloyd's example

- (a) (Figure 8.1) The population of amounts spent is normally distributed with a standard deviation of $\sigma = 20$. The sampling distribution of \bar{x} follows a normal distribution with a standard error of $\sigma_{\bar{x}} = \sigma/\sqrt{n} = 20/\sqrt{100} = 2$.



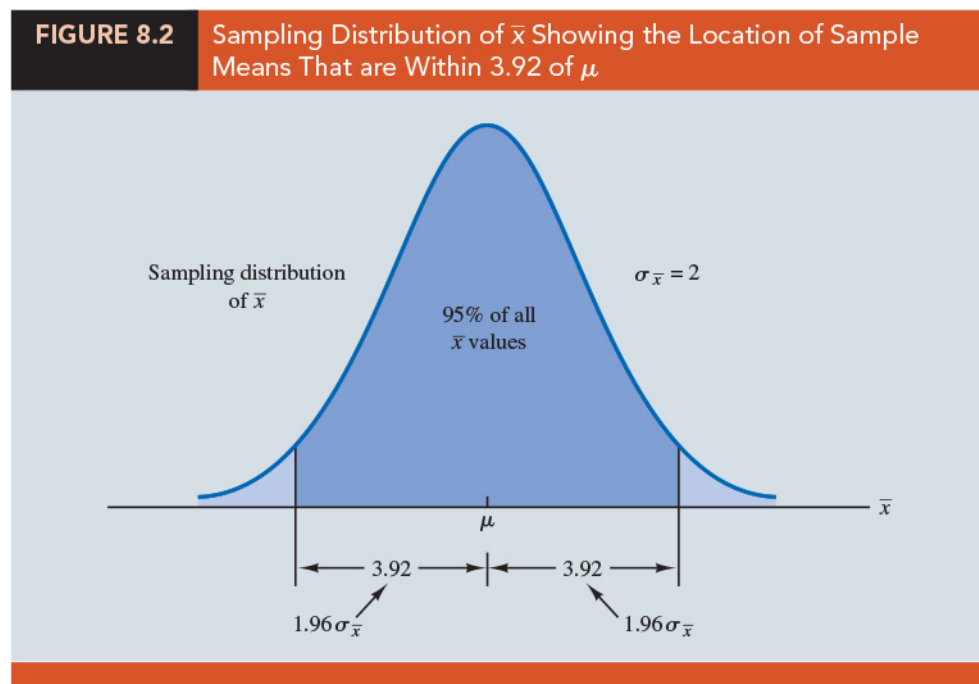
- (b) The sampling distribution shows how values of \bar{x} are distributed around the population mean μ ,
- (c) The sampling distribution of \bar{x} provides information about the possible differences between \bar{x} and μ .
3. (The standard normal probability table) The 95% of the values of any normally distributed random variable are within ± 1.96 standard deviations of the mean.

4. When the sampling distribution of \bar{x} is normally distributed, 95% of the \bar{x} values must be within $\pm 1.96\sigma_{\bar{x}}$ of the mean μ .

5. **Example** Lloyd's example:

(a) (Figure 8.2) The sampling distribution of \bar{x} is normally distributed with a standard error of $\sigma_{\bar{x}} = 2$.

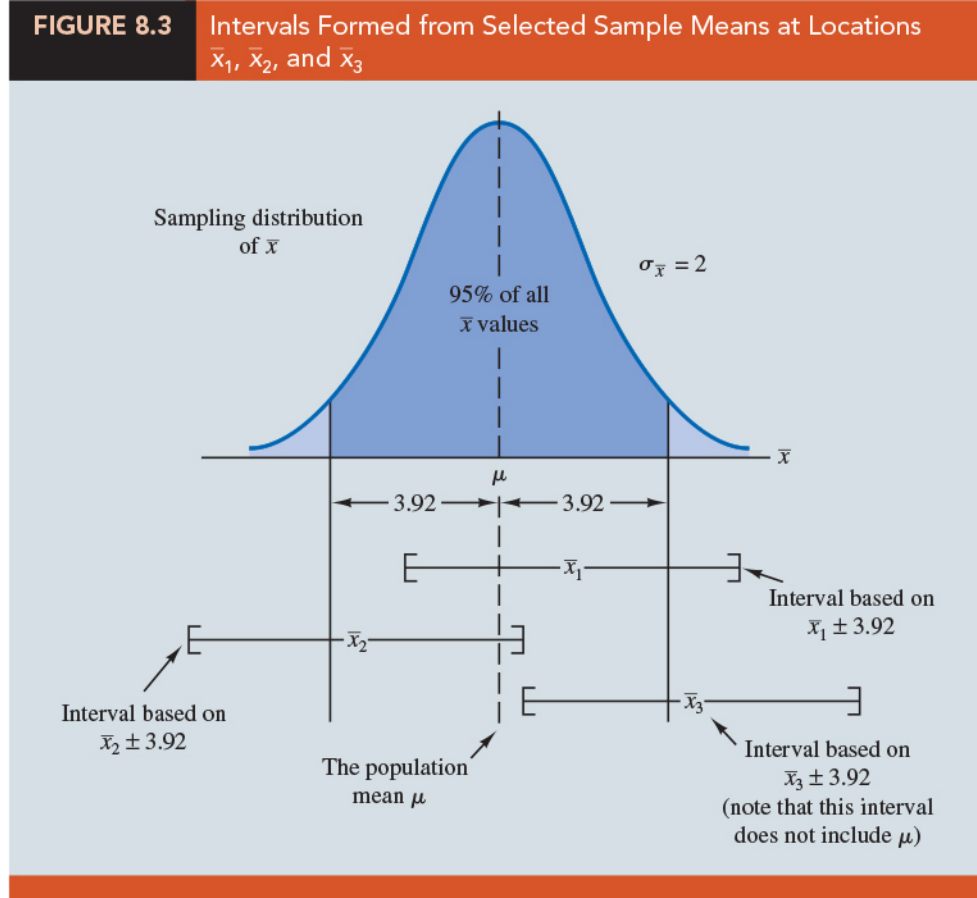
(b) Because $\pm 1.9\sigma_{\bar{x}} = \pm 1.96(2) = \pm 3.92$, we can conclude that 95% of all \bar{x} values obtained using a sample size of $n = 100$ will be within ± 3.92 of the population mean μ .



(c) Suppose we set the margin of error equal to 3.92 and compute the interval estimate of μ using $\bar{x} \pm 3.92$.

6. *An interpretation for this interval estimate:*

(a) (Figure 8.3) Consider the values of \bar{x} that could be obtained if we took three different simple random samples, each consisting of 100 Lloyd's customers.



- (b) The first sample mean might turn out to have the value shown as \bar{x}_1 in Figure 8.3. The interval formed by subtracting 3.92 from \bar{x}_1 and adding 3.92 to \bar{x}_1 includes the population mean μ .
- (c) The interval $\bar{x}_2 \pm 3.92$ obtained by the second sample mean \bar{x}_2 also includes the population mean μ .
- (d) The interval $\bar{x}_3 \pm 3.92$ formed by the third sample mean does not include the population mean μ .
- (e) Any sample mean \bar{x} that is within the darkly shaded region of Figure 8.3 will provide an interval that contains the population mean μ .
- (f) Because 95% of all possible sample means are in the darkly shaded region, 95% of all intervals formed by subtracting 3.92 from \bar{x} and adding 3.92 to \bar{x} will include the population mean μ .
- (g) The quality assurance team at Lloyd's surveyed 100 customers and obtained

a sample mean amount spent of $\bar{x} = 82$. Using $\bar{x} \pm 3.92$ to construct the interval estimate, we obtain $82 \pm 3.92 = (78.08, 85.92)$.

- (h) Because 95% of all the intervals constructed using $\bar{x} \pm 3.92$ will contain the population mean, we say that we are 95% confident that the interval 78.08 to 85.92 includes the population mean μ .
7. We say that this interval has been established at the 95% confidence level. The value 0.95 is referred to as the confidence coefficient, and the interval 78.08 to 85.92 is called the 95% confidence interval.

8. Interval Estimate of a Population Mean: σ Known:

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \quad (8.1)$$

where $(1-\alpha)$ is the confidence coefficient and $z_{\alpha/2}$ is the z value providing an area of $\alpha/2$ in the upper tail of the standard normal probability distribution.

9. Comparing the results for the 90%, 95%, and 99% confidence levels, we see that in order to have a higher degree of confidence, the margin of error and thus the width of the confidence interval must be larger.

 **Question** (p113)

Construct the 90%, 95% and 99% confidence intervals for the Lloyd's example.

sol:

Confidence Level	α	$\alpha/2$	$z_{\alpha/2}$
90%	.10	.05	1.645
95%	.05	.025	1.960
99%	.01	.005	2.576

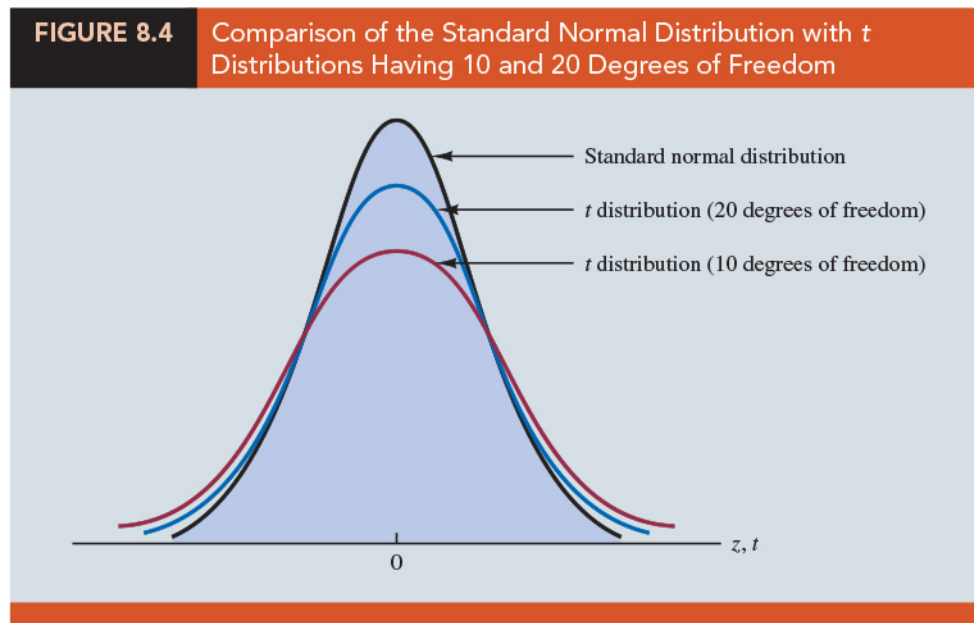
Practical Advice

1. If the population follows a normal distribution, the confidence interval provided by expression (8.1) is exact.
2. If expression (8.1) were used repeatedly to generate 95% confidence intervals, exactly 95% of the intervals generated would contain the population mean.
3. If the population does not follow a normal distribution, the confidence interval provided by expression (8.1) will be approximate. In this case, the quality of the approximation depends on both the distribution of the population and the sample size.
4. In most applications, a sample size of $n \geq 30$ is adequate when using expression (8.1) to develop an interval estimate of a population mean.
5. If the population is not normally distributed but is roughly symmetric, sample sizes as small as 15 can be expected to provide good approximate confidence intervals.
6. With smaller sample sizes, expression (8.1) should only be used if the analyst believes, or is willing to assume, that the population distribution is at least approximately normal.

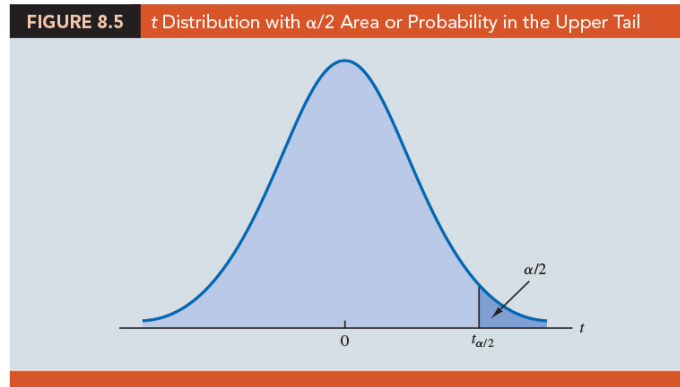
8.2 Population Mean: σ Unknown

1. When the population standard deviation σ is unknown, we must use the same sample to estimate both μ and σ to obtain an interval estimate of a population mean μ .
2. When s is used to estimate σ , the margin of error and the interval estimate for the population mean are based on a probability distribution known as the t distribution.

3. Although the mathematical development of the t distribution is based on the assumption of a normal distribution for the population we are sampling from, research shows that the t distribution can be successfully applied in many situations where the population deviates significantly from normal.
4. The t distribution is a family of similar probability distributions, with a specific t distribution depending on a parameter known as the degrees of freedom (df).
5. (Figure 8.4) The t distribution with one degree of freedom is unique, as is the t distribution with two degrees of freedom, with three degrees of freedom, and so on. As the number of degrees of freedom increases, the difference between the t distribution and the standard normal distribution becomes smaller and smaller.
6. A t distribution with more degrees of freedom exhibits less variability and more closely resembles the standard normal distribution.

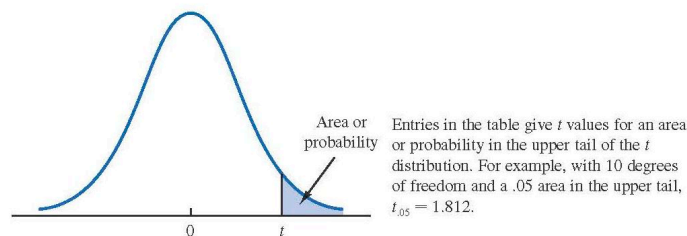


7. (Figure 8.5) The mean of the t distribution is zero. In general, we will use the notation $t_{\alpha/2}$ to represent a t value with an area of $\alpha/2$ in the upper tail of the t distribution.



8. (Table 2, Appendix B) Each row in the table corresponds to a separate t distribution with the degrees of freedom. (e.g., $t_{9,0.025} = 2.262$). As the degrees of freedom continue to increase, $t_{0.025}$ approaches $z_{0.025} = 1.96$. The standard normal distribution z values can be found in the infinite degrees of freedom (labeled ∞) of the t distribution table. For more than 100 degrees of freedom, the standard normal z value provides a good approximation to the t value.

TABLE 2 t Distribution



Degrees of Freedom	Area in Upper Tail					
	.20	.10	.05	.025	.01	.005
1	1.376	3.078	6.314	12.706	31.821	63.656
2	1.061	1.886	2.920	4.303	6.965	9.925
3	.978	1.638	2.353	3.182	4.541	5.841
4	.941	1.533	2.132	2.776	3.747	4.604
5	.920	1.476	2.015	2.571	3.365	4.032
6	.906	1.440	1.943	2.447	3.143	3.707
7	.896	1.415	1.895	2.365	2.998	3.499
8	.889	1.397	1.860	2.306	2.896	3.355
9	.883	1.383	1.833	2.262	2.821	3.250
10	.879	1.372	1.812	2.228	2.764	3.169
11	.876	1.363	1.796	2.201	2.718	3.106
12	.873	1.356	1.782	2.179	2.681	3.055
13	.870	1.350	1.771	2.160	2.650	3.012
14	.868	1.345	1.761	2.145	2.624	2.977

Margin of Error and the Interval Estimate

1. An interval estimate of a population mean for the σ known case: $\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

2. To compute an interval estimate of μ for the σ unknown case, the sample standard deviation s is used to estimate σ , and $z_{\alpha/2}$ is replaced by the t distribution value $t_{\alpha/2}$.
3. The margin of error is then given by $t_{\alpha/2} \frac{s}{\sqrt{n}}$.
4. **Interval Estimate of a Population Mean: σ Unknown:**

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} \quad (8.2)$$

where s is the sample standard deviation, $(1-\alpha)$ is the **confidence coefficient**, and $t_{\alpha/2}$ is the t value providing an area of $\alpha/2$ in the upper tail of the t distribution with $n-1$ degrees of freedom.

5. The reason the number of degrees of freedom associated with the t value in expression (8.2) is $n-1$ concerns the use of s as an estimate of the population standard deviation σ .

- (a) The expression for the sample standard deviation is

$$s = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

- (b) Degrees of freedom refer to the **number of independent** pieces of information that go into the computation of $\sum (x_i - \bar{x})^2$.

- (c) The n pieces of information involved in computing $\sum (x_i - \bar{x})^2$ are as follows: $x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x}$.

- (d) Note that $\sum (x_i - \bar{x}) = 0$ for any data set.

- (e) Thus, only $n-1$ of the $x_i - \bar{x}$ values are independent; that is, if we know $n-1$ of the values, the remaining value can be determined exactly by using the condition that $\sum (x_i - \bar{x})$ values must be 0 .

- (f) Thus, $n-1$ is the number of degrees of freedom associated with $\sum (x_i - \bar{x})^2$ and hence the number of degrees of freedom for the t distribution in expression (8.2).

 Question (p384)

(Table 8.3) To illustrate the interval estimation procedure for the σ unknown case, we will consider a study designed to estimate the mean credit card debt for the population of U.S. households. A sample of $n = 70$ households provided the credit card balances shown in Table 8.3. For this situation, no previous estimate of the population standard deviation σ is available. Thus, the sample data must be used to estimate both the population mean and the population standard deviation. Compute an interval estimate of the population mean credit card balance.

sol:

Compute: $\bar{x} = \$9312$, $s = 4007$.

With 95% confidence and $n-1 = 69$ degrees of freedom: $t_{69,0.025} = 1.995$.

An interval estimate of the population mean credit card balance:

$$9312 \pm 1.995 \frac{4007}{\sqrt{70}} = 9312 \pm 955.$$

The margin of error is \$955, and the 95% confidence interval is $9312 - 955 = \$8357$ to $9312 + 955 = \$10,267$. Thus, we are 95% confident that the mean credit card balance for the population of all households is between \$8357 and \$10,267.

9430	14661	7159	9071	9691	11032
7535	12195	8137	3603	11448	6525
4078	10544	9467	16804	8279	5239
5604	13659	12595	13479	5649	6195
5179	7061	7917	14044	11298	12584
4416	6245	11346	6817	4353	15415
10676	13021	12806	6845	3467	15917
1627	9719	4972	10493	6191	12591
10112	2200	11356	615	12851	9743
6567	10746	7117	13627	5337	10324
13627	12744	9465	12557	8372	
18719	5742	19263	6232	7445	

Practical Advice


1. If the population follows a normal distribution, the confidence interval provided by expression (8.2) is exact and can be used for any sample size.

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} \quad (8.2)$$

2. If the population does not follow a normal distribution, the confidence interval provided by expression (8.2) will be approximate. The quality of the approximation depends on both the distribution of the population and the sample size.
- (a) In most applications, a sample size of $n \geq 30$ is adequate when using expression (8.2) to develop an interval estimate of a population mean.
- (b) If the population distribution is highly skewed or contains outliers, most statisticians would recommend increasing the sample size to 50 or more.
- (c) If the population is not normally distributed but is roughly symmetric, sample sizes as small as 15 can be expected to provide good approximate confidence intervals.
- (d) With smaller sample sizes, expression (8.2) should only be used if the analyst believes, or is willing to assume, that the population distribution is at least approximately normal.

Using a Small Sample

1. Using a histogram of the sample data to learn about the distribution of a population is not always conclusive, but in many cases it provides the only information available. The histogram, along with judgment on the part of the analyst, can often be used to decide whether expression (8.2) can be used to develop the interval estimate.

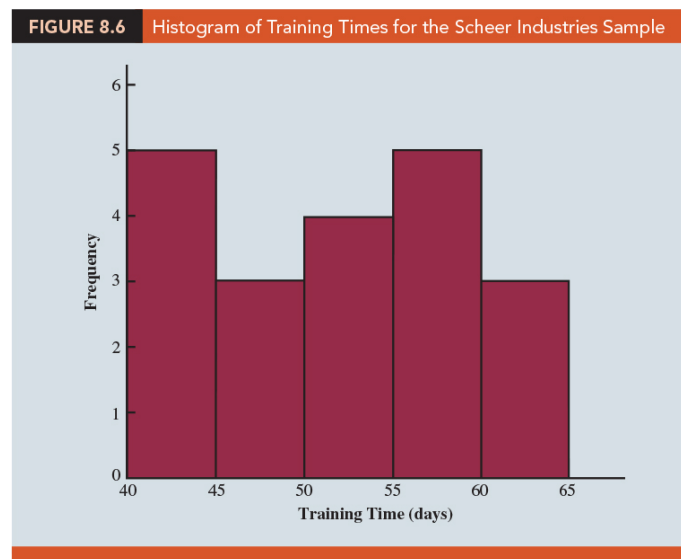
 **Question** (p385)

Scheer Industries is considering a new computer-assisted program to train maintenance employees to do machine repairs. In order to fully evaluate the program, the director of manufacturing requested an estimate of the population mean time required for maintenance employees to complete the computer-assisted training. A

sample of 20 employees is selected, with each employee in the sample completing the training program. Data on the training time in days for the 20 employees are shown in Table 8.4.

52	59	54	42
44	50	42	48
55	54	60	55
44	62	62	57
45	46	43	56

1. A histogram of the sample data is shown in Figure 8.6. What can we say about the distribution of the population based on this histogram?



2. With a 95% confidence, find the interval estimate for the population mean.

sol:

1. The sample data do not support the conclusion that the distribution of the population is normal, yet we do not see any evidence of skewness or outliers. Therefore, we conclude that an interval estimate based on the t distribution appears acceptable for the sample of 20 employees.
2. Compute the sample mean and sample standard deviation as

$$\bar{x} = \frac{\sum x_i}{n} = \frac{1030}{20} = 51.5 \text{ days} \quad s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n - 1}} = \sqrt{\frac{889}{20 - 1}} = 6.84 \text{ days}$$

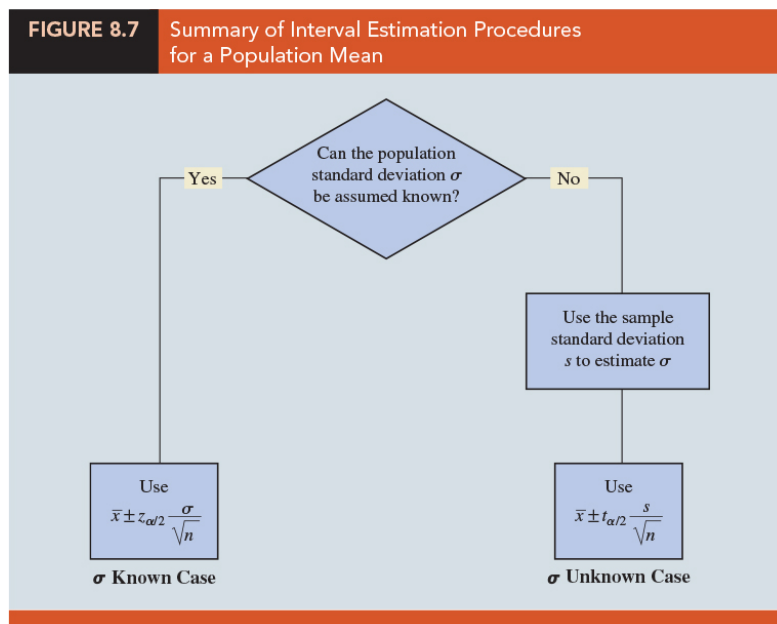
For a 95% confidence interval, $t_{19,0.025} = 2.093$. The interval estimate of the population mean is:

$$51.56 \pm 2.093 \left(\frac{6.84}{\sqrt{20}} \right) = 51.5 \pm 3.2$$

The 95% confidence interval is $51.5 - 3.2 = 48.3$ days to $51.5 + 3.2 = 54.7$ days.

Summary of Interval Estimation Procedures

1. Two approaches to develop an interval estimate of a population mean.
 - (a) For the σ known case, σ and the standard normal distribution are used in expression (8.1) to compute the margin of error and to develop the interval estimate.
 - (b) For the σ unknown case, the sample standard deviation s and the t distribution are used in expression (8.2) to compute the margin of error and to develop the interval estimate.



2. (Figure 8.7) In most applications, a sample size of $n \geq 30$ is adequate. If the population has a normal or approximately normal distribution, however, smaller sample sizes may be used. For the σ unknown case a sample size of $n \geq 50$ is recommended if the population distribution is believed to be highly skewed or has outliers.

8.3 Determining the Sample Size

1. How to choose a sample size large enough to provide a desired margin of error?
2. For the σ known case, the interval estimate is

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}},$$

The quantity $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ is the margin of error.

3. Once we select a confidence coefficient $1-\alpha$, $z_{\alpha/2}$ can be determined. Then, if we have a value for σ , we can determine the sample size n needed to provide any desired margin of error.

 **Question** (p390)

Suppose σ is known, using the interval estimate of a population mean at the chosen confidence level $1 - \alpha$ to develop the formula used to compute the required sample size n .

sol:


4. **Sample Size for an Interval Estimate of a Population Mean:**

$$n = \frac{(z_{\alpha/2})^2 \sigma^2}{E^2}, \quad E \text{ is the desired margin of error} \quad (8.3)$$

This sample size provides the desired margin of error at the chosen confidence level.

5. Although user preference must be considered, 95% confidence is the most frequently chosen value ($z_{0.025} = 1.96$).

6. However, even if σ is unknown, we can use equation (8.3) provided we have a preliminary or planning value for σ .
7. In practice, one of the following procedures can be chosen.
- Use the estimate of the population standard deviation computed from data of previous studies as the planning value for σ .
 - Use a pilot study to select a preliminary sample. The sample standard deviation from the preliminary sample can be used as the planning value for σ .
 - Use judgment or a "best guess" for the value of σ . For example, we might begin by estimating the largest and smallest data values in the population. The difference between the largest and smallest values provides an estimate of the range for the data. Finally, the range divided by 4 is often suggested as a rough approximation of the standard deviation and thus an acceptable planning value for σ .

 Question (p391)

A previous study that investigated the cost of renting automobiles in the United States found a mean cost of approximately \$55 per day for renting a midsize automobile. Suppose that the organization that conducted this study would like to conduct a new study in order to estimate the population mean daily rental cost for a midsize automobile in the United States. In designing the new study, the project director specifies that the population mean daily rental cost be estimated with a margin of error of \$2 and a 95% level of confidence. The project director specified a desired margin of error of $E = 2$, and the 95% level of confidence indicates $z_{0.025} = 1.96$. Thus, we only need a planning value for the population standard deviation σ in order to compute the required sample size. At this point, an analyst reviewed the sample data from the previous study and found that the sample standard deviation for the daily rental cost was \$9.65.

sol:

Using 9.65 as the planning value for σ , we obtain

$$n = \frac{(z_{\alpha/2})^2 \sigma^2}{E^2} = \frac{(1.96)^2 (9.65)^2}{2^2} = 89.43$$

Thus, the sample size for the new study needs to be at least 89.43 midsize automobile rentals in order to satisfy the project director's \$2 margin of error requirement. In cases where the computed n is not an integer, we round up to the next integer value; hence, the recommended sample size is 90 midsize automobile rentals.

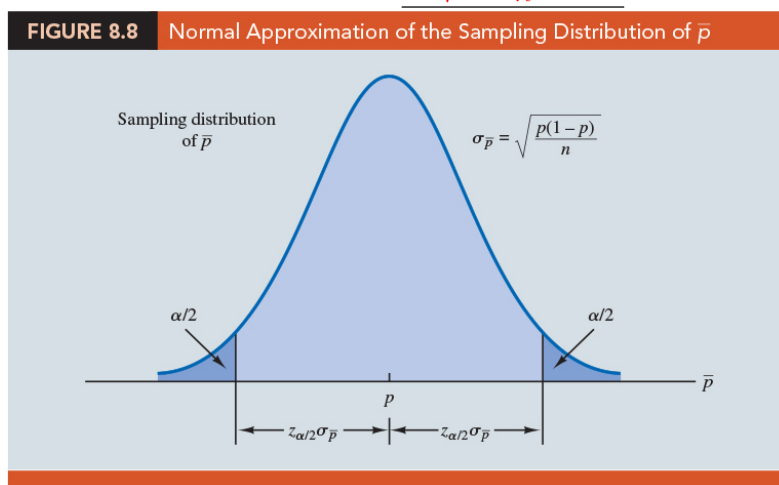
8.4 Population Proportion

1. The general form of an interval estimate of a population proportion p is

$$\underline{p \pm \text{Margin of error}}$$

2. (Chapter 7) Recall: the sampling distribution of \bar{p} can be approximated by a normal distribution whenever $np \geq 5$ and $n(1-p) \geq 5$.
3. (Figure 8.8) shows the normal approximation of the sampling distribution of \bar{p} . The mean of the sampling distribution of \bar{p} is the population proportion p , and the standard error of \bar{p} is

$$\sigma_{\bar{p}} = \sqrt{\frac{p(1-p)}{n}} \quad (8.4)$$



4. Because the sampling distribution of \bar{p} is normally distributed, if we choose $\frac{z_{\alpha/2}\sigma_{\bar{p}}}{n}$ as the margin of error in an interval estimate of a population proportion, we know that $100(1-\alpha)\%$ of the intervals generated will contain the true population proportion.
5. But $\sigma_{\bar{p}}$ cannot be used directly in the computation of the margin of error because p will not be known; p is what we are trying to estimate. So \bar{p} is substituted for p and the margin of error for an interval estimate of a population proportion is given by

$$\text{Margin of error} = \frac{z_{\alpha/2}\sqrt{\bar{p}(1-\bar{p})}}{n} \quad (8.5)$$

6. Interval Estimate of a Population Proportion:

$$\bar{p} \pm z_{\alpha/2}\sqrt{\frac{\bar{p}(1-\bar{p})}{n}} \quad (8.6)$$

where $1-\alpha$ is the confidence coefficient and $z_{\alpha/2}$ is the z value providing an area of $\alpha/2$ in the upper tail of the standard normal distribution.

Question (p394)

A national survey of 900 women golfers was conducted to learn how women golfers view their treatment at golf courses in the United States. The survey found that 396 of the women golfers were satisfied with the availability of tee times. Compute and interpret the 95% confidence interval estimate of the population proportion.

sol:

Interpretation: The margin of error is .0324 and the 95% confidence interval estimate of the population proportion is 0.4076 to 0.4724. Using percentages, the survey results enable us to state with 95% confidence that between 40.76% and 47.24% of all women golfers are satisfied with the availability of tee times.

Determining the Sample Size

1. How large the sample size should be to obtain an estimate of a population proportion at a specified level of precision.
2. Let E denote the desired margin of error:

$$E = \frac{z_{\alpha/2} \sqrt{\bar{p}(1 - \bar{p})}}{n}$$

3. Solving this equation for n provides a formula for the sample size that will provide a margin of error of size E :

$$n = \frac{(z_{\alpha/2})^2 \bar{p}(1 - \bar{p})}{E^2}$$

4. We cannot use this formula to compute the sample size that will provide the desired margin of error because \bar{p} will not be known until after we select the sample.
5. What we need is a planning value for \bar{p} (denoted by p^*) that can be used to make the computation.

6. Sample Size for an Interval Estimate of a Population Proportion:

$$n = \frac{(z_{\alpha/2})^2 p^*(1 - p^*)}{E^2}$$

7. In practice, the planning value p^* can be chosen by one of the following procedures:
 - (a) Use the sample proportion from a previous sample of the same or similar units.
 - (b) Use a pilot study to select a preliminary sample. The sample proportion from this sample can be used as the planning value, p^* .
 - (c) Use judgment or a "best guess" for the value of p^* .
 - (d) If none of the preceding alternatives applies, use a planning value of $p^* = 0.50$.

 Question (p395)

Consider the survey of women golfers and assume that the company is interested in conducting a new survey to estimate the current proportion of the population of women golfers who are satisfied with the availability of tee times. How large should the sample be if the survey director wants to estimate the population proportion with a margin of error of 0.025 at 95% confidence? Using the previous survey result of $p = 0.44$ as the planning value p^* .

sol:

8. The fourth alternative suggested for selecting a planning value $p^* = 0.50$ is frequently used when no other information is available.
9. The numerator of equation (8.7) shows that the sample size is proportional to the quantity $p^*(1-p^*)$. A larger value for the quantity $p^*(1-p^*)$ will result in a larger sample size.
10. (Table 8.5) in case of any uncertainty about an appropriate planning value, we know that $p^* = 0.50$ will provide the largest sample size recommendation. In effect, we play it safe by recommending the largest necessary sample size.

p^*	$p^*(1 - p^*)$
.10	$(.10)(.90) = .09$
.30	$(.30)(.70) = .21$
.40	$(.40)(.60) = .24$
.50	$(.50)(.50) = .25$
.60	$(.60)(.40) = .24$
.70	$(.70)(.30) = .21$
.90	$(.90)(.10) = .09$

← Largest value for $p^*(1 - p^*)$

- If the sample proportion turns out to be different from the 0.50 planning value, the margin of error will be smaller than anticipated. Thus, in using $p^* = 0.50$, we guarantee that the sample size will be sufficient to obtain the desired margin of error.
- In the survey of women golfers example, a planning value of $p^* = 0.50$ would have provided the sample size

$$n = \frac{(z_{\alpha/2})^2 p^*(1 - p^*)}{E^2} = \frac{(1.96)^2 (0.50)(1 - 0.50)}{(0.025)^2} = 1536.6$$

Thus, a slightly larger sample size of 1537 women golfers would be recommended.

8.5 Big Data and Confidence Intervals

- The confidence intervals are powerful tools for making inferences about population parameters.
- We now consider the ramifications of big data on confidence intervals for means and proportions.

Big Data and the Precision of Confidence Intervals

- The confidence intervals for the population mean μ and population proportion p become more narrow as the sample size increases. Therefore, the potential sampling error also decreases as the sample size increases.
- Example** Consider the online news service PenningtonDailyTimes.com (PDT). Prospective advertisers are willing to pay a premium to advertise on websites that have long visit times, so the time customers spend during their visits to PDT's website has a substantial impact on PDT's advertising revenues. Suppose PDT's management wants to develop a 95% confidence interval estimate of the mean amount of time customers spend during their visits to PDT's website.
- (Table 8.6) shows how the margin of error at the 95% confidence level decreases as the sample size increases when $s = 20$.

Sample Size n	Margin of Error $t_{\alpha/2} s_{\bar{x}}$
10	14.30714
100	3.96843
1,000	1.24109
10,000	.39204
100,000	.12396
1,000,000	.03920
10,000,000	.01240
100,000,000	.00392
1,000,000,000	.00124

- (Table 8.7) Suppose PDT would like to develop a 95% confidence interval estimate of the proportion of its website visitors that click on an ad. Table 8.7 shows how the margin of error for a 95% confidence interval estimate of the population proportion decreases as the sample size increases when the sample proportion is $\bar{p} = 0.51$.

TABLE 8.7 Margin of Error for Interval Estimates of the Population Proportion at the 95% Confidence Level for Various Sample Sizes n

Sample Size n	Margin of Error $z_{\alpha/2}\sigma_{\bar{p}}$
10	.30984
100	.09798
1,000	.03098
10,000	.00980
100,000	.00310
1,000,000	.00098
10,000,000	.00031
100,000,000	.00010
1,000,000,000	.00003

5. We see in Tables 8.6 and 8.7 that at a given confidence level, the margins of error decrease as the sample sizes increase.
- (a) If the sample mean time spent by customers when they visit PDT's website is 84.1 seconds, the 95% confidence interval estimate of the population mean time spent by customers when they visit PDT's website decreases from (69.79286, 98.40714) for a sample of $n = 10$ to (83.97604, 84.22396) for a sample of $n = 100,000$ to (84.09876, 84.10124) for a sample of $n = 1,000,000,000$.
- (b) Similarly, if the sample proportion of its website visitors who clicked on an ad is 0.51, the 95% confidence interval estimate of the population proportion of its website visitors who clicked on an ad decreases from (0.20016, 0.81984) for a sample of $n = 10$ to (0.50690, 0.51310) for a sample of $n = 100,000$ to (0.50997, 0.51003) for a sample of $n = 1,000,000,000$.
6. In both instances, as the sample size becomes extremely large, the margin of error becomes extremely small and the resulting confidence intervals become extremely narrow.

Implications of Big Data for Confidence Intervals

- Example** Last year the mean time spent by all visitors to PenningtonDailyTimes.com was 84 seconds. Suppose that PDT wants to assess whether the population mean time has changed since last year. PDT now collects a new sample of 1,000,000 visitors to its website and calculates the sample mean time spent by these visitors

to the PDT website to be $\bar{x} = 84.1$ seconds. The estimated population standard deviation is $s = 20$ seconds, so the standard error is $s_{\bar{x}} = s/\sqrt{n} = 0.02000$.

2. Furthermore, the sample is sufficiently large to ensure that the sampling distribution of the sample mean will be normally distributed. Thus, the 95% confidence interval estimate of the population mean is

$$\bar{x} \pm t_{\alpha/2}s_{\bar{x}} = 84.1 \pm 0.03925 = (84.06080, 84.13920)$$

3. What could PDT conclude from these results? There are three possible reasons that PDT's sample mean of 84.1 seconds differs from last year's population mean of 84 seconds: (1) sampling error, (2) nonsampling error, or (3) the population mean has changed since last year.

(a) The 95% confidence interval estimate of the population mean does not include the value for the mean time spent by all visitors to the PDT website for last year (84 seconds), suggesting that the difference between PDT's sample mean for the new sample (84.1 seconds) and the mean from last year (84 seconds) is not likely to be exclusively a consequence of sampling error.

(b) Nonsampling error is a possible explanation and should be investigated as the results of statistical inference become less reliable as nonsampling error is introduced into the sample data.

(c) If PDT determines that it introduced little or no nonsampling error into its sample data, the only remaining plausible explanation for a difference of this magnitude is that the population mean has changed since last year.

i. If PDT concludes that the sample has provided reliable evidence and the population mean has changed since last year, management must still consider the potential impact of the difference between the sample mean and the mean from last year.

ii. If a 0.1 second difference in the time spent by visitors to PenningtonDailyTimes.com has a consequential effect on what PDT can charge for advertising on its site, this result could have practical business implications for PDT.

- iii. Otherwise, there may be no practical significance of the 0.1 second difference in the time spent by visitors to PenningtonDailyTimes.com.
4. Confidence intervals are extremely useful, but as with any other statistical tool, they are only effective when properly applied.
5. Because interval estimates become increasingly precise as the sample size increases, extremely large samples will yield extremely precise estimates.
6. However, no interval estimate, no matter how precise, will accurately reflect the parameter being estimated unless the sample is relatively free of nonsampling error. Therefore, when using interval estimation, it is always important to carefully consider whether a random sample of the population of interest has been taken.

😊 EXERCISES

8.1 : 3, 8, 9

8.2 : 13, 16, 18, 20

8.3 : 24, 25, 29

8.4 : 31, 33, 37, 39, 43

8.5 : 47

SUP : 50, 51, 54, 56, 62

“你要搞清楚自己人生的劇本 – 不是你父母的續集，不是你子女的前傳，更不是你朋友的外篇。對待生命你不妨大膽冒險一點，因為好歹你要失去它。生命中最難的階段不是沒有人懂你，而是你不懂你自己。”

“You want to figure out your life’s script – not your parent’s sequel, not your child’s prequel, but not just a chapter in your friend’s book. Treat life like you want adventure and to be bold; because whatever the outcome, you can’t live forever. The hardest part of life is not that no one understands you, but you do not understand yourself.”

— *Friedrich Wilhelm Nietzsche (October 15, 1844 – August 25, 1900)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 9: Hypothesis Tests

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

系級: _____ 學號: _____ 姓名: _____

Overview

1. Statistical inference: how hypothesis testing can be used to determine whether a statement about the value of a population parameter should or should not be rejected.
2. The null hypothesis (H_0): making a tentative assumption about a population parameter.
3. The alternative hypothesis (H_a): the opposite of what is stated in H_0 .
4. The hypothesis testing procedure uses data from a sample to test the two competing statements indicated by H_0 and H_a .
5. This chapter shows how hypothesis tests can be conducted about a population mean and a population proportion.

9.1 Developing Null and Alternative Hypotheses

1. It is not always obvious how the null and alternative hypotheses should be formulated.
2. All hypothesis testing applications involve collecting a sample and using the sample results to provide evidence for drawing a conclusion.
3. In some situations it is easier to identify H_a first and then develop H_0 .
4. In other situations it is easier to identify H_0 first and then develop H_a .

The Alternative Hypothesis as a Research Hypothesis

1. Many applications of hypothesis testing involve an attempt to gather evidence in support of a research hypothesis. In these situations, it is often best to begin with the alternative hypothesis and make it the conclusion that the researcher hopes to support.
2. **Example** Consider a particular automobile that currently attains a fuel efficiency of 24 miles per gallon in city driving.
 - (a) *Goal:* A product research group has developed a new fuel injection system (燃料噴射系統) designed to increase the miles-per-gallon rating. The group will run controlled tests with the new fuel injection system looking for statistical support for the conclusion that the new fuel injection system provides more miles per gallon than the current system.
 - (b) Several new fuel injection units will be manufactured, installed in test automobiles, and subjected to research-controlled driving conditions.
 - (c) The sample mean miles per gallon for these automobiles will be computed and used in a hypothesis test to determine if it can be concluded that the new system provides more than 24 miles per gallon.
 - (d) In terms of the population mean miles per gallon μ , the research hypothesis $\mu > 24$ becomes the alternative hypothesis.

- (e) Since the current system provides an average or mean of 24 miles per gallon, we will make the tentative assumption that the new system is not any better than the current system and choose $\mu \leq 24$ as the null hypothesis.

$$\underline{H_0 : \mu \leq 24, \quad H_a : \mu > 24}$$

- (f) If the sample results lead to the conclusion to reject H_0 , the inference can be made that $H_a : \mu > 24$ is true.
- (g) The researchers have the statistical support to state that the new fuel injection system increases the mean number of miles per gallon.
- (h) If the sample results lead to the conclusion that H_0 cannot be rejected, the researchers cannot conclude that the new fuel injection system is better than the current system. Production of automobiles with the new fuel injection system on the basis of better gas mileage cannot be justified. Perhaps more research and further testing can be conducted.
3. Before adopting something new (e.g., products, methods, systems), it is desirable to conduct research to determine if there is statistical support for the conclusion that the new approach is indeed better. In such cases, the research hypothesis is stated as the alternative hypothesis.
- (a) **Example** A new teaching method is developed that is believed to be better than the current method.
- H_0 : the new method is no better than the old method.
 - H_a : the new method is better.
- (b) **Example** A new sales force bonus plan is developed in an attempt to increase sales.
- H_0 : the new bonus plan does not increase sales.
 - H_a : the new bonus plan increases sales.
- (c) **Example** A new drug is developed with the goal of lowering blood pressure more than an existing drug.
- H_0 : the new drug does not provide lower blood pressure than the existing drug.

- ii. H_a : the new drug lowers blood pressure more than the existing drug.
4. In each case, rejection of the null hypothesis H_0 provides statistical support for the research hypothesis.

The Null Hypothesis as an Assumption to Be Challenged

- The situations below that it is helpful to develop the null hypothesis first.
 - Consider applications of hypothesis testing where we begin with a belief or an assumption that a statement about the value of a population parameter is true.
 - We will then use a hypothesis test to challenge the assumption and determine if there is statistical evidence to conclude that the assumption is incorrect.
- The null hypothesis H_0 expresses the belief or assumption about the value of the population parameter. The alternative hypothesis H_a is that the belief or assumption is incorrect.
- Example** Consider the situation of a manufacturer of soft drink products.
 - The label on a soft drink bottle states that it contains 67.6 fluid ounces. We consider the label correct provided the population mean filling weight for the bottles is at least 67.6 fluid ounces.
 - We would begin with the assumption that the label is correct and state the null hypothesis as $\mu \geq 67.6$.
 - The challenge to this assumption would imply that the label is incorrect and the bottles are being under-filled. This challenge would be stated as the alternative hypothesis $\mu < 67.6$.

$$\underline{H_0 : \mu \geq 67.6 \quad H_a : \mu < 67.6}$$
 - A government agency with the responsibility for validating manufacturing labels could select a sample of soft drinks bottles, compute the sample mean filling weight, and use the sample results to test the preceding hypotheses.

- (e) If the sample results lead to the conclusion to reject H_0 , the inference that $H_a : \mu < 67.6$ is true can be made. With this statistical support, the agency is justified in concluding that the label is incorrect and underfilling of the bottles is occurring.
- (f) If the sample results indicate H_0 cannot be rejected, the assumption that the manufacturer's labeling is correct cannot be rejected. With this conclusion, no action would be taken.
4. **Example** Consider the soft drink bottle filling example from the manufacturer's point of view.
- (a) The bottle-filling operation has been designed to fill soft drink bottles with 67.6 fluid ounces as stated on the label.
- The company does not want to underfill the containers because that could result in an underfilling complaint from customers or, perhaps, a government agency.
 - However, the company does not want to overfill containers either because putting more soft drink than necessary into the containers would be an unnecessary cost.
- (b) The company's goal would be to adjust the bottle-filling operation so that the population mean filling weight per bottle is 67.6 fluid ounces as specified on the label.
- (c) In a hypothesis testing application, we would begin with the assumption that the production process is operating correctly and state the null hypothesis as $\mu = 67.6$ fluid ounces.
- (d) The alternative hypothesis that challenges this assumption is that $\mu \neq 67.6$, which indicates either overfilling or underfilling is occurring.

$$\underline{H_0 : \mu = 67.6 \quad H_a : \mu \neq 67.6.}$$

- (e) Suppose that the soft drink manufacturer uses a quality control procedure to periodically select a sample of bottles from the filling operation and computes the sample mean filling weight per bottle.

- i. If the sample results lead to the conclusion to reject H_0 , the inference is made that $H_a : \mu \neq 67.6$ is true. We conclude that the bottles are not being filled properly and the production process should be adjusted to restore the population mean to 67.6 fluid ounces per bottle.
 - ii. If the sample results indicate H_0 cannot be rejected, the assumption that the manufacturer's bottle filling operation is functioning properly cannot be rejected. In this case, no further action would be taken and the production operation would continue to run.
5. The two preceding forms of the soft drink manufacturing hypothesis test show that the null and alternative hypotheses may vary depending upon the point of view of the researcher or decision maker.
 6. To correctly formulate hypotheses it is important to understand the context of the situation and structure the hypotheses to provide the information the researcher or decision maker wants.

Summary Of Forms for Null and Alternative Hypotheses

1. Depending on the situation, hypothesis tests about a population parameter (the population mean and the population proportion) may take one of three forms:

$H_0 : \mu \geq \mu_0$	$H_0 : \mu \leq \mu_0$	$H_0 : \mu = \mu_0$
$H_a : \mu < \mu_0$	$H_a : \mu > \mu_0$	$H_a : \mu \neq \mu_0$
2. The first two forms are called one-tailed tests. The third form is called a two-tailed test.
3. The equality part of the expression (either \geq , \leq , or $=$) always appears in the null hypothesis.
4. In selecting the proper form of H_0 and H_a , keep in mind that the alternative hypothesis is often what the test is attempting to establish. Hence, asking whether the user is looking for evidence to support $\mu < \mu_0$, $\mu > \mu_0$, or $\mu \neq \mu_0$ will help determine H_a .

9.2 Type I and Type II Errors

1. Ideally the hypothesis testing procedure should lead to the acceptance of H_0 when H_0 is true and the rejection of H_0 when H_a is true.
2. (Table 9.1) The correct conclusions are not always possible.

		Population Condition	
		H_0 True	H_a True
Conclusion	Accept H_0	Correct Conclusion	Type II Error
	Reject H_0	Type I Error	Correct Conclusion

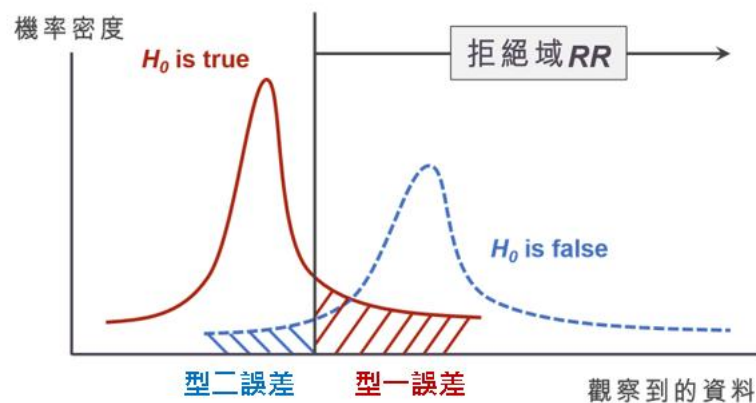
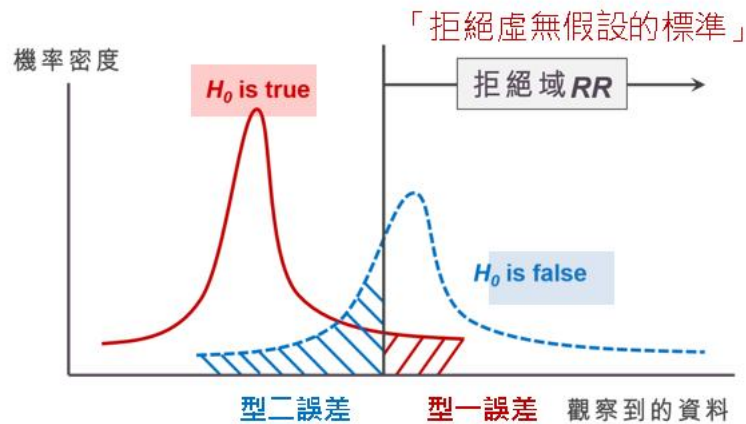
- (a) We reject H_0 if H_0 is true, we make a Type I error.
 - (b) If H_a is true, the conclusion is correct when we reject H_0 .
 - (c) If H_0 is true, the conclusion is correct when we accept H_0 .
 - (d) If H_0 is false (H_a is true), we make a Type II error when we accept H_0 .
3. **Example** An automobile product research group developed a new fuel injection system designed to increase the miles-per-gallon rating of a particular automobile.
 - (a) With the current model obtaining an average of 24 miles per gallon, the hypothesis test was formulated as follows.

$$\underline{H_0 : \mu \leq 24 \text{ against } H_a : \mu > 24}$$

- (b) The alternative hypothesis, $H_a : \mu > 24$, indicates that the researchers are looking for sample evidence to support the conclusion that the population mean miles per gallon with the new fuel injection system is greater than 24.
- (c) *Type I error*: rejecting H_0 when it is true corresponds to the researchers claiming that the new system improves the miles-per-gallon rating ($\mu > 24$) when in fact the new system is not any better than the current system.

- (d) *Type II error*: accepting H_0 when it is false corresponds to the researchers concluding that the new system is not any better than the current system ($\mu \leq 24$) when in fact the new system improves miles-per-gallon performance.
4. **Level of Significance**: The level of significance is the probability of making a Type I error when the null hypothesis is true as an equality.
- (a) **Example** For the miles-per-gallon rating hypothesis test, the null hypothesis is $H_0 : \mu \leq 24$. Suppose the null hypothesis is true as an equality; that is, $\mu = 24$. The level of significance is the probability of rejecting $H_0 : \mu \leq 24$ when $\mu = 24$.
- (b) The Greek symbol α (alpha) is used to denote the level of significance, and common choices for α are 0.05 and 0.01.
- (c) In practice, the person responsible for the hypothesis test specifies the level of significance. By selecting α , that person is controlling the probability of making a Type I error.
- (d) If the cost of making a Type I error is high (not too high), small (larger) values of α are preferred.
5. **The significance tests**: Applications of hypothesis testing that only control for the Type I error are called significance tests.
6. Although most applications of hypothesis testing control for the probability of making a Type I error, they do not always control for the probability of making a Type II error.
- (a) Hence, if we decide to accept H_0 , we cannot determine how confident we can be with that decision. Because of the uncertainty associated with making a Type II error when conducting significance tests, statisticians usually recommend that we use the statement "do not reject H_0 " instead of "accept H_0 ".
- (b) Using the statement "do not reject H_0 " carries the recommendation to withhold both judgment and action. In effect, by not directly accepting H_0 , the statistician avoids the risk of making a Type II error.

- (c) Whenever the probability of making a Type II error has not been determined and controlled, we will not make the statement "accept H_0 ." In such cases, only two conclusions are possible: do not reject H_0 or reject H_0 .
- (d) Although controlling for a Type II error in hypothesis testing is not common, it can be done. In Sections 9.7 and 9.8 we will illustrate procedures for determining and controlling the probability of making a Type II error. If proper controls have been established for this error, action based on the "accept H_0 " conclusion can be appropriate.



9.3 Population Mean: σ Known

One-tailed Test

1. One-tailed tests about a population mean take one of the following two forms:

Lower Tail Test	Upper Tail Test
$H_0 : \mu \geq \mu_0$	$H_0 : \mu \leq \mu_0$
$H_a : \mu < \mu_0$	$H_a : \mu > \mu_0$

2. Example The Federal Trade Commission (FTC) periodically conducts statistical studies designed to test the claims that manufacturers make about their products. For example, the label on a large can of Hilltop Coffee states that the can contains 3 pounds of coffee. The FTC knows that Hilltop's production process cannot place exactly 3 pounds of coffee in each can, even if the mean filling weight for the population of all cans filled is 3 pounds per can ($\mu_0 = 3$). However, as long as the population mean filling weight is at least 3 pounds per can, the rights of consumers will be protected. Thus, the FTC interprets the label information on a large can of coffee as a claim by Hilltop that the population mean filling weight is at least 3 pounds per can. We will show how the FTC can check Hilltop's claim by conducting a lower tail hypothesis test.

- (a) *Develop the null and alternative hypotheses for the test.* If the population mean filling weight is at least 3 pounds per can, Hilltop's claim is correct:

$$H_0 : \underline{\mu \geq 3} \quad H_a : \underline{\mu < 3}$$

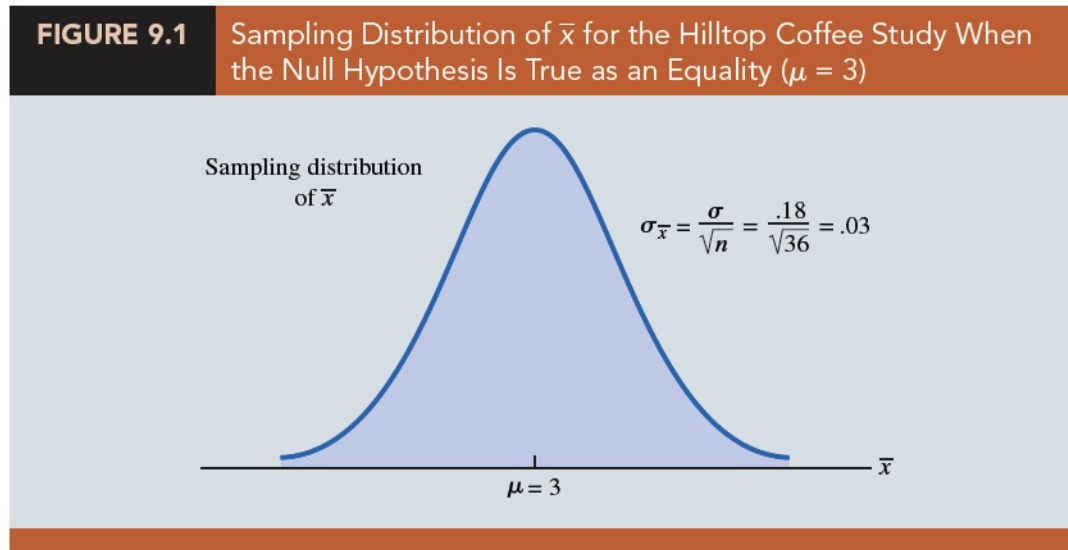
- i. If the sample data indicate that H_0 cannot be rejected, no action should be taken against Hilltop.
- ii. If the sample data indicate H_0 can be rejected, $H_a : \mu < 3$, is true. A conclusion of underfilling and a charge of a label violation against Hilltop would be justified.
- iii. Suppose a sample of $n = 36$ cans of coffee is selected and the sample mean \bar{x} is computed as an estimate of the population mean μ . If the value of the sample mean \bar{x} is less than 3 pounds, the sample results will cast doubt on the null hypothesis.

- (b) *Specifying the level of significance, α :*
- i. (Recall) The level of significance is the probability of making a Type I error by rejecting H_0 when H_0 is true as an equality.
 - ii. If the cost of making a Type I error is high (not high), a small (larger) value should be chosen for the level of significance.
 - iii. In the Hilltop Coffee study, the director of the FTC's testing program made the following statement: "If the company is meeting its weight specifications at $\mu = 3$, I do not want to take action against them. But, I am willing to risk a 1% chance of making such an error." From the director's statement, we set the level of significance for the hypothesis test at $\alpha = 0.01$.
 - iv. Thus, we must design the hypothesis test so that the probability of making a Type I error when $\mu = 3$ is 0.01.
3. By developing the null and alternative hypotheses and specifying the level of significance for the test, we carry out the first two steps required in conducting every hypothesis test. We are now ready to perform the third step of hypothesis testing: collect the sample data and compute the value of what is called a test statistic.

Test statistic

1. Example For the Hilltop Coffee study, previous FTC tests show that the population standard deviation can be assumed known with a value of $\sigma = 0.18$. These tests also show that the population of filling weights can be assumed to have a normal distribution.
2. The sampling distribution of \bar{x} is normally distributed with a known value of $\sigma = 0.18$ and a sample size of $n = 36$.
3. (Figure 9.1) the sampling distribution of \bar{x} when the null hypothesis is true as an equality ($\mu = \mu_0 = 3$). The standard error of \bar{x} is given by $\sigma_{\bar{x}} = \underline{\sigma/\sqrt{n} = 0.18/\sqrt{36} = 0.03}$.
The sampling distribution of

$$z = \frac{\bar{x} - \mu_0}{\sigma_{\bar{x}}} = \frac{\bar{x} - 3}{0.03} \text{ is a standard normal distribution.}$$



4. A value of $z = -1$ ($z = -3$) means that the value of \bar{x} is one (three) standard error below the hypothesized value of the mean.
5. The lower tail area at $z = -3.00$ is 0.0013. Hence, the probability of obtaining a value of z that is three or more standard errors below the mean is 0.0013.
6. The probability of obtaining a value of \bar{x} that is 3 or more standard errors below the hypothesized population mean $\mu_0 = 3$ is also 0.0013. Such a result is unlikely if the null hypothesis is true.
7. For hypothesis tests about a population mean in the σ known case, we use the standard normal random variable z as a test statistic to determine whether \bar{x} deviates from the hypothesized value of μ enough to justify rejecting the null hypothesis.
8. **Test Statistic for Hypothesis Tests About a Population Mean: σ Known**

$$z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

9. The key question for a lower tail test is, How small must the test statistic z be before we choose to reject H_0 ? Two approaches: the p -value approach and the critical value approach.

***p*-value approach**

1. ***p*-value:** A *p*-value is a probability that provides a measure of the evidence against H_0 provided by the sample.
 - (a) The *p*-value is used to determine whether H_0 should be rejected.
 - (b) A small *p*-value indicates the value of the test statistic is unusual given the assumption that H_0 is true.
 - (c) Smaller *p*-values indicate more evidence against H_0 .
2. The value of the test statistic is used to compute the *p*-value.
 - (a) For a lower tail test, the *p*-value is the probability of obtaining a value for the test statistic as small as or smaller than that provided by the sample.
 - (b) To compute the *p*-value for the lower tail test in the σ known case, we use the standard normal distribution to find the probability that z is less than or equal to the value of the test statistic.
 - (c) After computing the *p*-value, we must then decide whether it is small enough to reject the null hypothesis; as we will show, this decision involves comparing the *p*-value to the level of significance.

 **Question** (p427)

Suppose the sample of 36 Hilltop coffee cans provides a sample mean of $\bar{x} = 2.92$ pounds. Is $\bar{x} = 2.92$ small enough to cause us to reject H_0 ? Compute the *p*-value for the Hilltop Coffee lower tail test.

sol:

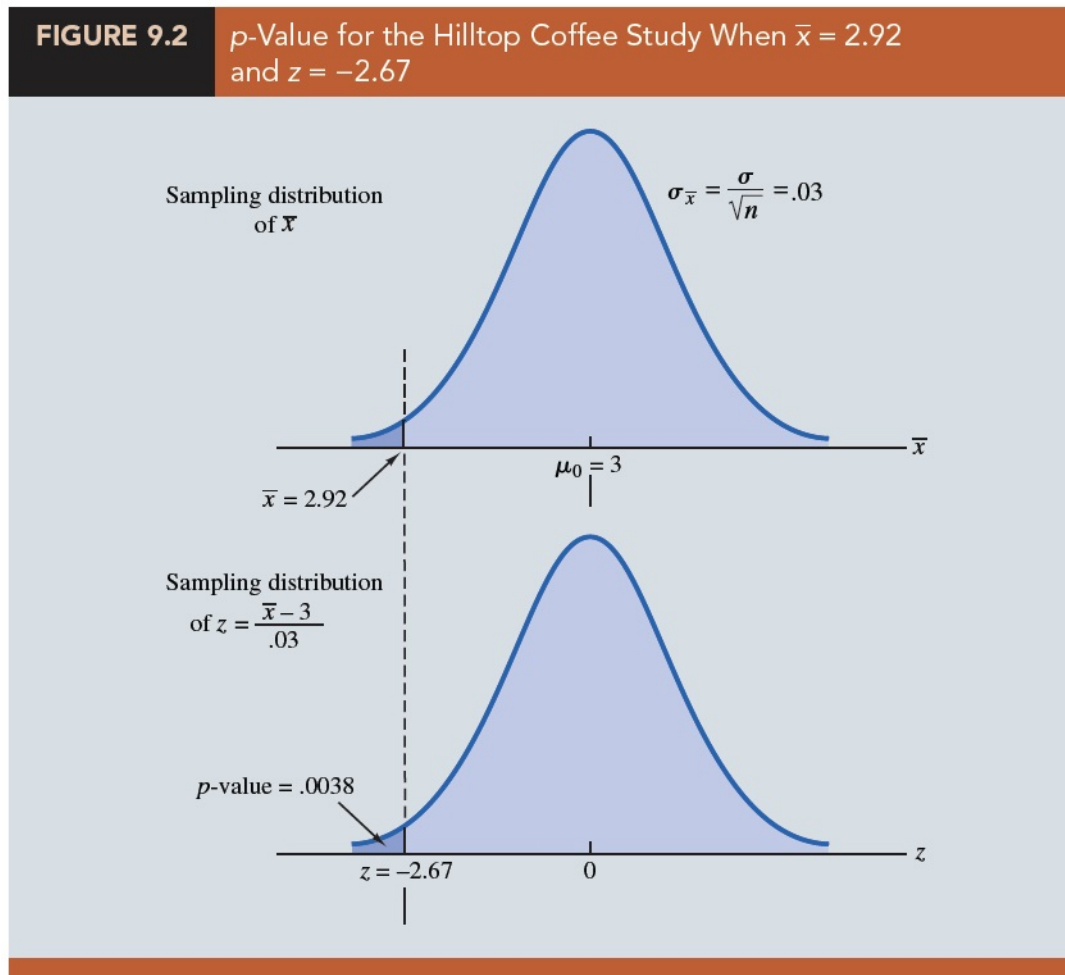
- Because this is a lower tail test, the *p*-value is the area under the standard normal curve for values of $z \leq$ the value of the test statistic.

– Using $\bar{x} = 2.92$, $\sigma = 0.18$, and $n = 36$, the value of the test statistic z :

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{2.92 - 3}{0.18/\sqrt{36}} = -2.67$$

– p -value $= P(z \leq -2.67) = 0.0038$.

– (Figure 9.2) $\bar{x} = 2.92$ corresponds to $z = -2.67$ and a p -value = 0.0038.



3. This p -value (0.0038) indicates a small probability of obtaining a sample mean of $\bar{x} = 2.92$ (and a test statistic of -2.67) or smaller when sampling from a population with $\mu = 3$.

4. This p -value does not provide much support for the null hypothesis, but is it small enough to cause us to reject H_0 ? The answer depends upon α for the test.
5. As noted previously, the director of the FTC's testing program selected a value of 0.01 for the level of significance means that the director is willing to tolerate a probability of 0.01 of rejecting H_0 when it is true as an equality ($\mu_0 = 3$).
6. The sample of 36 coffee cans in the Hilltop Coffee study resulted in a p -value = 0.0038, which means that the probability of obtaining a value of $\bar{x} = 2.92$ or less when H_0 is true as an equality is 0.0038.
7. Because $p\text{-value} = 0.0038 \leq \alpha = 0.01$, we reject H_0 . Therefore, we find sufficient statistical evidence to reject the null hypothesis at the 0.01 level of significance.
8. **Rejection Rule Using p -value.** For a level of significance α , the rejection rule using the p -value approach is:

$$\underline{\text{Reject } H_0 \text{ if } p\text{-value} \leq \alpha}$$

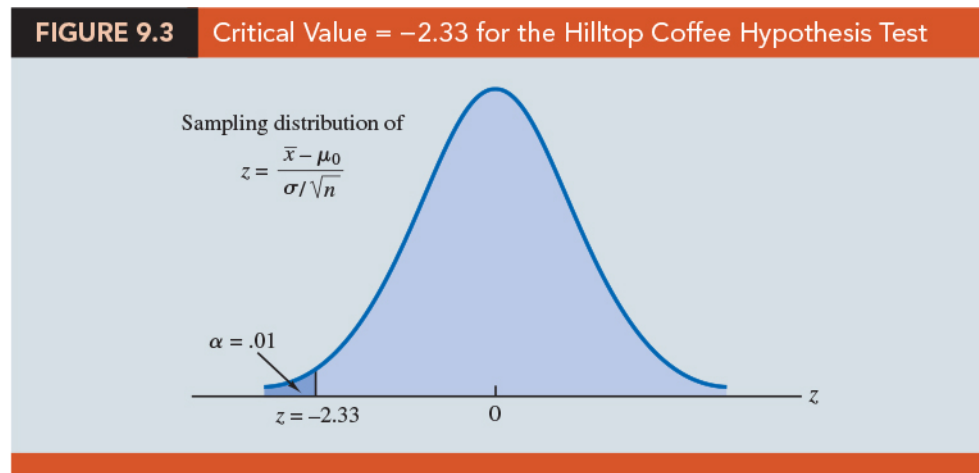
9. **Example** In the Hilltop Coffee test, the p -value of 0.0038 resulted in the rejection of H_0 . The observed p -value of 0.0038 means that we would reject H_0 for any value of $\alpha \geq 0.0038$. For this reason, the p -value is also called the observed level of significance.
10. Different decision makers may express different opinions concerning the cost of making a Type I error and may choose a different α .

Critical value approach

1. The critical value is the value of the test statistic that corresponds to an area of α in the lower tail of the sampling distribution of the test statistic.
2. The critical value is the largest value of the test statistic that will result in the rejection of the null hypothesis.
3. For a lower tail test, the critical value serves as a benchmark for determining whether the value of the test statistic is small enough to reject the null hypothesis.

4. **Example** the Hilltop Coffee example.

- (a) In the σ known case, the sampling distribution for the test statistic z is a standard normal distribution. Therefore, the critical value is the value of the test statistic that corresponds to an area of $\alpha = 0.01$ in the lower tail of a standard normal distribution.
- (b) (Figure 9.3) We find that $-z_{0.01} = -2.33$ provides an area of 0.01 in the lower tail, $P(z \leq -2.33) = 0.01$.
- (c) If the sample results in a value of the test statistic that is less than or equal to -2.33 , the corresponding p -value will be less than or equal to 0.01; in this case, we should reject H_0 .



- (d) Hence, for the Hilltop Coffee study the critical value rejection rule for a level of significance of 0.01 is

$$\text{Reject } H_0 \text{ if } \underline{z \leq -2.33}$$

- (e) In the Hilltop Coffee example, $\bar{x} = 2.92$ and the test statistic is $z = -2.67$. Because $z = -2.67 < -2.33$, we can reject H_0 and conclude that Hilltop Coffee is underfilling cans.

5. Rejection Rule for a Lower Tail Test: Critical Value Approach. We can generalize the rejection rule for the critical value approach to handle any level of significance. The rejection rule for a lower tail test follows.

$$\text{Reject } H_0 \text{ if } \underline{z \leq -z_\alpha}$$

where $-z_\alpha$ is the critical value; that is, the z value that provides an area of α in the lower tail of the standard normal distribution.

Summary

1. The p -value approach to hypothesis testing and the critical value approach will always lead to the same rejection decision.
2. The advantage of the p -value approach is that the p -value tells us how significant the results are (the observed level of significance).
3. If we use the critical value approach, we only know that the results are significant at the stated α .
4. We can use the same general approach to conduct an upper tail test. The test statistic z is still computed using equation (9.1). But, for an upper tail test, the p -value is the probability of obtaining a value for the test statistic as large as or larger than that provided by the sample.
5. To compute the p -value for the upper tail test in the σ known case, we must use the standard normal distribution to find the probability that z is greater than or equal to the value of the test statistic.

6. Computation of p -Values for One-Tailed Tests

- (a) Compute the value of the test statistic using equation (9.1):

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

- (b) *Lower tail test*: Using the standard normal distribution, compute the probability that z is less than or equal to the value of the test statistic (area in the lower tail).
- (c) *Upper tail test*: Using the standard normal distribution, compute the probability that z is greater than or equal to the value of the test statistic (area in the upper tail).

Two-tailed Test

1. The general form for a two-tailed test about a population mean:

$$\underline{H_0 : \mu = \mu_0, \quad H_a : \mu \neq \mu_0}$$

2. **Example** The U.S. Golf Association (USGA) establishes rules that manufacturers of golf equipment must meet if their products are to be acceptable for use in USGA events. MaxFlight Inc. uses a high-technology manufacturing process to produce golf balls with a mean driving distance of 295 yards. Sometimes, however, the process gets out of adjustment and produces golf balls with a mean driving distance different from 295 yards. When the mean distance falls below 295 yards, the company worries about losing sales because the golf balls do not provide as much distance as advertised. When the mean distance passes 295 yards, MaxFlight's golf balls may be rejected by the USGA for exceeding the overall distance standard concerning carry and roll. MaxFlight's quality control program involves taking periodic samples of 50 golf balls to monitor the manufacturing process. For each sample, a hypothesis test is conducted to determine whether the process has fallen out of adjustment.

- (a) We begin by assuming that the process is functioning correctly; that is, the golf balls being produced have a mean distance of 295 yards. This assumption establishes the null hypothesis. The alternative hypothesis is that the mean distance is not equal to 295 yards.

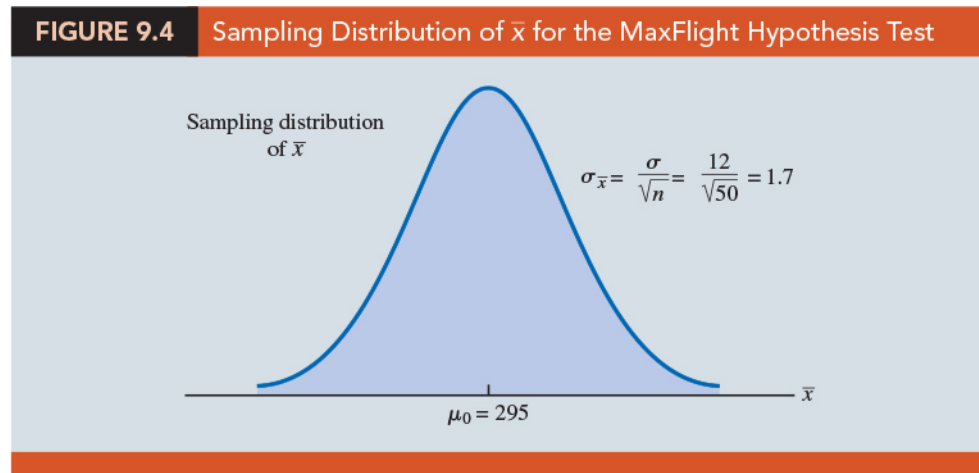
$$\underline{H_0 : \mu = 295, \quad H_a : \mu \neq 295} .$$

- (b) If the sample mean \bar{x} is significantly less than 295 yards or significantly greater than 295 yards, we will reject H_0 . In this case, corrective action will be taken to adjust the manufacturing process.
- (c) If \bar{x} does not deviate from the hypothesized mean $\mu_0 = 295$ by a significant amount, H_0 will not be rejected and no action will be taken to adjust the manufacturing process.
- (d) The quality control team selected $\alpha = 0.05$ as the level of significance for the test. Data from previous tests conducted when the process was known to be

in adjustment show that the population standard deviation can be assumed known with a value of $\sigma = 12$. Thus, with a sample size of $n = 50$, the standard error of \bar{x} is

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{12}{\sqrt{50}} = 1.7$$

- (e) Because the sample size is large, the central limit theorem allows us to conclude that the sampling distribution of \bar{x} can be approximated by a normal distribution.
- (f) (Figure 9.4) the sampling distribution of \bar{x} for the MaxFlight hypothesis test with a hypothesized population mean of $\mu_0 = 295$.



3. Suppose that a sample of 50 golf balls is selected and that the sample mean is $\bar{x} = 297.6$ yards. This sample mean provides support for the conclusion that the population mean is larger than 295 yards. Is this value of \bar{x} enough larger than 295 to cause us to reject H_0 at the 0.05 level of significance?

***p*-value approach**

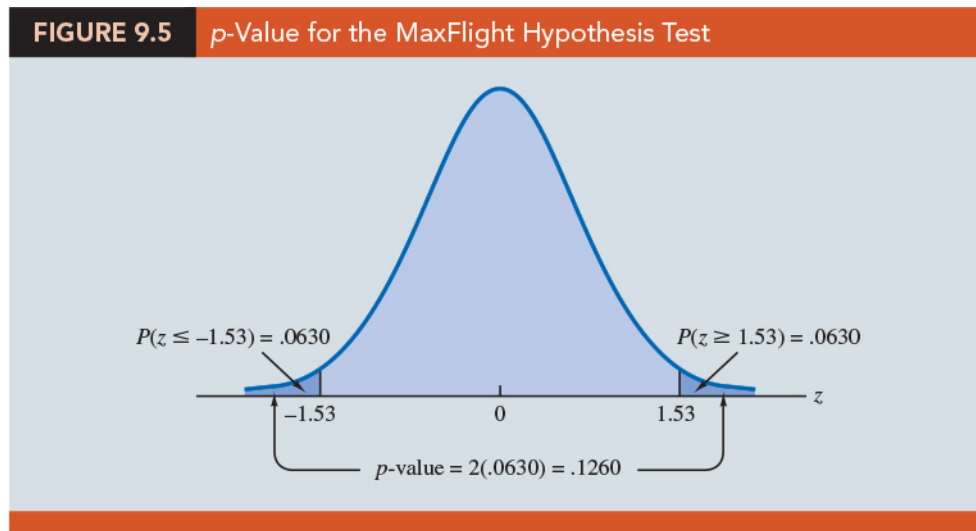
- (Recall) the *p*-value is a probability used to determine whether the null hypothesis should be rejected.
- For a two-tailed test, values of the test statistic in either tail provide evidence against the null hypothesis.

3. For a two-tailed test, the p -value is the probability of obtaining a value for the test statistic as unlikely as or more unlikely than that provided by the sample.
4. **Example** the MaxFlight hypothesis test example.

(a) *Compute the value of the test statistic.* For the σ known case, the test statistic z is a standard normal random variable. Using equation (9.1) with $\bar{x} = 297.6$, the value of the test statistic is

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{297.6 - 295}{12/\sqrt{50}} = 1.53$$

- (b) *Compute the p -value.* Find the probability of obtaining a value for the test statistic at least as unlikely as $z = 1.53$. Clearly values of $z \geq 1.53$ are at least as unlikely.
- (c) But, because this is a two-tailed test, values of $z \leq -1.53$ are also at least as unlikely as the value of the test statistic provided by the sample.
- (d) (Figure 9.5) the two-tailed p -value: $P(z \leq -1.53) + P(z \geq 1.53)$.
- (e) Because the normal curve is symmetric, $P(z < 1.53) = 0.9370$. Thus, the upper tail area is $P(z \geq 1.53) = \underline{1.0000 - 0.9370 = 0.0630}$.



(f) The p -value for the MaxFlight two-tailed hypothesis test is

$$p\text{-value} = \underline{2(0.0630) = 0.1260}$$

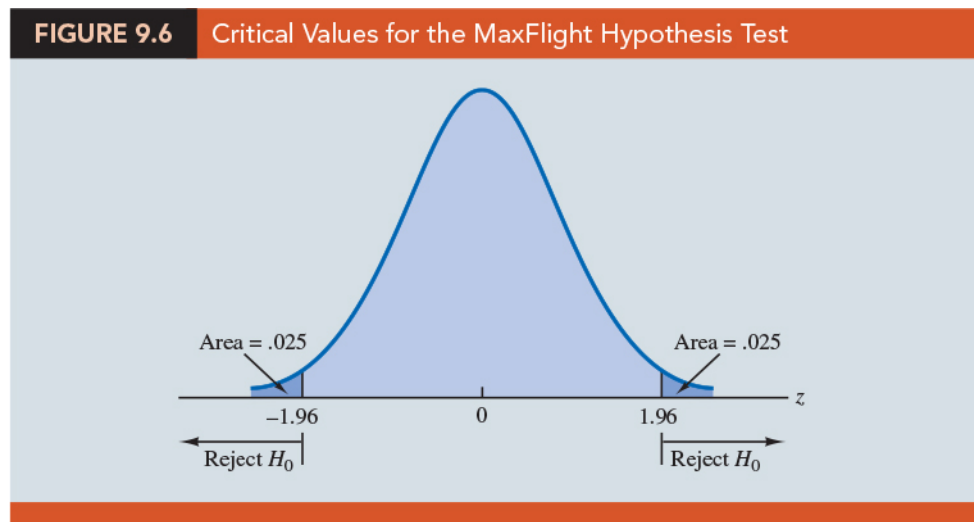
- (g) With a level of significance of $\alpha = 0.05$, we do not reject H_0 because the $p\text{-value} = 0.1260 > 0.05$. Because the null hypothesis is not rejected, no action will be taken to adjust the MaxFlight manufacturing process.

5. Computation of p -Values for Two-Tailed Tests.

- (a) Compute the value of the test statistic using equation (9.1).
- (b) If the value of the test statistic is in the upper tail, compute the probability that z is greater than or equal to the value of the test statistic (the upper tail area).
- (c) If the value of the test statistic is in the lower tail, compute the probability that z is less than or equal to the value of the test statistic (the lower tail area).
- (d) Double the probability (or tail area) from step (b) or (c) to obtain the p -value.

Critical value approach

1. (Figure 9.6) the critical values for the test will occur in both the lower and upper tails of the standard normal distribution. With a level of significance of $\alpha = 0.05$, the area in each tail corresponding to the critical values is $\alpha/2 = 0.05/2 = 0.025$.



2. The critical values for the test statistic are $-z_{0.025} = -1.96$ and $z_{0.025} = 1.96$.

3. The two-tailed rejection rule is

$$\text{Reject } H_0 \text{ if } z \leq -1.96 \text{ or if } z \geq 1.96$$

4. Because the value of the test statistic for the MaxFlight study is $z = 1.53$, the statistical evidence will not permit us to reject the null hypothesis at the 0.05 level of significance.

Summary and Practical Advice

1. Summary of the hypothesis testing procedures about a population mean for the σ known case. Note that μ_0 is the hypothesized value of the population mean.

	Lower Tail Test	Upper Tail Test	Two-Tailed Test
Hypotheses	$H_0: \mu \geq \mu_0$ $H_a: \mu < \mu_0$	$H_0: \mu \leq \mu_0$ $H_a: \mu > \mu_0$	$H_0: \mu = \mu_0$ $H_a: \mu \neq \mu_0$
Test Statistic	$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$	$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$	$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$
Rejection Rule: p-Value Approach	Reject H_0 if $p\text{-value} \leq \alpha$	Reject H_0 if $p\text{-value} \leq \alpha$	Reject H_0 if $p\text{-value} \leq \alpha$
Rejection Rule: Critical Value Approach	Reject H_0 if $z \leq -z_\alpha$	Reject H_0 if $z \geq z_\alpha$	Reject H_0 if $z \leq -z_{\alpha/2}$ or if $z \geq z_{\alpha/2}$

2. Steps of Hypothesis Testing

Step 1. Develop the null and alternative hypotheses (H_0, H_a).

Step 2. Specify the level of significance (α).

Step 3. Collect the sample data ($\mathbf{x} = \{x_1, x_2, \dots, x_n\}$) and compute the value of the test statistic ($T(\mathbf{x}) = z$).

p-value Approach

Step 4. Use the value of the test statistic to compute the p -value.

Step 5. Reject H_0 if the $p\text{-value} \leq \alpha$.

Step 6. Interpret the statistical conclusion in the context of the application.

Critical Value Approach

Step 4. Use α to determine the critical value (z_α or $z_{\alpha/2}$) and the rejection rule.

Step 5. Use the value of the test statistic and the rejection rule to determine whether to reject H_0 .

Step 6. Interpret the statistical conclusion in the context of the application.

3. Practical advice about the sample size for hypothesis tests is similar to the advice we provided about the sample size for interval estimation in Chapter 8.

(a) In most applications, a sample size of $n \geq 30$ is adequate when using the hypothesis testing procedure described in this section.

(b) If the population is normally distributed, the hypothesis testing procedure that we described is exact and can be used for any sample size.

(c) If the population is not normally distributed but is at least roughly symmetric, sample sizes as small as 15 can be expected to provide acceptable results.

Relationship Between Interval Estimation and Hypothesis Testing

1. (Recall, Chapter 8) For the σ known case, the $(1-\alpha)\%$ confidence interval estimate of a population mean is given by

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

2. (Recall, Chapter 9) a two-tailed hypothesis test about a population mean:

$$H_0 : \mu = \mu_0, \quad H_a : \mu \neq \mu_0$$

where μ_0 is the hypothesized value for the population mean.

3. Constructing a $100(1-\alpha)\%$ confidence interval for the population mean: $100(1-\alpha)\%$ of the confidence intervals generated will contain the population mean and $100\alpha\%$ of the confidence intervals generated will not contain the population mean.
4. If we reject H_0 whenever the confidence interval does not contain μ_0 , we will be rejecting H_0 when it is true ($\mu = \mu_0$) with probability α .
5. Recall that α is the probability of rejecting the null hypothesis when it is true.
6. So constructing a $100(1-\alpha)\%$ confidence interval and rejecting H_0 whenever the interval does not contain μ_0 is equivalent to conducting a two-tailed hypothesis test with α as the level of significance.
7. A Confidence Interval Approach to Testing a Hypothesis of the Form:

$$H_0 : \mu = \mu_0, \quad H_a : \mu \neq \mu_0$$

- (a) Select a simple random sample from the population and use the value of the sample mean \bar{x} to develop the confidence interval for the population mean μ .

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

- (b) If the confidence interval contains the hypothesized value μ_0 , do not reject H_0 . Otherwise, reject H_0 .
8. Note that this discussion and example pertain to two-tailed hypothesis tests about a population mean. However, the same confidence interval and two-tailed hypothesis testing relationship exists for other population parameters.
9. The relationship can also be extended to one-tailed tests about population parameters. Doing so, however, requires the development of one-sided confidence intervals, which are rarely used in practice.

 Question (p435)

The MaxFlight hypothesis test takes the following form:

$$H_0 : \mu = 295, \quad H_a : \mu \neq 295.$$

Conducting the MaxFlight hypothesis test with a level of significance of $\alpha = 0.05$ using the confidence interval approach.

sol:

- We sampled $n = 50$ golf balls and found a sample mean distance of $\bar{x} = 297.6$ yards. Recall that the population standard deviation is $\sigma = 12$.

- The 95% confidence interval estimate of the population mean is

$$\begin{aligned} & \frac{\bar{x} \pm z_{0.025} \frac{\sigma}{\sqrt{n}}}{297.6 \pm 1.96 \frac{12}{\sqrt{50}}} \\ & 297.6 \pm 3.3 \quad \text{or} \quad (294.3, 300.9). \end{aligned}$$

- With 95% confidence, the mean distance for the population of golf balls is between 294.3 and 300.9 yards.
- Because the interval contains the hypothesized value for the population mean, $\mu_0 = 295$, the hypothesis testing conclusion is that the null hypothesis, $H_0 : \mu = 295$, cannot be rejected.

9.4 Population Mean: σ Unknown

1. To conduct a hypothesis test about a population mean for the σ unknown case, the sample mean \bar{x} is used as an estimate of μ and the sample standard deviation s is used as an estimate of σ .

2. (Recall) For the σ known case, the sampling distribution of the test statistic has a standard normal distribution. For the σ unknown case, however, the sampling distribution of the test statistic follows the t distribution; it has slightly more variability because the sample is used to develop estimates of both μ and σ .
3. **Test Statistic for Hypothesis Tests about a Population Mean: σ Unknown**
the test statistic has a t distribution with $n-1$ degrees of freedom:

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \quad (9.2)$$

4. The t distribution is based on an assumption that the population from which we are sampling has a normal distribution. However, research shows that this assumption can be relaxed considerably when the sample size is large enough.

One-tailed Test

1. **Example** A business travel magazine wants to classify transatlantic gateway airports according to the mean rating for the population of business travelers. A rating scale with a low score of 0 and a high score of 10 will be used, and airports with a population mean rating greater than 7 will be designated as superior service airports. The magazine staff surveyed a sample of 60 business travelers at each airport to obtain the ratings data. The sample for London's Heathrow Airport provided a sample mean rating of $\bar{x} = 7.25$ and a sample standard deviation of $s = 1.052$. Do the data indicate that Heathrow should be designated as a superior service airport?
- (a) We want to develop a hypothesis test for which the decision to reject H_0 will lead to the conclusion that the population mean rating for the Heathrow Airport is greater than 7.
- (b) The null and alternative hypotheses for this upper tail test:

$$H_0 : \mu \leq 7, \quad H_a : \mu > 7$$

- (c) Use $\alpha = 0.05$, with $\bar{x} = 7.25$, $\mu_0 = 7$, $s = 1.052$, and $n = 60$, the value of the test statistic:

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{7.25 - 7}{1.052/\sqrt{60}} = 1.84$$

- (d) The sampling distribution of t has $n-1 = 60-1 = 59$ degrees of freedom. Because the test is an upper tail test, the p -value is $P(t \geq 1.84)$, that is, the upper tail area corresponding to the value of the test statistic.
- (e) (Table 2 in Appendix B) the t distribution with 59 degrees of freedom provides the following information.

Area in Upper Tail	0.20	0.10	0.05	0.025	0.01	0.005
t -Value (59 df)	0.848	1.296	1.671	2.001	2.391	2.662

- i. We see that $t = 1.84$ is between 1.671 and 2.001. The values in the "Area in Upper Tail" row show that the p -value must be less than 0.05 and greater than 0.025.
- ii. With a level of significance of $\alpha = 0.05$, this placement is all we need to know to make the decision to reject the null hypothesis and conclude that Heathrow should be classified as a superior service airport.
- (f) (Using software) $t = 1.84$ provides the upper tail p -value of 0.0354 for the Heathrow Airport hypothesis test. With $0.0354 < 0.05$, we reject the null hypothesis and conclude that Heathrow should be classified as a superior service airport.
- (g) The critical value corresponding to an area of $\alpha = 0.05$ in the upper tail of a t distribution with 59 degrees of freedom is $t_{59;0.05} = 1.671$.
- (h) The rejection rule using the critical value approach is to reject H_0 if $t \geq 1.671$. Because $t = 1.84 > 1.671$, H_0 is rejected. Heathrow should be classified as a superior service airport.

Two-tailed Test

- Example** Consider the hypothesis testing situation facing Holiday Toys. The company manufactures and distributes its products through more than 1000 retail outlets. In planning production levels for the coming winter season, Holiday must decide how many units of each product to produce prior to knowing the actual demand at the retail level. For this year's most important new toy, Holiday's marketing director is expecting demand to average 40 units per retail outlet. Prior to making the final production decision based upon this estimate, Holiday decided

to survey a sample of 25 retailers in order to develop more information about the demand for the new product. Each retailer was provided with information about the features of the new toy along with the cost and the suggested selling price. Then each retailer was asked to specify an anticipated order quantity. With μ denoting the population mean order quantity per retail outlet, the sample data will be used to conduct the following two-tailed hypothesis test:

$$\underline{H_0 : \mu = 40, \quad H_a : \mu \neq 40}$$

- (a) If H_0 cannot be rejected, Holiday will continue its production planning based on the marketing director's estimate that the population mean order quantity per retail outlet will be $\mu = 40$ units.
- (b) If H_0 is rejected, Holiday will immediately reevaluate its production plan for the product.
- (c) A two-tailed hypothesis test is used because Holiday wants to reevaluate the production plan if the population mean quantity per retail outlet is less than anticipated or greater than anticipated.
- (d) Because no historical data are available (it's a new product), the population mean μ and the population standard deviation must both be estimated using \bar{x} and s from the sample data.
- (e) The sample of 25 retailers provided a mean of $\bar{x} = 37.4$ and a standard deviation of $s = 11.79$ units.
- (f) (*Check on the form of the population distribution*). The histogram of the sample data showed no evidence of skewness or any extreme outliers, so the analyst concluded that the use of the t distribution with $n-1 = 24$ degrees of freedom was appropriate.
- (g) Using equation (9.2) with $\bar{x} = 37.4$, $\mu_0 = 40$, $s = 11.79$, and $n = 25$, the value of the test statistic is

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{37.4 - 40}{11.79/\sqrt{25}} = -1.10 .$$

- (h) The t distribution table only contains positive t values. Because the t distribution is symmetric, however, the upper tail area at $t = 1.10$ is the same as the lower tail area at $t = -1.10$.

(i) (Table 2 in Appendix B)

Area in Upper Tail	0.20	0.10	0.05	0.025	0.01	0.005
t-Value (24 <i>df</i>)	0.857	1.318	1.711	2.064	2.492	2.797

(j) We see that $t = 1.10$ is between 0.857 and 1.318. From the "Area in Upper Tail" row, we see that the area in the upper tail at $t = 1.10$ is between 0.20 and 0.10.

(k) When we double these amounts, we see that the p -value must be between 0.40 and 0.20. With a level of significance of $\alpha = 0.05$, we now know that the p -value is greater than α . Therefore, H_0 cannot be rejected. Sufficient evidence is not available to conclude that Holiday should change its production plan for the coming season.

(l) (Software) The p -value obtained is 0.2822. With a level of significance of $\alpha = 0.05$, we cannot reject H_0 because $0.2822 > 0.05$.

(m) With $\alpha = 0.05$ and the t distribution with 24 degrees of freedom, $-t_{24,0.025} = -2.064$ and $t_{24,0.025} = 2.064$ are the critical values for the two-tailed test. The rejection rule using the test statistic is

$$\underline{\text{Reject } H_0 \text{ if } t \leq -2.064 \text{ or if } t \geq 2.064}.$$

(n) Based on the test statistic $t = -1.10$, H_0 cannot be rejected. This result indicates that Holiday should continue its production planning for the coming season based on the expectation that $\mu = 40$.

Summary and Practical Advice

- (Table 9.3) A summary of the hypothesis testing procedures about a population mean for the σ unknown case.

TABLE 9.3 Summary of Hypothesis Tests About a Population Mean: σ Unknown Case

	Lower Tail Test	Upper Tail Test	Two-Tailed Test
Hypotheses	$H_0: \mu \geq \mu_0$ $H_a: \mu < \mu_0$	$H_0: \mu \leq \mu_0$ $H_a: \mu > \mu_0$	$H_0: \mu = \mu_0$ $H_a: \mu \neq \mu_0$
Test Statistic	$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$	$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$	$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$
Rejection Rule: p-Value Approach	Reject H_0 if $p\text{-value} \leq \alpha$	Reject H_0 if $p\text{-value} \leq \alpha$	Reject H_0 if $p\text{-value} \leq \alpha$
Rejection Rule: Critical Value Approach	Reject H_0 if $t \leq -t_\alpha$	Reject H_0 if $t \geq t_\alpha$	Reject H_0 if $t \leq -t_{\alpha/2}$ or if $t \geq t_{\alpha/2}$

2. The applicability of the hypothesis testing procedures of this section is dependent on the distribution of the population being sampled from and the sample size.

9.5 Population Proportion

1. Using p_0 to denote the hypothesized value for the population proportion, the three forms for a hypothesis test about a population proportion:

lower tail test	upper tail test	two-tailed test
$H_0 : p \geq p_0$	$H_0 : p \leq p_0$	$H_0 : p = p_0$
$H_a : p < p_0$	$H_a : p > p_0$	$H_a : p \neq p_0$

2. Hypothesis tests about a population proportion are based on the difference between the sample proportion \bar{p} and the hypothesized population proportion p_0 .
3. The sampling distribution of \bar{p} , the point estimator of the population parameter p , is the basis for developing the test statistic.

4. When the null hypothesis is true as an equality, the expected value of \bar{p} equals the hypothesized value p_0 ; that is, $E(\bar{p}) = p_0$. The standard error of \bar{p} is given by

$$\sigma_{\bar{p}} = \sqrt{\frac{p_0(1-p_0)}{n}}$$

5. (Recall, Chapter 7) we said that if $np \geq 5$ and $n(1-p) \geq 5$, the sampling distribution of \bar{p} can be approximated by a normal distribution. Under these conditions, which usually apply in practice, the quantity

$$z = \frac{\bar{p} - p_0}{\sigma_{\bar{p}}} \quad (9.3)$$

has a standard normal probability distribution.

6. Test Statistic for Hypothesis Tests About a Population Proportion

With $\sigma_{\bar{p}} = \sqrt{p_0(1-p_0)/n}$, the standard normal random variable z is the test statistic used to conduct hypothesis tests about a population proportion.

$$z = \frac{\bar{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

7. **Example** (Pine Creek golf course example). Over the past year, 20% of the players at Pine Creek were women. In an effort to increase the proportion of women players, Pine Creek implemented a special promotion designed to attract women golfers. One month after the promotion was implemented, the course manager requested a statistical study to determine whether the proportion of women players at Pine Creek had increased.

- (a) Because the objective of the study is to determine whether the proportion of women golfers increased, an upper tail test with $H_a : p > 0.20$ is appropriate:

$$H_0 : p \leq 0.20, \quad H_a : p > 0.20$$

- (b) If H_0 can be rejected, the test results will give statistical support for the conclusion that the proportion of women golfers increased and the promotion was beneficial.

- (c) The course manager specified that a level of significance of $\alpha = 0.05$ be used in carrying out this hypothesis test.
- (d) The next step of the hypothesis testing procedure is to select a sample and compute the value of an appropriate test statistic. Suppose a random sample of $n = 400$ players was selected, and that $x = 100$ of the players were women. The proportion of women golfers in the sample is

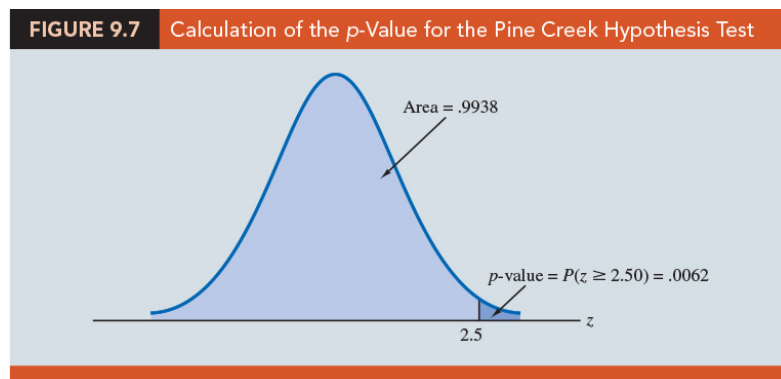
$$\bar{p} = \frac{100}{400} = 0.25$$

Using equation (9.4), the value of the test statistic is

$$z = \frac{\bar{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.25 - 0.20}{\sqrt{\frac{0.20(1-0.20)}{400}}} = \frac{0.05}{0.02} = 2.50.$$

(e) *The p-value approach.*

- i. The p -value is the probability that z is greater than or equal to $z = 2.50$. $P(z \geq 2.50) = 0.9938$, the p -value for the Pine Creek test is $1.0000 - 0.9938 = 0.0062$.
- ii. (Figure 9.7) Recall that the course manager specified a level of significance of $\alpha = 0.05$. A $p\text{-value} = 0.0062 < 0.05$ gives sufficient statistical evidence to reject H_0 at the 0.05 level of significance.
- iii. The test provides statistical support for the conclusion that the special promotion increased the proportion of women players at the Pine Creek golf course.



- (f) *The critical value approach.* The critical value corresponding to an area of 0.05 in the upper tail of a normal probability distribution is $z_{0.05} = 1.645$.

Thus, the rejection rule using the critical value approach is to reject H_0 if $z \geq 1.645$. Because $z = 2.50 > 1.645$, H_0 is rejected.

- (g) The p -value approach provides more information. With a p -value = 0.0062, the null hypothesis would be rejected for any level of significance greater than or equal to 0.0062.

Summary

- The procedure used to conduct a hypothesis test about a population proportion is similar to the procedure used to conduct a hypothesis test about a population mean.
- Although we only illustrated how to conduct a hypothesis test about a population proportion for an upper tail test, similar procedures can be used for lower tail and two-tailed tests.
- (Table 9.4) A summary of the hypothesis tests about a population proportion. We assume that $np \geq 5$ and $n(1-p) \geq 5$; thus the normal probability distribution can be used to approximate the sampling distribution of \bar{p} .

TABLE 9.4 Summary of Hypothesis Tests About a Population Proportion			
	Lower Tail Test	Upper Tail Test	Two-Tailed Test
Hypotheses	$H_0: p \geq p_0$ $H_a: p < p_0$	$H_0: p \leq p_0$ $H_a: p > p_0$	$H_0: p = p_0$ $H_a: p \neq p_0$
Test Statistic	$z = \frac{\bar{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$	$z = \frac{\bar{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$	$z = \frac{\bar{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$
Rejection Rule: p-Value Approach	Reject H_0 if $p\text{-value} \leq \alpha$	Reject H_0 if $p\text{-value} \leq \alpha$	Reject H_0 if $p\text{-value} \leq \alpha$
Rejection Rule: Critical Value Approach	Reject H_0 if $z \leq -z_\alpha$	Reject H_0 if $z \geq z_\alpha$	Reject H_0 if $z \leq -z_{\alpha/2}$ or if $z \geq z_{\alpha/2}$

9.6 Hypothesis Testing and Decision Making

1. The hypothesis testing applications are considered as significance tests :
 - (a) formulate the null and alternative hypotheses, H_0, H_a .
 - (b) specify the level of significance, α .
 - (c) select a sample, $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$.
 - (d) compute the value of a test statistic, $T(\mathbf{x})$.
 - (e) compute the associated p -value.
 - (f) compare the p -value to α .
 - (g) conclude "reject H_0 " and declare the results significant if $p\text{-value} \leq \alpha$; otherwise, we made the conclusion "do not reject H_0 ."
2. With a significance test, we control the probability of making the Type I error, but not the Type II error. Thus, we recommended the conclusion "do not reject H_0 " rather than "accept H_0 " because the latter puts us at risk of making the Type II error of accepting H_0 when it is false.
3. With the conclusion "do not reject H_0 ," the statistical evidence is considered inconclusive and is usually an indication to postpone a decision or action until further research and testing can be undertaken.
4. If the purpose of a hypothesis test is to make a decision when H_0 is true and a different decision when H_a is true, the decision maker may want to, and in some cases be forced to, take action with both the conclusion do not reject H_0 and the conclusion reject H_0 . If this situation occurs, statisticians generally recommend controlling the probability of making a Type II error .
5. With the probabilities of both the Type I and Type II error controlled, the conclusion from the hypothesis test is either to accept H_0 or reject H_0 .
6. Example (lot-acceptance example) A quality control manager must decide to accept a shipment of batteries from a supplier or to return the shipment because of poor quality.

- (a) Assume that design specifications require batteries from the supplier to have a mean useful life of at least 120 hours. To evaluate the quality of an incoming shipment, a sample of 36 batteries will be selected and tested.
- (b) On the basis of the sample, a decision must be made to accept the shipment of batteries or to return it to the supplier because of poor quality.
- (c) Let μ denote the mean number of hours of useful life for batteries in the shipment. The null and alternative hypotheses about the population mean:

$$\underline{H_0 : \mu \geq 120, \quad H_a : \mu < 120}$$

- i. If H_0 is rejected, the alternative hypothesis is concluded to be true. This conclusion indicates that the appropriate action is to return the shipment to the supplier.
- ii. If H_0 is not rejected, the decision maker must still determine what action should be taken. Thus, without directly concluding that H_0 is true, but merely by not rejecting it, the decision maker will have made the decision to accept the shipment as being of satisfactory quality.
- (d) In such decision-making situations, it is recommended that the hypothesis testing procedure be extended to control the probability of making a Type II error.
7. Because a decision will be made and action taken when we do not reject H_0 , knowledge of the probability of making a Type II error will be helpful.

9.7 Calculating The Probability of Type II Errors

1. **Example** (lot-acceptance example) The null and alternative hypotheses about the mean number of hours of useful life for a shipment of batteries:

$$H_0 : \mu \geq 120, \quad H_a : \mu < 120.$$

- (a) If H_0 is rejected, the decision will be to return the shipment to the supplier because the mean hours of useful life are less than the specified 120 hours.
- (b) If H_0 is not rejected, the decision will be to accept the shipment.

2. Suppose a level of significance of $\alpha = 0.05$, the test statistic in the σ known case:

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{\bar{x} - 120}{\sigma/\sqrt{n}}$$

3. The rejection rule for the lower tail test:

$$\text{Reject } H_0 \text{ if } z \leq -1.645$$

4. Suppose a sample of $n = 36$ batteries will be selected and based upon previous testing the population standard deviation can be assumed known with a value of $\sigma = 12$ hours.

5. The rejection rule indicates that we will reject H_0 if

$$z = \frac{\bar{x} - 120}{12/\sqrt{36}} \leq -z_{0.05} = -1.645$$

6. Solving for \bar{x} in the preceding expression indicates that we will reject H_0 if

$$\bar{x} \leq 120 - 1.645 \left(\frac{12}{\sqrt{36}} \right) = 116.71$$

7. Rejecting H_0 when $\bar{x} \leq 116.71$ means that we will make the decision to accept the shipment whenever $\bar{x} > 116.71$.

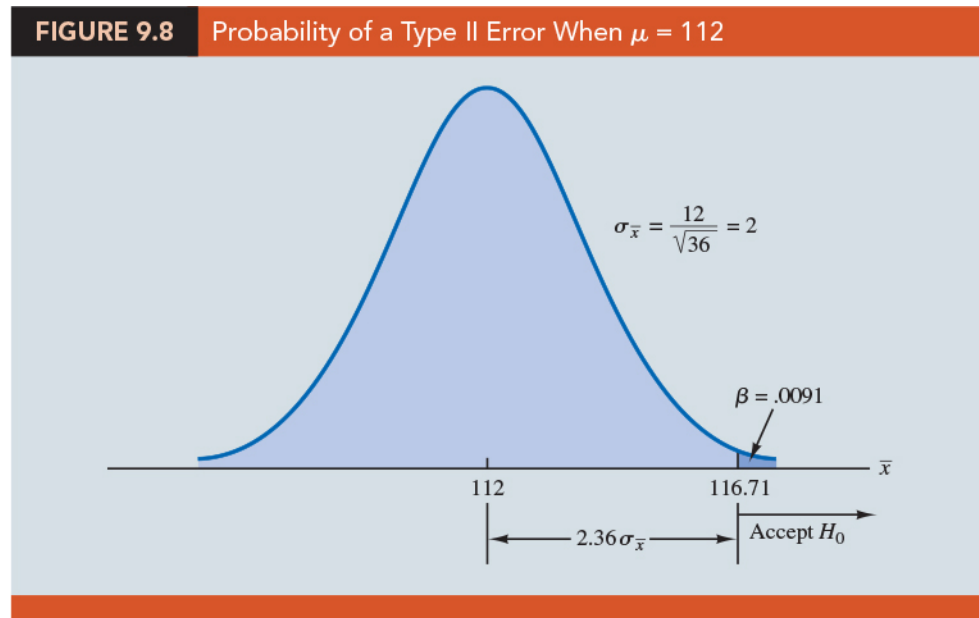
8. Compute probabilities associated with making a Type II error.

- (a) (Recall) we make a Type II error whenever the true shipment mean is less than 120 hours and we make the decision to accept $H_0 : \mu \geq 120$.
- (b) Hence, to compute the probability of making a Type II error, we must select a value of μ less than 120 hours.
- (c) For example, suppose the shipment is considered to be of poor quality if the batteries have a mean life of $\mu = 112$ hours.

- (d) If $\mu = 112$ is really true, what is the probability of accepting $H_0 : \mu \geq 120$ and hence committing a Type II error?

$$\underline{P(\bar{x} \geq 116.71), \text{ when } \mu = 112} .$$

- (e) (Figure 9.8) the sampling distribution of \bar{x} when the mean is $\mu = 112$. The shaded area in the upper tail gives the probability of obtaining $\underline{\bar{x} > 116.71}$.



- (f) Using the standard normal distribution, we see that at $\bar{x} = 116.71$

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{116.71 - 112}{12/\sqrt{36}} = 2.36$$

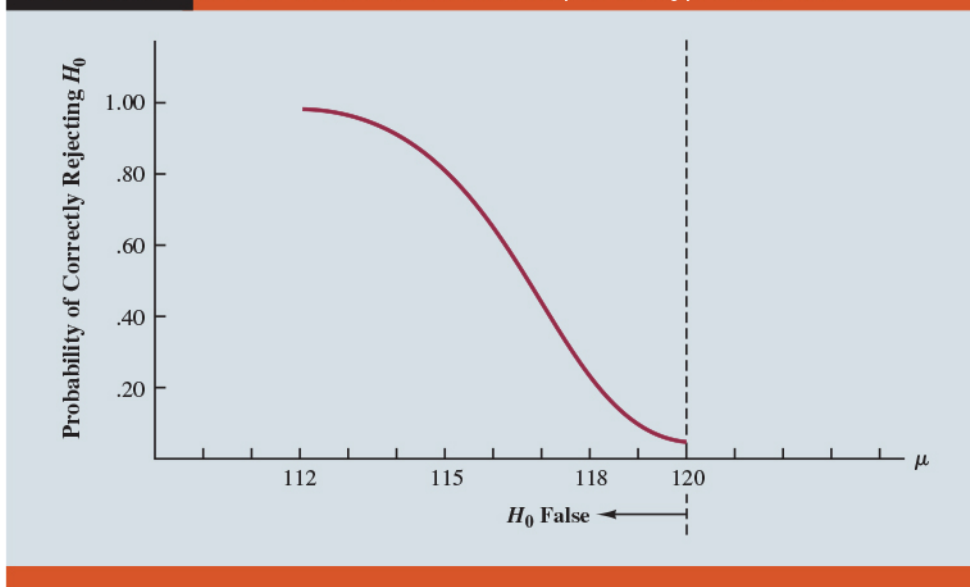
- (g) The probability of making a Type II error when $\mu = 112$ is $\underline{\beta := P(z \geq 2.36) = 0.0091}$.
- (h) Therefore, we can conclude that if the mean of the population is 112 hours, the probability of making a Type II error is only 0.0091.
- (i) We can repeat these calculations for other values of μ less than 120.

9. (Table 9.5) we show the probability of making a Type II error for a variety of values of μ less than 120. Note that as μ increases toward 120, the probability of making a Type II error increases toward an upper bound of 0.95. However, as μ decreases to values farther below 120, the probability of making a Type II error diminishes.

TABLE 9.5 Probability of Making a Type II Error for the Lot-Acceptance Hypothesis Test

Value of μ	$z = \frac{116.71 - \mu}{12/\sqrt{36}}$	Probability of a Type II Error (β)	Power ($1 - \beta$)
112	2.36	.0091	.9909
114	1.36	.0869	.9131
115	.86	.1949	.8051
116.71	.00	.5000	.5000
117	-.15	.5596	.4404
118	-.65	.7422	.2578
119.999	-1.645	.9500	.0500

10. When the true population mean μ is close to (far below) the null hypothesis value of $\mu = 120$, the probability is high (low) that we will make a Type II error.
11. For any particular value of μ , the power is $1 - \beta$; that is, the probability of correctly rejecting H_0 is 1 minus the probability of making a Type II error.
12. (Figure 9.9) *Power curve*: the power associated with each value of μ :

FIGURE 9.9 Power Curve for the Lot-Acceptance Hypothesis Test

- (a) Note that the power curve extends over the values of μ for which the H_0 is false.

- (b) The height of the power curve at any value of μ indicates the probability of correctly rejecting H_0 when H_0 is false.
13. **The step-by-step procedure to compute the probability of making a Type II error in hypothesis tests about a population mean**
- (a) Formulate the null and alternative hypotheses.
- (b) Use the level of significance α and the critical value approach to determine the critical value and the rejection rule for the test.
- (c) Use the rejection rule to solve for the value of the sample mean corresponding to the critical value of the test statistic.
- (d) Use the results from step (c) to state the values of the sample mean that lead to the acceptance of H_0 . These values define the acceptance region for the test.
- (e) Use the sampling distribution of \bar{x} for a value of μ satisfying H_a , and the acceptance region from step (d), to compute the probability that the sample mean will be in the acceptance region.
- (f) This probability is the probability of making a Type II error at the chosen value of μ .

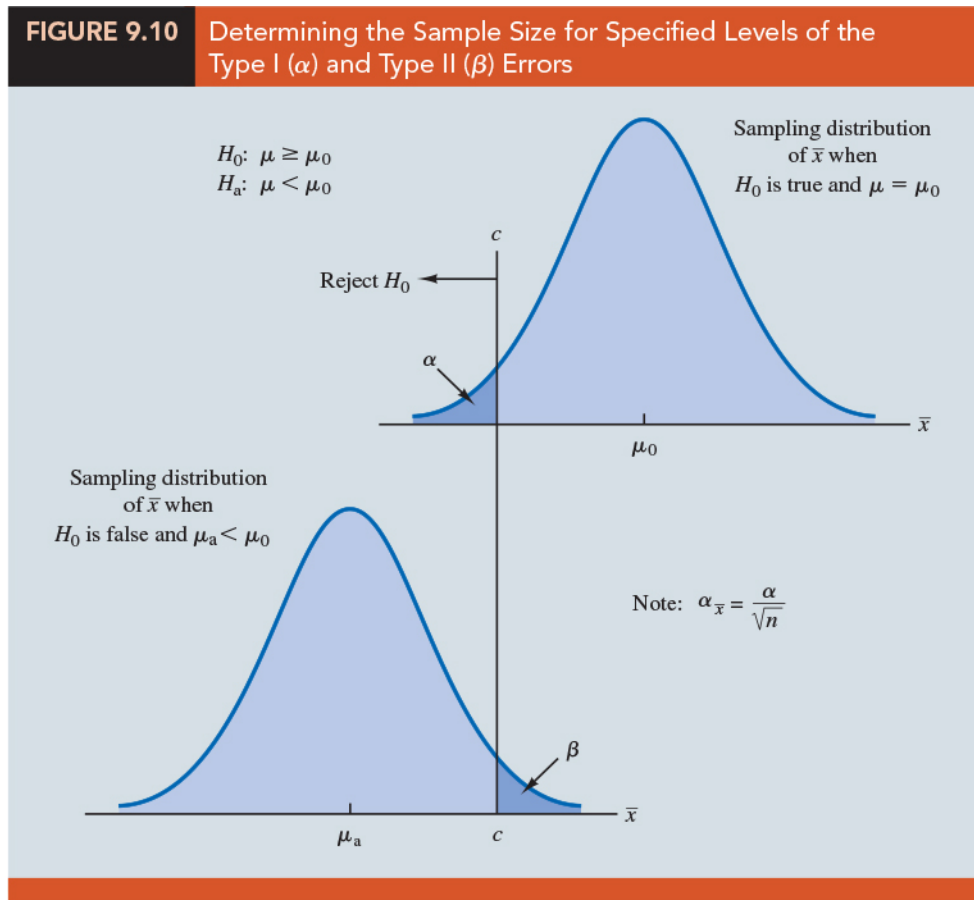
9.8 Determining The Sample Size for a Hypothesis Test about a Population Mean

1. Assume that a hypothesis test is to be conducted about the value of a population mean. The level of significance specified by the user determines the probability of making a Type I error for the test. By controlling the sample size, the user can also control the probability of making a Type II error.

2. Let us show how a sample size can be determined for the following lower tail test about a population mean.

$$H_0 : \mu \geq \mu_0, \quad H_a : \mu < \mu_0$$

3. (Figure 9.10) The upper panel is the sampling distribution of \bar{x} when H_0 is true with $\mu = \mu_0$.



4. For a lower tail test, the critical value of the test statistic is denoted $-z_\alpha$. In the upper panel of the figure the vertical line, labeled c , is the corresponding value of \bar{x} .
5. If we reject H_0 when $\bar{x} \leq c$, the probability of a Type I error will be α : $P(\bar{x} \leq c \mid H_0 T) = \alpha$.
6. With z_α representing the z value corresponding to an area of α in the upper tail of the standard normal distribution, we compute c using the following formula:

$$c = \mu_0 - z_\alpha \frac{\sigma}{\sqrt{n}}$$

7. (Figure 9.10) The lower panel is the sampling distribution of \bar{x} when the alternative hypothesis is true with $\mu = \mu_a < \mu_0$. The shaded region shows β , the probability of a Type II error that the decision maker will be exposed to if the null hypothesis is accepted when $\bar{x} > c$.
8. With z_β representing the z value corresponding to an area of β in the upper tail of the standard normal distribution, we compute c using the following formula:

$$c = \mu_a + z_\beta \frac{\sigma}{\sqrt{n}} \quad (9.6)$$

9. Now what we want to do is to select a value for c so that when we reject H_0 and accept H_a , the probability of a Type I error is equal to the chosen value of α and the probability of a Type II error is equal to the chosen value of β . Therefore, both equations (9.5) and (9.6) must provide the same value for c :

$$\begin{aligned} \mu_0 - z_\alpha \frac{\sigma}{\sqrt{n}} &= \mu_a + z_\beta \frac{\sigma}{\sqrt{n}} \\ \mu_0 - \mu_a &= z_\alpha \frac{\sigma}{\sqrt{n}} + z_\beta \frac{\sigma}{\sqrt{n}} \\ \mu_0 - \mu_a &= \frac{(z_\alpha + z_\beta)\sigma}{\sqrt{n}} \\ \sqrt{n} &= \frac{(z_\alpha + z_\beta)\sigma}{(\mu_0 - \mu_a)} \end{aligned}$$

10. Sample Size for a One-Tailed Hypothesis Test About a Population Mean

$$n = \frac{(z_\alpha + z_\beta)^2 \sigma^2}{(\mu_0 - \mu_a)^2} \quad (9.7)$$

where

- $z_\alpha = z$ value providing an area of α in the upper tail of a standard normal distribution.
- $z_\beta = z$ value providing an area of β in the upper tail of a standard normal distribution.
- $\sigma =$ the population standard deviation.
- $\mu_0 =$ the value of the population mean in the null hypothesis.
- $\mu_a =$ the value of the population mean used for the Type II error.

Note: In a two-tailed hypothesis test, use (9.7) with $z_{\alpha/2}$ replacing z_{α} .

11. **Example** lot-acceptance example

- (a) The design specification for the shipment of batteries indicated a mean useful life of at least 120 hours for the batteries. Shipments were rejected if $H_0 : \mu \geq 120$ was rejected.
- (b) Let us assume that the quality control manager makes the following statements about the allowable probabilities for the Type I and Type II errors.
- Type I error statement:* If the mean life of the batteries in the shipment is $\mu = 120$, I am willing to risk an $\alpha = 0.05$ probability of rejecting the shipment.
 - Type II error statement:* If the mean life of the batteries in the shipment is 5 hours under the specification (i.e., $\mu = 115$), I am willing to risk a $\beta = 0.10$ probability of accepting the shipment.
 - These statements are based on the judgment of the manager. Someone else might specify different restrictions on the probabilities. However, statements about the allowable probabilities of both errors must be made before the sample size can be determined.
- (c) In the example, $\alpha = 0.05$ and $\beta = 0.10$. Using the standard normal probability distribution, we have $z_{0.05} = 1.645$ and $z_{0.10} = 1.28$. From the statements about the error probabilities, we note that $\mu_0 = 120$ and $\mu_a = 115$. Finally, the population standard deviation was assumed known at $\sigma = 12$.

- (d) The recommended sample size for the lot-acceptance example is

$$n = \frac{(z_{\alpha} + z_{\beta})^2 \sigma^2}{(\mu_0 - \mu_a)^2} = \frac{(1.645 + 1.28)^2 (12)^2}{(120 - 115)^2} = 49.3$$

Rounding up, we recommend a sample size of 50.

- (e) Because both the Type I and Type II error probabilities have been **controlled** at allowable levels with $n = 50$, the quality control manager is now justified in using the accept H_0 and reject H_0 statements for the hypothesis test.

12. We can make three observations about the relationship among α , β , and the sample size n .

- (a) Once two of the three values are known, the other can be computed.
- (b) For a given α , increasing n will reduce β .
- (c) For a given n , decreasing α will increase β , whereas increasing α will decrease β .
13. The third observation should be kept in mind when the probability of a Type II error is not being controlled. It suggests that one should not choose unnecessarily small values for the level of significance α .
14. For a given sample size, choosing a smaller level of significance means more exposure to a Type II error. Inexperienced users of hypothesis testing often think that smaller values of α are always better. They are better if we are concerned only about making a Type I error. However, smaller values of α have the disadvantage of increasing the probability of making a Type II error.

9.9 Big Data And Hypothesis Testing*

☺ EXERCISES

9.1 : 2, 3

9.2 : 6, 7

9.3 : 9, 11, 15, 18, 22

9.4 : 23, 24, 27, 33

9.5 : 35, 36, 43, 44

9.7 : 46, 50, 53

9.8 : 55, 59

SUP : 69, 79, 83

“成功是歷經一個又一個的失敗卻不失去熱情。”

“Success is the ability to go from failure to failure without losing your enthusiasm.”

— *Winston Churchill (November 30, 1874 – January 24, 1965)*

統計學 (一)

Anderson's Statistics for Business & Economics (14/E)

Chapter 10: Inference About Means and Proportions with Two Populations

上課時間地點: 四 D56, 研究 250105

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

系級: _____ 學號: _____ 姓名: _____

Overview

1. Discuss the statistical inference (interval estimates and hypothesis tests) for two population means (three situations: population standard deviations known, unknown; match samples) and the two population proportions.
2. *Examples:*
 - (a) Develop an interval estimate of the difference between the mean starting salary for a population of men and the mean starting salary for a population of women.
 - (b) Conduct a hypothesis test to determine whether any difference is present between the proportion of defective parts in a population of parts produced by supplier *A* and the proportion of defective parts in a population of parts produced by supplier *B*.

10.1 Inferences About the Difference Between Two Population Means: σ_1 and σ_2 Known

1. μ_1 (μ_2) denote the mean of population 1 (2), we will focus on inferences about the difference between the means: $\mu_1 - \mu_2$.
2. A simple random sample of n_1 (n_2) units from population 1 (2). The two samples, taken separately and independently, are referred to as independent simple random samples.
3. Assume the two population standard deviations, σ_1 and σ_2 , can be assumed known prior to collecting the samples.
4. Question: how to compute a margin of error and develop an interval estimate of the difference between the two population means when σ_1 and σ_2 are known.

Interval Estimation of $\mu_1 - \mu_2$

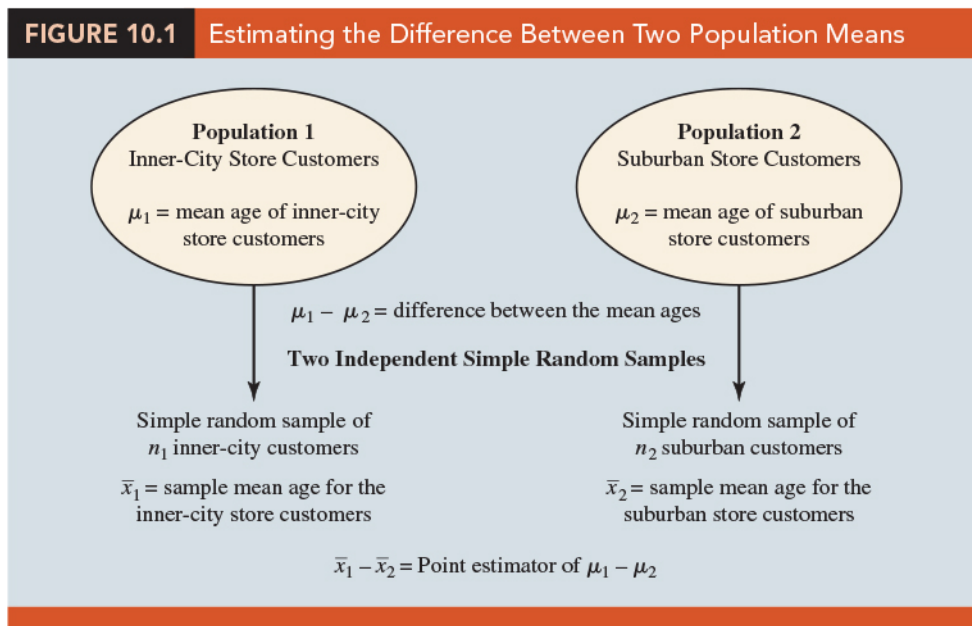
1. **Example** Greystone Department Stores, Inc., operates two stores in Buffalo, New York: One is in the inner city and the other is in a suburban shopping center. The regional manager noticed that products that sell well in one store do not always sell well in the other. The manager believes this situation may be attributable to differences in customer demographics at the two locations. Customers may differ in age, education, income, and so on. (對觀察事物提出問題)
2. Suppose the manager asks us to investigate the difference between the mean ages of the customers who shop at the two stores. (針對問題收集資料)
3. Let us define population 1 as all customers who shop at the inner-city store and population 2 as all customers who shop at the suburban store.
 - (a) μ_1 : mean of population 1 (i.e., the mean age of all customers who shop at the inner-city store)
 - (b) μ_2 : mean of population 2 (i.e., the mean age of all customers who shop at the suburban store)

4. The difference between the two population means is $\mu_1 - \mu_2$. To estimate $\mu_1 - \mu_2$, we will select a simple random sample of n_1 customers from population 1 and a simple random sample of n_2 customers from population 2.
5. We then compute the two sample means.
 - (a) \bar{x}_1 : sample mean age for the simple random sample of n_1 inner-city customers
 - (b) \bar{x}_2 : sample mean age for the simple random sample of n_2 suburban customers
6. The point estimator of the difference between the two population means is the difference between the two sample means.

7. Point Estimator of the Difference Between Two Population Means

$$\hat{\mu}_1 - \hat{\mu}_2 = \bar{x}_1 - \bar{x}_2 \quad (10.1)$$

8. (Figure 10.1) the process used to estimate the difference between two population means based on two independent simple random samples.



9. The point estimator $\bar{x}_1 - \bar{x}_2$ has a standard error that describes the variation in the sampling distribution of the estimator.

10. **Standard Error of $\bar{x}_1 - \bar{x}_2$** With two independent simple random samples, the standard error of $\bar{x}_1 - \bar{x}_2$ is as follows:

$$\sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad (10.2)$$

(証明如下:) (Hint: $Var(aX + bY) = a^2Var(X) + b^2Var(Y) + 2abCov(X, Y)$).

11. If both populations have a normal distribution, or if the sample sizes are large enough that the central limit theorem enables us to conclude that the sampling distributions of \bar{x}_1 and \bar{x}_2 can be approximated by a normal distribution, the sampling distribution of $\bar{x}_1 - \bar{x}_2$ will have a normal distribution with mean given by $\mu_1 - \mu_2$. (Denoted by $\bar{X}_1 - \bar{X}_2 \sim N(\mu_1 - \mu_2, \sigma_{\bar{x}_1 - \bar{x}_2}^2)$)
12. In general, an interval estimate is given by a point estimate \pm a margin of error. In the case of estimation of the difference between two population means, an interval estimate will take the following form:

$$\underline{(\bar{x}_1 - \bar{x}_2) \pm \text{Margin of error}}$$

13. With the sampling distribution of $\bar{x}_1 - \bar{x}_2$ having a normal distribution, we can write the margin of error as follows:


$$\text{Margin of error} = \underline{z_{\alpha/2} \sigma_{\bar{x}_1 - \bar{x}_2}} = \underline{z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad (10.3)$$

14. **Interval Estimate of the Difference Between Two Population Means: σ_1 and σ_2 Known**

$$\underline{(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad (10.4)$$

where $1 - \alpha$ is the confidence coefficient.

(公式說明如下:)

 **Question** (p485)

Example Greystone example. Based on data from previous customer demographic studies, the two population standard deviations are known with $\sigma_1 = 9$ years and $\sigma_2 = 10$ years. The data collected from the two independent simple random samples of Greystone customers provided the following results.

	Inner City Store	Suburban Store
Sample Size	$n_1 = 36$	$n_2 = 49$
Sample Mean	$\bar{x}_1 = 40$ years	$\bar{x}_2 = 35$ years

Find the margin of error and the 95% confidence interval estimate of the difference between the two population means.

sol:

Hypothesis Tests About $\mu_1 - \mu_2$

- Let us consider hypothesis tests about the difference between two population means. Using D_0 to denote the hypothesized difference between μ_1 and μ_2 , the three

forms for a hypothesis test are as follows:

Left-tailed test	Right-tailed test	Two-tailed test
$H_0 : \underline{\mu_1 - \mu_2 \geq D_0}$	$H_0 : \underline{\mu_1 - \mu_2 \leq D_0}$	$H_0 : \underline{\mu_1 - \mu_2 = D_0}$
$H_a : \underline{\mu_1 - \mu_2 < D_0}$	$H_a : \underline{\mu_1 - \mu_2 > D_0}$	$H_a : \underline{\mu_1 - \mu_2 \neq D_0}$

- In many applications, $\underline{D_0 = 0}$. Using the two-tailed test as an example, when $D_0 = 0$ the null hypothesis is $H_0 : \mu_1 - \mu_2 = 0$.
- In this case, the null hypothesis is that μ_1 and μ_2 are equal. Rejection of H_0 leads to the conclusion that $H_a : \mu_1 - \mu_2 \neq 0$ is true; that is, μ_1 and μ_2 are not equal.
- The general steps for conducting hypothesis tests: choose a level of significance, compute the value of the test statistic, and find the p-value to determine whether the null hypothesis should be rejected.
- With two independent simple random samples, we showed that the point estimator $\bar{x}_1 - \bar{x}_2$ has a standard error $\sigma_{\bar{x}_1 - \bar{x}_2}$ given by expression (10.2) and, when the sample sizes are large enough, the distribution of $\bar{x}_1 - \bar{x}_2$ can be described by a normal distribution.

6. Test Statistic for Hypothesis Tests About $\mu_1 - \mu_2$: σ_1 and σ_2 Known

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad (10.5)$$

- We demonstrated a two-tailed hypothesis test about the difference between two population means. Lower tail and upper tail tests can also be considered. These tests use the same test statistic as given in equation (10.5). The procedure for computing the p -value and the rejection rules for these one-tailed tests are the same as those for hypothesis tests involving a single population mean and single population proportion.

 Question (p486)

As part of a study to evaluate differences in education quality between two training centers, a standardized examination is given to individuals who are trained at the centers. The difference between the mean examination scores is used to assess quality differences between the centers. The population means for the two centers are as follows. μ_1 is the mean examination score for the population of individuals trained at center A , μ_2 is the mean examination score for the population of individuals trained at center B . We begin with the tentative assumption that no difference exists between the training quality provided at the two centers. The standardized examination given previously in a variety of settings always resulted in an examination score standard deviation near 10 points. Thus, we will use this information to assume that the population standard deviations are known with $\sigma_1 = 10$ and $\sigma_2 = 10$. An $\alpha = 0.05$ level of significance is specified for the study. Independent simple random samples of $n_1 = 30$ individuals from training center A and $n_2 = 40$ individuals from training center B are taken. The respective sample means are $\bar{x}_1 = 82$ and $\bar{x}_2 = 78$. Do these data suggest a significant difference between the population means at the two training centers? State the null and alternative hypotheses for this two-tailed test, compute the test statistic, and state the decision rules based on the p -value approach and the critical value approach and make the decision.

sol:

Practical Advice

1. In most applications of the interval estimation and hypothesis testing procedures presented in this section, random samples with $n_1 \geq 30$ and $n_2 \geq 30$ are adequate.
2. In cases where either or both sample sizes are less than 30, the distributions of the populations become important considerations.
3. In general, with smaller sample sizes, it is more important for the analyst to be satisfied that it is reasonable to assume that the distributions of the two populations are at least approximately normal.

10.2 Inferences About The Difference Between Two Population Means: σ_1 and σ_2 Unknown

1. Extend the discussion of inferences about the difference between two population means to the case when the two population standard deviations, σ_1 and σ_2 , are unknown.
2. In this case, we will use the sample standard deviations, s_1 and s_2 , to estimate the unknown population standard deviations.
3. When we use the sample standard deviations, the interval estimation and hypothesis testing procedures will be based on the t distribution rather than the standard normal distribution.

Interval Estimation of $\mu_1 - \mu_2$

1. Let us develop the margin of error and an interval estimate of the difference between these two population means. (Recall) The interval estimate for the case when the

population standard deviations, σ_1 and σ_2 , are known.

$$(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

2. With σ_1 and σ_2 unknown, we will use the sample standard deviations s_1 and s_2 to estimate σ_1 and σ_2 and replace $z_{\alpha/2}$ with $t_{\alpha/2}$.

3. Interval Estimate of the Difference Between Two Population Means: σ_1 and σ_2 Unknown

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \tag{10.6}$$


where $1 - \alpha$ is the confidence coefficient.

4. In this expression, the use of the t distribution is an approximation, but it provides excellent results and is relatively easy to use. The only difficulty that we encounter in using expression (10.6) is determining the appropriate degrees of freedom for $t_{\alpha/2}$.

5. Statistical software packages compute the appropriate degrees of freedom automatically. The formula used is as follows:

Degrees of Freedom: t Distribution With Two Independent Random Samples

$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{1}{n_1-1} \left(\frac{s_1^2}{n_1}\right)^2 + \frac{1}{n_2-1} \left(\frac{s_2^2}{n_2}\right)^2} \tag{10.7}$$

 **Question** (p490)

(Clearwater National Bank example) Clearwater National Bank is conducting a study designed to identify differences between checking account practices by customers at two of its branch banks. A simple random sample of 28 checking accounts is selected from the Cherry Grove Branch and an independent simple random sample of 22 checking accounts is selected from the Beechmont Branch. The current checking account balance is recorded for each of the checking accounts. A summary of the account balances follows:

	Cherry Grove	Beechmont
Sample Size	$n_1 = 28$	$n_2 = 22$
Sample Mean	$\bar{x}_1 = \$1025$	$\bar{x}_2 = \$910$
Sample Standard Deviation	$s_1 = \$150$	$s_2 = \$125$

Clearwater National Bank would like to estimate the difference between the mean checking account balance maintained by the population of Cherry Grove customers and the population of Beechmont customers. Compute a 95% confidence interval estimate of the difference between the population mean checking account balances at the two branch banks.

sol:

Hypothesis Tests About $\mu_1 - \mu_2$

1. (Recall) Letting D_0 denote the hypothesized difference between μ_1 and μ_2 , the test statistic used for the case where σ_1 and σ_2 are known is as follows.

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$


The test statistic, z , follows the standard normal distribution.

2. When σ_1 and σ_2 are unknown, we use s_1 as an estimator of σ_1 and s_2 as an estimator of σ_2 . Substituting these sample standard deviations for σ_1 and σ_2 provides the following test statistic when σ_1 and σ_2 are unknown.

3. Test Statistic for Hypothesis Tests About $\mu_1 - \mu_2$: σ_1 and σ_2 Unknown

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (10.8)$$

The degrees of freedom for t are given by equation (10.7).

 Question (p491)

Consider a new computer software package developed to help systems analysts reduce the time required to design, develop, and implement an information system. To evaluate the benefits of the new software package, a random sample of 24 systems analysts is selected. Each analyst is given specifications for a hypothetical information system. Then 12 of the analysts are instructed to produce the information system by using current technology. The other 12 analysts are trained in the use of the new software package and then instructed to use it to produce the information system. This study involves two populations: a population of systems analysts using the current technology and a population of systems analysts using the new software package. In terms of the time required to complete the information system design project, the population means are as follows. μ_1 is the mean project completion time for systems analysts using the current technology and μ_2 is the mean project completion time for systems analysts using the new software package. The researcher in charge of the new software evaluation project hopes to show that the new software package will provide a shorter mean project completion time. Thus, the researcher is looking for evidence to conclude that μ_2 is less than μ_1 ; in this case, the difference between the two population means, $\mu_1 - \mu_2$, will be greater than zero. Suppose that the 24 analysts complete the study with the results shown in Table 10.1.

TABLE 10.1 Completion Time Data and Summary Statistics for the Software Testing Study

	Current Technology	New Software
	300	274
	280	220
	344	308
	385	336
	372	198
	360	300
	288	315
	321	258
	376	318
	290	310
	301	332
	283	263
Summary Statistics		
Sample size	$n_1 = 12$	$n_2 = 12$
Sample mean	$\bar{x}_1 = 325$ hours	$\bar{x}_2 = 286$ hours
Sample standard deviation	$s_1 = 40$	$s_2 = 44$

Let the level of significance be $\alpha = 0.05$. State the null and the alternative hypothesis, the test statistic, p -value, the rejection rule, make a decision and conclusion.

sol:

(Software Output)

TABLE 10.2 Output for the Hypothesis Test on the Difference Between the Current and New Software Technology

	Current	New
Mean	325	286
Variance	1600	1936
Observations	12	12
<hr/>		
Hypothesized Mean Difference	0	
Degrees of Freedom	21	
Test Statistic	2.272	
One-Tail p-value	0.017	
One-Tail Critical Value	1.717	

Practical Advice

1. The interval estimation and hypothesis testing procedures presented in this section are robust and can be used with relatively small sample sizes.
2. In most applications, equal or nearly equal sample sizes such that the total sample size $n_1 + n_2 \geq 20$ can be expected to provide very good results even if the populations are not normal.
3. Larger sample sizes are recommended if the distributions of the populations are highly skewed or contain outliers.
4. Smaller sample sizes should only be used if the analyst is satisfied that the distributions of the populations are at least approximately normal.

Notes + Comments

1. How to make inferences about the difference between two population means when σ_1 and σ_2 are equal and unknown ($\sigma_1 = \sigma_2 = \sigma$)?
2. Based on above assumption, the two sample standard deviations are combined to provide the following pooled sample variance:

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

3. The t test statistic becomes

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

and has $n_1 + n_2 - 2$ degrees of freedom. At this point, the computation of the p -value and the interpretation of the sample results are identical to the procedures discussed earlier in this section.

4. A difficulty with this procedure is that the assumption that the two population standard deviations are equal is usually difficult to verify. Unequal population standard deviations are frequently encountered.
5. Using the pooled procedure may not provide satisfactory results, especially if the sample sizes n_1 and n_2 are quite different.
6. The t procedure that we presented in this section does not require the assumption of equal population standard deviations and can be applied whether the population standard deviations are equal or not. It is a more general procedure and is recommended for most applications.

10.3 Inferences About The Difference Between Two Population Means: Matched Samples

1. Example Matched.

- (a) Suppose employees at a manufacturing company can use two different methods to perform a production task. To maximize production output, the company wants to identify the method with the smaller population mean completion time.
- (b) Let μ_1 denote the population mean completion time for production method 1 and μ_2 denote the population mean completion time for production method 2.

- (c) With no preliminary indication of the preferred production method, we begin by tentatively assuming that the two production methods have the same population mean completion time. Thus, the null hypothesis is $H_0 : \mu_1 - \mu_2 = 0$.
- (d) If this hypothesis is rejected, we can conclude that the population mean completion times differ. In this case, the method providing the smaller mean completion time would be recommended.
- (e) The null and alternative hypotheses are written as follows.

$$H_0 : \mu_1 - \mu_2 = 0, \quad H_a : \mu_1 - \mu_2 \neq 0$$

2. In choosing the sampling procedure that will be used to collect production time data and test the hypotheses, we consider two alternative designs. One is based on independent samples and the other is based on matched samples.

- (a) Independent sample design: A simple random sample of workers is selected and each worker in the sample uses method 1. A second independent simple random sample of workers is selected and each worker in this sample uses method 2. The test of the difference between population means is based on the procedures in Section 10.2.
- (b) Matched sample design: One simple random sample of workers is selected. Each worker first uses one method and then uses the other method. The order of the two methods is assigned randomly to the workers, with some workers performing method 1 first and others performing method 2 first. Each worker provides a pair of data values, one value for method 1 and another value for method 2.

3. In the matched sample design the two production methods are tested under similar conditions (i.e., with the same workers); hence this design often leads to a smaller sampling error than the independent sample design. The primary reason is that in a matched sample design, variation between workers is eliminated because the same workers are used for both production methods.

4. Assuming the analysis of a matched sample design is the method used to test the difference between population means for the two production methods. The key

to the analysis of the matched sample design is to realize that we consider only the column of differences, d_i .

5. Therefore, we have six data values (0.6, -0.2, 0.5, 0.3, 0.0, and 0.6) that will be used to analyze the difference between population means of the two production methods.
6. Let μ_d is the mean of the difference in values for the population of workers. With this notation, the null and alternative hypotheses are rewritten as follows.

$$\underline{H_0 : \mu_d = 0, \quad H_a : \mu_d \neq 0}$$


7. Assume the population of differences has a normal distribution. This assumption is necessary so that we may use the t distribution for hypothesis testing and interval estimation procedures. Based on this assumption, the following test statistic has a t distribution with $n-1$ degrees of freedom.

8. **Test Statistic for Hypothesis Tests Involving Matched Samples**

$$t = \frac{\bar{d} - \mu_d}{s_d / \sqrt{n}}$$

where

$$\underline{\bar{d} = \frac{\sum d_i}{n}}, \quad \text{and} \quad \underline{s_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n - 1}}} \quad (10.9)$$

 **Question** (p498)

(Table 10.3) (Matched Example). A random sample of six workers is used. The data on completion times for the six workers are given in Table 10.3. Note that each worker provides a pair of data values, one for each production method. Also note that the last column contains the difference in completion times d_i for each worker in the sample. Assume that the population of differences has a normal distribution. Test the hypotheses $H_0 : \mu_d = 0$ and $H_a : \mu_d \neq 0$, using $\alpha = 0.05$. Compute the test statistic, the p -value and draw a conclusion. Compute the 95% confidence interval for the difference between the population means of the two production methods. If H_0 is rejected, we can conclude that the population mean completion times differ.

TABLE 10.3 Task Completion Times for a Matched Sample Design

Worker	Completion Time for Method 1 (minutes)	Completion Time for Method 2 (minutes)	Difference in Completion Times (d_i)
1	6.0	5.4	.6
2	5.0	5.2	-.2
3	7.0	6.5	.5
4	6.2	5.9	.3
5	6.0	6.0	.0
6	6.4	5.8	.6

sol:

Area in Upper Tail	0.20	0.10	0.05	0.025	0.01	0.005
t -Value (5 df)	0.920	1.476	2.015	2.571	3.365	4.032

10.4 Inferences About The Difference Between Two Population Proportions

1. Letting p_1 denote the proportion for population 1 and p_2 denote the proportion for population 2.
2. Consider inferences about the difference between the two population proportions: $p_1 - p_2$.
3. To make an inference about this difference, we will select two independent random samples consisting of n_1 units from population 1 and n_2 units from population 2.

Interval Estimation of $p_1 - p_2$

1. Example Tax Preparation Firm

A tax preparation firm is interested in comparing the quality of work at two of its regional offices. By randomly selecting samples of tax returns prepared at each office and verifying the sample returns' accuracy, the firm will be able to estimate the proportion of erroneous returns prepared at each office. Of particular interest is the difference between these proportions.

- (a) p_1 : proportion of erroneous returns for population 1 (office 1)
- (b) p_2 : proportion of erroneous returns for population 2 (office 2)
- (c) \bar{p}_1 : sample proportion for a simple random sample from population 1
- (d) \bar{p}_2 : sample proportion for a simple random sample from population 2

2. Point Estimator of the Difference Between Two Population Proportions

$$\hat{p}_1 - \hat{p}_2 = \bar{p}_1 - \bar{p}_2 \quad (10.10)$$

3. Thus, the point estimator of the difference between two population proportions is the difference between the sample proportions of two independent simple random samples.
4. As with other point estimators, the point estimator $\bar{p}_1 - \bar{p}_2$ has a sampling distribution that reflects the possible values of $\bar{p}_1 - \bar{p}_2$ if we repeatedly took two independent

random samples. The mean of this sampling distribution is $\underline{p_1 - p_2}$ and the standard error of $\bar{p}_1 - \bar{p}_2$ is:

Standard Error of $\bar{p}_1 - \bar{p}_2$

$$\sigma_{\bar{p}_1 - \bar{p}_2} = \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}} \quad (10.11)$$


- If the sample sizes are large enough that $\underline{n_1 p_1}$, $\underline{n_1(1-p_1)}$, $\underline{n_2 p_2}$, and $\underline{n_2(1-p_2)}$ are all greater than or equal to $\underline{5}$, the sampling distribution of $\bar{p}_1 - \bar{p}_2$ can be approximated by a normal distribution.
- With the sampling distribution of $\bar{p}_1 - \bar{p}_2$ approximated by a normal distribution, we would like to use $\underline{z_{\alpha/2} \sigma_{\bar{p}_1 - \bar{p}_2}}$ as the margin of error.
- However, $\sigma_{\bar{p}_1 - \bar{p}_2}$ given by equation (10.11) cannot be used directly because the two population proportions, p_1 and p_2 , are unknown. Using the sample proportion \bar{p}_1 to estimate p_1 and the sample proportion \bar{p}_2 to estimate p_2 , the margin of error is:

$$\text{Margin of error} = \underline{z_{\alpha/2} \sqrt{\frac{\bar{p}_1(1 - \bar{p}_1)}{n_1} + \frac{\bar{p}_2(1 - \bar{p}_2)}{n_2}}} \quad (10.12)$$

8. Interval Estimate of the Difference Between Two Population Proportions

$$\underline{(\bar{p}_1 - \bar{p}_2) \pm z_{\alpha/2} \sqrt{\frac{\bar{p}_1(1 - \bar{p}_1)}{n_1} + \frac{\bar{p}_2(1 - \bar{p}_2)}{n_2}}} \quad (10.13)$$

where $1-\alpha$ is the confidence coefficient.

 **Question** (p504)

(Tax Preparation Example) We find that independent simple random samples from the two offices provide the following information.

	Office 1	2
n_i	250	300
Number of returns with errors	35	27

Find a margin of error and interval estimate of the difference between the two population proportions. and 90% confidence interval.

sol:

Hypothesis Tests About $p_1 - p_2$

1. Let us now consider hypothesis tests about no difference between the proportions of two populations. In this case, the three forms for a hypothesis test are as follows:

$$\begin{array}{lll} H_0 : p_1 - p_2 \geq 0, & H_0 : p_1 - p_2 \leq 0, & H_0 : p_1 - p_2 = 0 \\ H_a : p_1 - p_2 < 0 & H_a : p_1 - p_2 > 0 & H_a : p_1 - p_2 \neq 0 \end{array}$$

2. When we assume H_0 is true as an equality, we have $p_1 - p_2 = 0$, which is the same as saying that the population proportions are equal, $p_1 = p_2$.
3. Under the assumption H_0 is true as an equality, the population proportions are equal and $p_1 = p_2 = p$. In this case, $\sigma_{\bar{p}_1 - \bar{p}_2}$ becomes **Standard Error of $\bar{p}_1 - \bar{p}_2$ when $p_1 = p_2 = p$**

$$\sigma_{\bar{p}_1 - \bar{p}_2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} = \sqrt{p(1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \quad (10.14)$$

4. With p unknown, we pool, or combine, the point estimators from the two samples (\bar{p}_1 and \bar{p}_2) to obtain a single point estimator of p as follows:

Pooled Estimator of p When $p_1 = p_2 = p$


$$\bar{p} = \frac{n_1 \bar{p}_1 + n_2 \bar{p}_2}{n_1 + n_2} \quad (10.15)$$

This pooled estimator of p is a weighted average of \bar{p}_1 and \bar{p}_2 .

5. Substituting \bar{p} for p in equation (10.14), we obtain an estimate of the standard error of $\bar{p}_1 - \bar{p}_2$. This estimate of the standard error is used in the test statistic.
6. The general form of the test statistic for hypothesis tests about the difference between two population proportions is the point estimator divided by the estimate of $\sigma_{\bar{p}_1 - \bar{p}_2}$.
7. **Test Statistic for Hypothesis Tests About $p_1 - p_2$**

$$z = \frac{\bar{p}_1 - \bar{p}_2}{\sqrt{\bar{p}(1 - \bar{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad (10.16)$$

This test statistic applies to large sample situations where n_1p_1 , $n_1(1-p_1)$, n_2p_2 , and $n_2(1-p_2)$ are all greater than or equal to 5.

 **Question** (p506)

(Tax Preparation Firm Example) Assume that the firm wants to use a hypothesis test to determine whether the error proportions differ between the two offices. A two-tailed test is required. Use $\alpha = 0.10$ as the level of significance. State the null and alternative hypotheses. Compute the test statistic, and the p -value for this two-tailed test. State the decision rule and draw a conclusion.

sol:

☺ **EXERCISES**

10.1 : 1, 2, 4, 6

10.2 : 9, 10, 13, 14, 15

10.3 : 19, 23, 24

10.4 : 28, 29, 31, 34

SUP : 38, 39, 44

“成功人士與其他人的區別，不是在於他們的力量或知識，而是在於他們有多堅持”

“The difference between a successful person and others is not a lack of strength, not a lack of knowledge, but rather a lack of will.”

— *Vince Lombardi (June 11, 1913 – September 3, 1970)*

國立政治大學 111 學年度第 1 學期 小考 (1) 考試命題紙

考試科目：統計學 (一)

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：10 月 18 日 (四) 14:10-15:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共3頁，印刷份數：90 份

計算機	課本	筆記	字典	手機平板筆電
O	×	×	×	×

備註：注意事項要看!! (範圍: §ch1~4)

注意事項: (1) 答案卷請寫上科目、系級、學號及姓名。(2) 請按題號順序書寫。(3) 每一題號需置於答案卷最左邊。(4) 中英文作答皆可。(5) 建議用深色原子筆。(6) 需要計算過程。(7) 總分共 120 分。(8) 請複寫下列宣誓詞至答案卷第一頁最上頭。

本人姓名 重視榮譽，以認真負責的態度，參與於本次 (線上遠距) 考試，恪遵各項考試規則，無任何不法或舞弊情事，如違誓言，願受校方最嚴厲之處罰，謹誓。

1. (15%) 名詞解釋 (不能只列出公式，需說明所使用符號的意思及公式代表的意義):

- (a) 統計學 (statistics)
- (b) 統計推論 (statistical inference)
- (c) 樣本空間 (sample space)
- (d) 聯合機率 (joint probability)
- (e) 獨立事件 (independent events)

2. (45%) 簡答題/問答題

- (a) (10%) 依課本所述，測量尺度 (scales of measurement) 有哪四種類型？每一種類型請各舉 2 個例子。
- (b) (5%) 收集樣本 (sample) 的目的為何？
- (c) (5%) 何謂「以相對次數法 (relative frequency method) 給予隨機實驗每一個結果一個機率值 (assigning probabilities)」？
- (d) (5%) 當我們在進行敘述統計，想了解一數值資料的分佈形狀 (distribution shape) 時，可採用 Chebyshev's Theorem 或 Empirical rule。請問何謂「Empirical rule」？適用於哪類型的資料？
- (e) (20%) 大學新生入學，校方紀錄了新生的資料，假設此資料中僅有性別 (男、女、未告知)、入學管道 (繁星推薦、個人申請、考試入學)、身高及體重等變數，若您想對此資料進行敘述統計 (descriptive statistics)，請依課本所述，應如何進行。(提示: 進行 VVV 處理，或計算 OOO(數值/統計量)，或繪製 XXX 圖表，以了解 YYY。)

考試日期：10 月 18 日 (四) 14:10-15:50

※准帶項目打「O」· 否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳· 試題一律採雙面印刷· 如有特殊印製需求· 請註記：

本試題共3頁· 印刷份數: 90 份

計算機	課本	筆記	字典	手機平板筆電
-----	----	----	----	--------

備註：注意事項要看!! (範圍: §ch1~4)

O	×	×	×	×
---	---	---	---	---

3. (20%) **Golf Course Complaints.** Recently, management at Oak Tree Golf Course received a few complaints about the condition of the greens. Several players complained that the greens are too fast. Rather than react to the comments of just a few, the Golf Association conducted a survey of 100 male and 100 female golfers. The survey results are summarized here.

Male Golfers			Female Golfers		
Greens Condition			Greens Condition		
Handicap	Too Fast	Fine	Handicap	Too Fast	Fine
Under 15	10	40	Under 15	1	9
15 or more	25	25	15 or more	39	51

- (a) (5%) Combine these two crosstabulations into one with Male and Female as the row labels and Too Fast and Fine as the column labels. Which group shows the highest percentage saying that the greens are too fast?
- (b) (5%) Refer to the initial crosstabulations. For those players with low handicaps (better players), which group (male or female) shows the higher percentage saying the greens are too fast?
- (c) (5%) Refer to the initial crosstabulations. For those players with higher handicaps, which group (male or female) shows the higher percentage saying the greens are too fast?
- (d) (5%) What conclusions can you draw about the preferences of men and women concerning the speed of the greens? Are the conclusions you draw from part (a) as compared with parts (b) and (c) consistent? Explain any apparent inconsistencies.

國立政治大學 111 學年度第 1 學期 小考 (1) 考試命題紙

考試科目：統計學 (一)

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：10 月 18 日 (四) 14:10-15:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共3頁，印刷份數：90 份

計算機	課本	筆記	字典	手機平板筆電
O	×	×	×	×

備註：注意事項要看!! (範圍：§ch1~4)

4. (20%) **Apple iPads in Schools.** The New York Times reported that Apple has unveiled a new iPad marketed specifically to school districts for use by students (The New York Times website). The 9.7-inch iPads will have faster processors and a cheaper price point in an effort to take market share away from Google Chromebooks in public school districts. Suppose that the following data represent the percentages of students currently using Apple iPads for a sample of 18 U.S. public school districts ($x_i, i = 1, \dots, 18$). ($\sum x_i = 444, \sum x_i^2 = 13346$.)

15 22 12 21 26 18 42 29 64 20 15 22 18 24 27 24 26 19

- (a) (5%) Compare the first and second quartiles for these data.
 - (b) (5%) Compute the variance and standard deviation for these data.
 - (c) (5%) Are there any outliers in these data?
 - (d) (5%) Based on your calculated values, what can we say about the percentage of students using iPads in public school districts?
5. (20%) **Treatment-Caused Injuries.** A study of 31,000 hospital admissions in New York State found that 4% of the admissions led to treatment-caused injuries. One-seventh of these treatment-caused injuries resulted in death, and one-fourth were caused by negligence. Malpractice claims were filed in one out of 7.5 cases involving negligence, and payments were made in one out of every two claims.
- (a) (5%) Represent the events indicated in this problem using letters or symbols.
 - (b) (5%) What is the probability a person admitted to the hospital will suffer a treatment-caused injury due to negligence?
 - (c) (5%) What is the probability a person admitted to the hospital will die from a treatment-caused injury?
 - (d) (5%) In the case of a negligent treatment-caused injury, what is the probability a malpractice claim will be paid?

注意：1、考試求公平及公正，請同學務必自律，維護學校與學生之榮譽。

2、考試時不得有交談、窺視、夾帶、抄襲、傳遞、代考或其它作弊等舞弊行為，考畢務必交卷，不得攜卷出場，違者依考場規則議處。

國立政治大學 111 學年度第 1 學期 小考 (2) 考試命題紙

考試科目：統計學 (一)

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：12 月 20 日 (二) 14:10-15:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共2頁，印刷份數：90 份

計算機	課本	筆記	字典	手機平板筆電
-----	----	----	----	--------

備註：注意事項要看!! (範圍: ch6~ch7)

O	×	×	×	×
---	---	---	---	---

注意事項: (1) 答案卷請寫上科目、系級、學號及姓名。(2) 請按題號順序書寫。(3) 每一題號需置於答案卷最左邊。(4) 中英文作答皆可。(5) 建議用深色原子筆。(6) 需要計算過程。(7) 總分共 100+20 分。(8) 請複寫下列宣誓詞至答案卷第一頁最上頭 (或空白之處)。

本人姓名 重視榮譽，以認真負責的態度，參與於本次 (線上遠距) 考試，恪遵各項考試規則，無任何不法或舞弊情事，如違誓言，願受校方最嚴厲之處罰，謹誓。

1. (15%) 名詞解釋 (說明中若有列出公式，需說明所使用符號的意思及公式代表的意義):

- (a) 樣本統計量 (Sample Statistic)
- (b) 抽樣分佈 (Sampling Distribution)
- (c) 中央極限定理 (Central Limit Theorem)

2. (45%) 簡答題/問答題

- (a) (10%) The percentage of values in some commonly used intervals: _____ (_____, _____) of the values of a normal random variable are within plus or minus one (two, three) standard deviation of its mean. (整數百分比，或小數點以下一位百分比，皆可)
- (b) (5%) 機率分佈的重要性為何? (What is the important role of the probability distributions?) (不可回答: 「可以算機率。」)
- (c) (10%) 令 X 是二項式隨機變數，代表 n 次伯努力 (Bernoulli) 試驗中，成功的次數。令 p 為一次試驗成功的機率。請問何謂「常態逼近二項式機率 (Normal Approximation of Binomial Probabilities)」? 以計算 $P(X \leq k)$ 來說明 (k 為小於等於 n 之整數)。(註: 需說明此逼近的先決條件 (例如: $np \geq \dots$)，並令 $\mu = np$ and \dots 及說明 Continuity correction)
- (d) (10%) 卜瓦松分佈與指數分佈的關係為何? (Relationship Between the Poisson and Exponential Distributions)
- (e) (10%) 依課本所述，評估一個點估計量的優劣，有哪些準則? 各是什麼意思?

國立政治大學 111 學年度第 1 學期 小考(2) 考試命題紙

考試科目：統計學(一)

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：12月20日(二) 14:10-15:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共2頁，印刷份數：90份

計算機	課本	筆記	字典	手機平板筆電
-----	----	----	----	--------

備註：注意事項要看!! (範圍: ch6~ch7)

O	×	×	×	×
---	---	---	---	---

3. (10%) **Americans Who Believe Global Warming Is Occurring.** According to a Yale program on climate change communication survey, 71% of Americans think global warming is happening (American Psychological Association website). For a sample of 150 Americans, what is the probability that at least 100 believe global warming is occurring? Use the normal approximation of the binomial distribution to answer this question. ($\text{pnorm}(1.26) = 0.8962$, $\text{pnorm}(0.22) = 0.5871$)
4. (10%) **Arrival of Vehicles at an Intersection.** The time between arrivals of vehicles at a particular intersection follows an exponential probability distribution with a mean of 12 seconds. What is the probability that the arrival time between vehicles is 12 seconds or less? (用 e 表示即可)
5. (10%) **Income Tax Return Preparation Fees.** The CPA Practice Advisor reports that the mean preparation fee for 2017 federal income tax returns was \$273. Use this price as the population mean and assume the population standard deviation of preparation fees is \$100. What is the probability that the mean price for a sample of 50 federal income tax returns is within \$16 of the population mean? ($\text{pnorm}(1.13) = 0.8708$, $\text{pnorm}(0.16) = 0.5636$)
6. (10%) **Product Labeling.** The Grocery Manufacturers of America reported that 76% of consumers read the ingredients listed on a product's label. Assume the population proportion is $p = 0.76$ and a sample of 400 consumers is selected from the population. What is the probability that the sample proportion will be within ± 0.03 of the population proportion? ($\text{pnorm}(1.40) = 0.9192$, $\text{pnorm}(0.79) = 0.7852$)
7. (加分: 20%) 證明樣本比例 \bar{p} (sample proportion) 的抽樣分佈是常態分佈，具有平均數 $\mu = p$ 及標準差 $\sigma_{\bar{p}} = \sqrt{\frac{p(1-p)}{n}}$ 。其中 p 為母體比例， n 為樣本數。(假設 sample size is large enough)

註 1: 若您所需的機率值，考卷中並無提供，請直接註明「未提供機率值，無法計算出答案」。

註 2: $\text{pnorm}(x) = P(Z \leq x)$, $Z \sim N(0, 1)$ 。

注意：1、考試求公平及公正，請同學務必自律，維護學校與學生之榮譽。

2、考試時不得有交談、窺視、夾帶、抄襲、傳遞、代考或其它作弊等舞弊行為，考畢務必交卷，不得攜卷出場，違者依考場規則議處。

國立政治大學 111 學年度第 1 學期 期中考 考試命題紙

考試科目：統計學（一）

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：11 月 15 日（二）13:10-14:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共2頁，印刷份數：85 份

計算機	課本	筆記	字典	手機平板筆電
-----	----	----	----	--------

備註：注意事項要看!! (範圍: §1~§5)

O	×	×	×	×
---	---	---	---	---

注意事項: (1) 答案卷請寫上科目、系級、學號及姓名。(2) 請按題號順序書寫。(3) 每一題號需置於答案卷最左邊。(4) 中英文作答皆可。(5) 可用鉛筆，建議用深色原子筆。(6) 需要計算過程。(7) 小數請計算至小數點以下 4 位。(8) 交回題目卷及答案卷。(9) 總分共 100 分。(8) 請複寫下列宣誓詞至答案卷第一頁最上頭或空白處 (沒有寫的扣 10 分)。

本人姓名 重視榮譽，以認真負責的態度，參與於本次 (線上遠距) 考試，恪遵各項考試規則，無任何不法或舞弊情事，如違誓言，願受校方最嚴厲之處罰，謹誓。

1. (25%; 5% each) 統計名詞解釋 (不能只列出公式，需說明所使用符號的意思及公式代表的意義):

- (a) Covariance (共變異數)
- (b) z -score (z -分數)
- (c) Discrete random variables (離散型隨機變數)(禁止寫: 隨機變數是離散型的)
- (d) Empirical discrete distribution (經驗離散型分佈)
- (e) Poisson probability distribution (卜瓦松機率分佈)

2. (25%; 5% each) 簡答題/問答題

- (a) 何謂貝氏定理 (Bayes' theorem)? (提示: A_1, A_2, \dots, A_n ，貝氏定理的條件是什麼? 貝氏定理的結果是什麼? 若有使用符號，需解釋符號的意思。)
- (b) 依教科書所述，二項式實驗需具有哪四項特性?
- (c) 要應用卜瓦松分佈需有什麼假設?
- (d) 令 x 是 n 個試驗 (trials) 中成功 (success) 的次數，請問超幾何分佈和二項式分佈有何相同及不相同的地方?
- (e) 依教科書所述，以統計學的觀點，要如何評估一個金融投資組合 (financial portfolio)?

注意：1、考試求公平及公正，請同學務必自律，維護學校與學生之榮譽。

2、考試時不得有交談、窺視、夾帶、抄襲、傳遞、代考或其它作弊等舞弊行為，考畢務必交卷，不得攜卷出場，違者依考場規則議處。

國立政治大學 111 學年度第 1 學期 期中考 考試命題紙

考試科目：統計學（一）

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：11 月 15 日（二）13:10-14:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共2頁，印刷份數：85 份

計算機	課本	筆記	字典	手機平板筆電
O	×	×	×	×

備註：注意事項要看!! (範圍：§1~§5)

3. (20%; 5% each) **Emails Received (modified)**. According to a 2017 survey conducted by the technology market research firm The Radicati Group, U.S. office workers receive an average of 120 emails per day (Entrepreneur magazine website). Assume the number of emails received per hour follows a Poisson distribution.

- (a) What is the probability of receiving no emails during an hour?
- (b) What is the probability of receiving at least three emails during an hour?
- (c) What is the expected number of emails received during 15 minutes?
- (d) What is the probability that no emails are received during 15 minutes?

4. (30%; 5% each) **Investment Portfolio of Index Fund and Core Bonds Fund**. J.P. Morgan Asset Management publishes information about financial investments. Between 2002 and 2011, the expected return for the S&P 500 was 5.04% with a standard deviation of 19.45% and the expected return over that same period for a core bonds fund was 5.78% with a standard deviation of 2.13% (J.P. Morgan Asset Management, Guide to the Markets). The publication also reported that the correlation between the S&P 500 and core bonds is -0.32 . You are considering portfolio investments that are composed of an S&P 500 index fund and a core bonds fund.

- (a) Using the information provided, determine the covariance between the S&P 500 and core bonds.
- (b) Construct a portfolio that is 50% invested in an S&P 500 index fund and 50% in a core bonds fund. In percentage terms, what are the expected return and standard deviation for such a portfolio?
- (c) Construct a portfolio that is 20% invested in an S&P 500 index fund and 80% invested in a core bonds fund. In percentage terms, what are the expected return and standard deviation for such a portfolio?
- (d) Construct a portfolio that is 80% invested in an S&P 500 index fund and 20% invested in a core bonds fund. In percentage terms, what are the expected return and standard deviation for such a portfolio?
- (e) Which of the portfolios in parts (b), (c), and (d) has the largest expected return? Which has the smallest standard deviation? Which of these portfolios is the best investment?
- (f) Discuss the advantages and disadvantages of investing in the three portfolios in parts (b), (c), and (d). Would you prefer investing all your money in the S&P 500 index, the core bonds fund, or one of the three portfolios? Why?

國立政治大學 111 學年度第 1 學期 期末考 考試命題紙

考試科目：統計學（一）

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：01 月 10 日（二）13:10-14:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共2頁，印刷份數：85 份

計算機	課本	筆記	字典	手機平板筆電
-----	----	----	----	--------

備註：注意事項要看!! (範圍: §6~§9)

O	×	×	×	×
---	---	---	---	---

注意事項: (1) 答案卷請寫上科目、系級、學號及姓名。(2) 請按題號順序書寫。(3) 每一題號需置於答案卷最左邊。(4) 中英文作答皆可。(5) 可用鉛筆，建議用深色原子筆。(6) 需要計算過程。(7) 小數請計算至小數點以下 4 位。(8) 交回題目卷及答案卷。(9) 總分共 100 分。(8) 請複寫下列宣誓詞至答案卷第一頁最上頭或空白處 (沒有寫的扣 10 分)。

本人姓名 重視榮譽，以認真負責的態度，參與於本次 (線上遠距) 考試，恪遵各項考試規則，無任何不法或舞弊情事，如違誓言，願受校方最嚴厲之處罰，謹誓。

1. (25%; 5% each) 統計名詞解釋 (不能只列出公式，需說明所使用符號的意思及公式代表的意義):

- (a) 誤差邊際 (margin of error)。
- (b) 95% 信賴區間 (95% confidence interval)。
- (c) 假設檢定 (hypothesis testing)。
- (d) 顯著水準 α 值為 0.05 (level of significance, $\alpha = 0.05$)。
- (e) p 值 (p -value)。

2. (25%; 5% each) 簡答題/問答題

- (a) 區間估計 (interval estimate) 的目的為何?
- (b) 當在進行母體平均 (Population Mean) 的區間估計時，母體標準差 (σ) 已知及未知的差異為何? (禁止直接講「使用公式不同」而不做任何說明。)
- (c) 教科書中有提到「The sampling distributions of \bar{x} and \bar{p} play key roles in computing these interval estimates」。為什麼樣本平均數 \bar{x} 和樣本比例 \bar{p} 的抽樣分佈對於計算母體參數 (μ, p) 的區間估計扮演著重要的角色?
- (d) 教科書中指出要計算母體比例的信賴區間，所需要的樣本數公式為： $n = \frac{(z_{\alpha/2})^2 p^*(1 - p^*)}{E^2}$ 。若分析者對資料不十分了解，且無法事先獲得 p^* 的估計，則教科書建議使用 $p^* = 0.5$ ，請問理由為何?
- (e) 假設母體標準差 (σ) 已知，要進行母體平均 (μ) 的區間估計和雙尾檢定。請以此例子，說明區間估計和假設檢定的關係。

3. (10%) 證明樣本比例 \bar{p} (sample proportion) 的抽樣分佈是常態分佈，具有平均數 μ 及標準差 $\sigma_{\bar{p}}$ 。(其中母體比例記為 p ，樣本數為 n 且足夠大。需寫出 μ 及 $\sigma_{\bar{p}}$)

國立政治大學 111 學年度第 1 學期 期末考 考試命題紙

考試科目：統計學（一）

開課班別：統計學整合開課

命題教授：吳漢銘

考試日期：01 月 10 日（二）13:10-14:50

※准帶項目打「O」，否則打「×」

1. 需加發計算紙或答案紙請在試題內封袋備註。
2. 為環保節能減碳，試題一律採雙面印刷，如有特殊印製需求，請註記：

本試題共2頁，印刷份數：85 份

計算機	課本	筆記	字典	手機平板筆電
-----	----	----	----	--------

備註：注意事項要看!! (範圍：§6~§9)

O	×	×	×	×
---	---	---	---	---

** 假設檢定若有要求「conclusion」時，不能只有「Reject 或 Accept H_0 」，要回到題目去做結論。

4. (20%; 5% each) **Scholarship Examination Scores.** At Western University the historical mean of scholarship examination scores for freshman applications is 900. A historical population standard deviation $s = 180$ is assumed known. Each year, the assistant dean uses a sample of applications to determine whether the mean examination score for the new freshman applications has changed.

- (a) State the hypotheses.
- (b) What is the 95% confidence interval estimate of the population mean examination score if a sample of 200 applications provided a sample mean $\bar{x} = 935$?
- (c) Use the confidence interval to conduct a hypothesis test. Using $\alpha = 0.05$, what is your conclusion?
- (d) What is the p -value?

pnorm(c(0.19, 2, 2.55, 2.65, 2.75, 2.85, 2.95)):
0.5753 0.9772 0.9946 0.9960 0.9970 0.9978 0.9984

5. (10%) **Starting Salaries for Business Graduates.** Michigan State University's Collegiate Employment Research Institute found that starting salary for recipients of bachelor's degrees in business was \$50,032 in 2017. The results for a sample of 100 business majors receiving a bachelor's degree in 2018 showed a mean starting salary of \$51,276 with a sample standard deviation of \$5200. Conduct a hypothesis test to determine whether the mean starting salary for business majors in 2018 is greater than the mean starting salary in 2017. Use $\alpha = 0.01$ as the level of significance.

Area in Lower Tail	0.20	0.10	0.05	0.025	0.01	0.005
t -Value (99 df)	-0.845	-1.290	-1.660	-1.984	-2.365	-2.626

6. (10%; 5% each) **Construction Worker Idle Time.** Shorney Construction Company bids on projects assuming that the mean idle time per worker is 72 or fewer minutes per day. A sample of 30 construction workers will be used to test this assumption. Assume that the population standard deviation is 20 minutes.

- (a) State the hypotheses to be tested.
- (b) What is the probability of making a Type II error when the population mean idle time is 80 minutes? (註：若 $Z \sim N(0,1)$ ，則 $P(Z \leq z_\beta) = \beta$ 。請以 $P(Z \leq z_\beta)$ 表示型二誤差的機率。)

注意：1、考試求公平及公正，請同學務必自律，維護學校與學生之榮譽。

2、考試時不得有交談、窺視、夾帶、抄襲、傳遞、代考或其它作弊等舞弊行為，考畢務必交卷，不得攜卷出場，違者依考場規則議處。

“人生路上遇到的挫折只是為了要塑造你、讓你達成目標”

“The struggles along the way are only meant to shape you for your purpose.”

— *Chadwick Boseman (November 29, 1976 – August 28, 2020)*