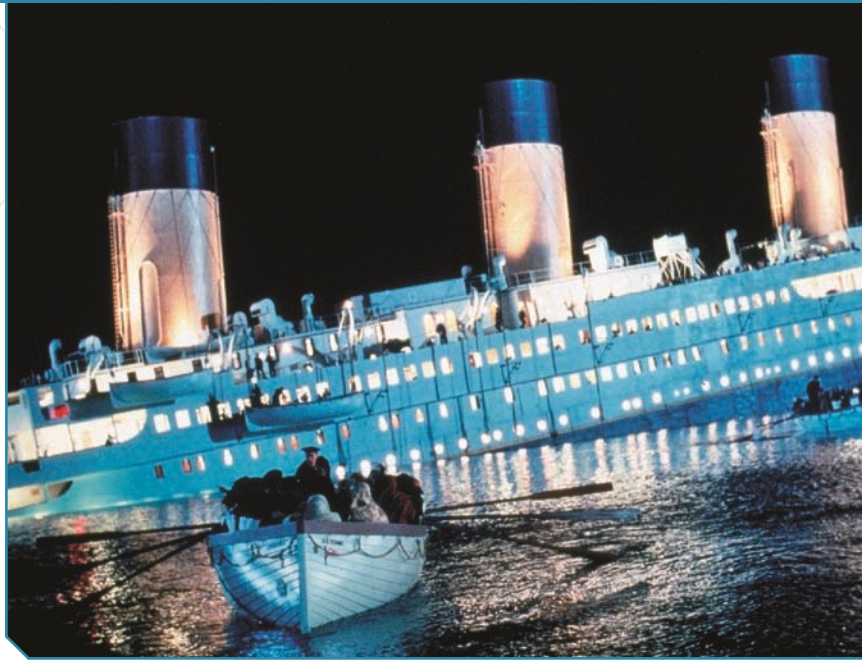


Displaying and Describing Categorical Data



WHO	People on the <i>Titanic</i>
WHAT	Survival status, age, sex, ticket class
WHEN	April 14, 1912
WHERE	North Atlantic
HOW	A variety of sources and Internet sites
WHY	Historical interest

What happened on the *Titanic* at 11:40 on the night of April 14, 1912, is well known. Frederick Fleet’s cry of “Iceberg, right ahead” and the three accompanying pulls of the crow’s nest bell signaled the beginning of a nightmare that has become legend. By 2:15 a.m., the *Titanic*, thought by many to be unsinkable, had sunk, leaving more than 1500 passengers and crew members on board to meet their icy fate.

Here are some data about the passengers and crew aboard the *Titanic*. Each case (row) of the data table represents a person on board the ship. The variables are the person’s *Survival* status (Dead or Alive), *Age* (Adult or Child), *Sex* (Male or Female), and ticket *Class* (First, Second, Third, or Crew).

The problem with a data table like this—and in fact with all data tables—is that you can’t *see* what’s going on. And seeing is just what we want to do. We need ways to show the data so that we can see patterns, relationships, trends, and exceptions.

A S **Video: The Incident** tells the story of the *Titanic*, and includes rare film footage.

Survival	Age	Sex	Class
Dead	Adult	Male	Third
Dead	Adult	Male	Crew
Dead	Adult	Male	Third
Dead	Adult	Male	Crew
Dead	Adult	Male	Crew
Dead	Adult	Male	Crew
Alive	Adult	Female	First
Dead	Adult	Male	Third
Dead	Adult	Male	Crew

Table 3.1

Part of a data table showing four variables for nine people aboard the *Titanic*.

The Three Rules of Data Analysis



FIGURE 3.1 A Picture to Tell a Story

Florence Nightingale (1820–1910), a founder of modern nursing, was also a pioneer in health management, statistics, and epidemiology. She was the first female member of the British Statistical Society and was granted honorary membership in the newly formed American Statistical Association.

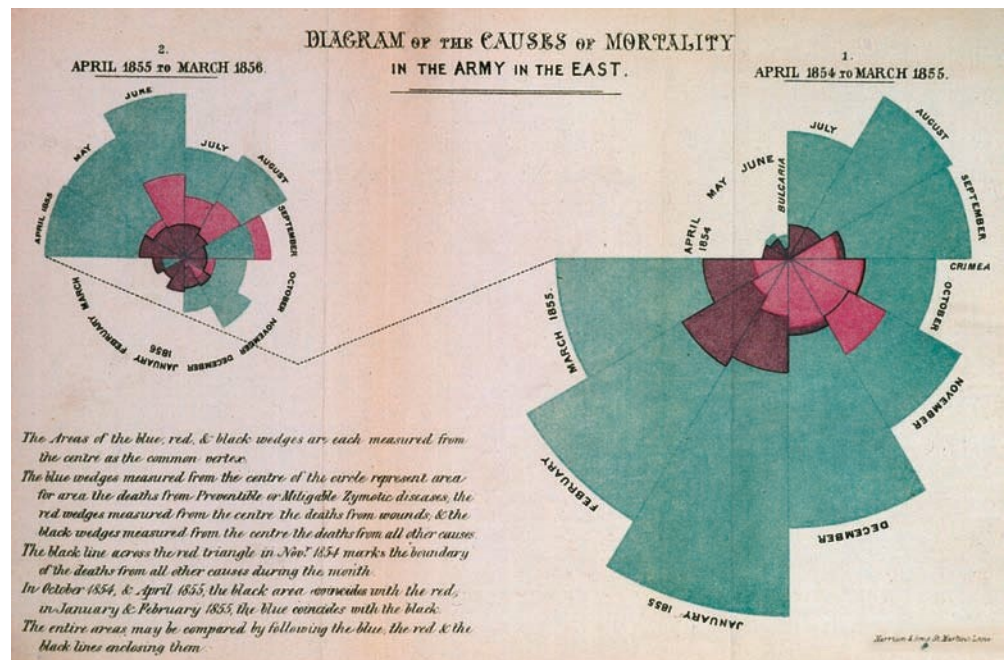
To argue forcefully for better hospital conditions for soldiers, she and her colleague, Dr. William Farr, invented this display, which showed that in the Crimean War, far more soldiers died of illness and infection than of battle wounds. Her campaign succeeded in improving hospital conditions and nursing for soldiers.

Florence Nightingale went on to apply statistical methods to a variety of important health issues and published more than 200 books, reports, and pamphlets during her long and illustrious career.

So, what should we do with data like these? There are three things you should always do first with data:

1. **Make a picture.** A display of your data will reveal things you are not likely to see in a table of numbers and will help you to *Think* clearly about the patterns and relationships that may be hiding in your data.
2. **Make a picture.** A well-designed display will *Show* the important features and patterns in your data. A picture will also show you the things you did not expect to see: the extraordinary (possibly wrong) data values or unexpected patterns.
3. **Make a picture.** The best way to *Tell* others about your data is with a well-chosen picture.

These are the three rules of data analysis. There are pictures of data throughout the book, and new kinds keep showing up. These days, technology makes drawing pictures of data easy, so there is no reason not to follow the three rules.



Frequency Tables: Making Piles

AS **Activity:** Make and examine a table of counts. Even data on something as simple as hair color can reveal surprises when you organize it in a data table.

Class	Count
First	325
Second	285
Third	706
Crew	885

Table 3.2

A frequency table of the *Titanic* passengers.

To make a picture of data, the first thing we have to do is to make piles. Making piles is the beginning of understanding about data. We pile together things that seem to go together, so we can see how the cases distribute across different categories. For categorical data, piling is easy. We just count the number of cases corresponding to each category and pile them up.

One way to put all 2201 people on the *Titanic* into piles is by ticket *Class*, counting up how many had each kind of ticket. We can organize these counts into a frequency table, which records the totals and the category names.

Even when we have thousands of cases, a variable like ticket *Class*, with only a few categories, has a frequency table that's easy to read. A frequency table with dozens or hundreds of categories would be much harder to read. We use the names of the categories to label each row in the frequency table. For ticket *Class*, these are "First," "Second," "Third," and "Crew."

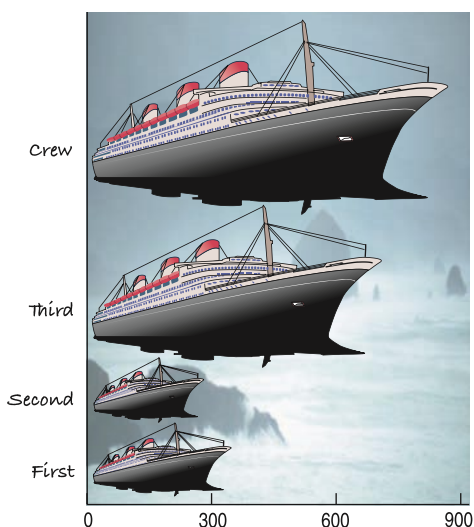
Class	%
First	14.77
Second	12.95
Third	32.08
Crew	40.21

Table 3.3

A relative frequency table for the same data.

Counts are useful, but sometimes we want to know the fraction or **proportion** of the data in each category, so we divide the counts by the total number of cases. Usually we multiply by 100 to express these proportions as **percentages**. A **relative frequency table** displays the *percentages*, rather than the counts, of the values in each category. Both types of tables show how the cases are distributed across the categories. In this way, they describe the **distribution** of a categorical variable because they name the possible categories and tell how frequently each occurs.

The Area Principle

**FIGURE 3.2**

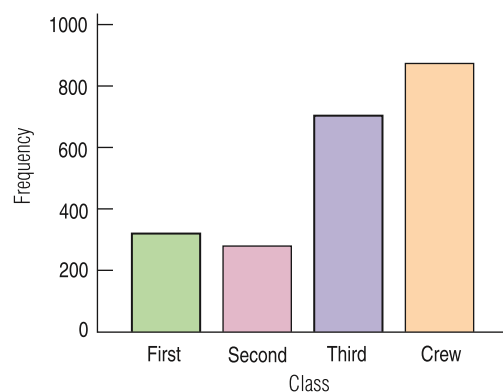
How many people were in each class on the *Titanic*? From this display, it looks as though the service must have been great, since most aboard were crew members. Although the length of each ship here corresponds to the correct number, the impression is all wrong. In fact, only about 40% were crew.

Now that we have the frequency table, we're ready to follow the three rules of data analysis and make a picture of the data. But a bad picture can distort our understanding rather than help it. Here's a graph of the *Titanic* data. What impression do you get about who was aboard the ship?

It sure looks like most of the people on the *Titanic* were crew members, with a few passengers along for the ride. That doesn't seem right. What's wrong? The lengths of the ships *do* match the totals in the table. (You can check the scale at the bottom.) However, experience and psychological tests show that our eyes tend to be more impressed by the *area* than by other aspects of each ship image. So, even though the *length* of each ship matches up with one of the totals, it's the associated *area* in the image that we notice. Since there were about 3 times as many crew as second-class passengers, the ship depicting the number of crew is about 3 times longer than the ship depicting second-class passengers, but it occupies about 9 times the area. As you can see from the frequency table (Table 3.2), that just isn't a correct impression.

The best data displays observe a fundamental principle of graphing data called the **area principle**. The area principle says that the area occupied by a part of the graph should correspond to the magnitude of the value it represents. Violations of the area principle are a common way to lie (or, since most mistakes are unintentional, we should say err) with Statistics.

Bar Charts

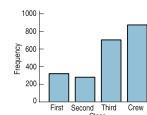
**FIGURE 3.3 People on the Titanic by Ticket Class**

With the area principle satisfied, we can see the true distribution more clearly.

Here's a chart that obeys the area principle. It's not as visually entertaining as the ships, but it does give an *accurate* visual impression of the distribution. The height of each bar shows the count for its category. The bars are the same width, so their heights determine their areas, and the areas are proportional to the counts in each class. Now it's easy to see that the majority of people on board were *not* crew, as the ships picture led us to believe. We can also see that there were about 3 times as many crew as second-class passengers. And there were more than twice as many third-class passengers as either first- or second-class passengers, something you may have missed in the frequency table. Bar charts make these kinds of comparisons easy and natural.

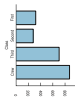
A **bar chart** displays the distribution of a categorical variable, showing the counts for each category next to each other for easy comparison. Bar charts should have small spaces between the bars to indicate that these are freestanding bars that could be rearranged into any order. The bars are lined up along a common base.

Usually they stick up like this



but sometimes they run

sideways like this



If we really want to draw attention to the relative *proportion* of passengers falling into each of these classes, we could replace the counts with percentages and use a **relative frequency bar chart**.

AS

Activity: Bar Charts.

Watch bar charts grow from data; then use your statistics package to create some bar charts for yourself.

For some reason, some computer programs give the name “bar chart” to any graph that uses bars. And others use different names according to whether the bars are horizontal or vertical. Don’t be misled. “Bar chart” is the term for a *display of counts of a categorical variable with bars*.

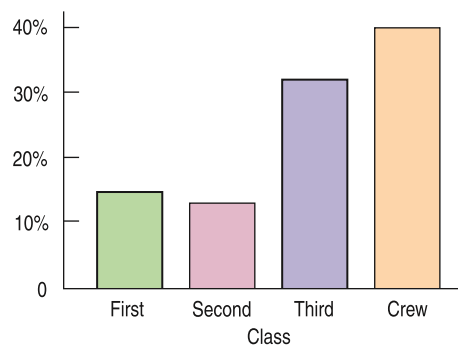


FIGURE 3.4

The relative frequency bar chart looks the same as the bar chart (Figure 3.3) but shows the *proportion* of people in each category rather than the counts.

Pie Charts

Another common display that shows how a whole group breaks into several categories is a pie chart. **Pie charts** show the whole group of cases as a circle. They slice the circle into pieces whose sizes are proportional to the fraction of the whole in each category.

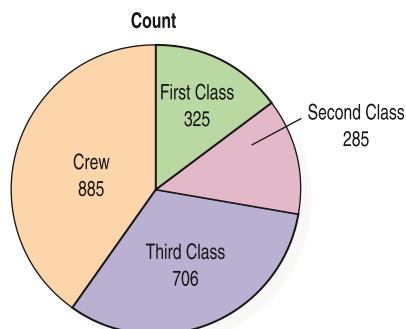


FIGURE 3.5 Number of Titanic passengers in each class

Pie charts give a quick impression of how a whole group is partitioned into smaller groups. Because we’re used to cutting up pies into 2, 4, or 8 pieces, pie charts are good for seeing relative frequencies near $1/2$, $1/4$, or $1/8$. For example, you may be able to tell that the pink slice, representing the second-class passengers, is very close to $1/8$ of the total. It’s harder to see that there were about twice as many third-class as first-class passengers. Which category had the most passengers? Were there more crew or more third-class passengers? Comparisons such as these are easier in a bar chart.

Think before you draw. Our first rule of data analysis is *Make a picture*. But what kind of picture? We don’t have a lot of options—yet. There’s more to Statistics than pie charts and bar charts, and knowing when to use each type of graph is a critical first step in data analysis. That decision depends in part on what type of data we have.

It’s important to check that the data are appropriate for whatever method of analysis you choose. Before you make a bar chart or a pie chart, always check the

Categorical Data Condition: The data are counts or percentages of individuals in categories.

If you want to make a relative frequency bar chart or a pie chart, you'll need to also make sure that the categories don't overlap so that no individual is counted twice. If the categories do overlap, you can still make a bar chart, but the percentages won't add up to 100%. For the *Titanic* data, either kind of display is appropriate because the categories don't overlap.

Throughout this course, you'll see that doing Statistics right means selecting the proper methods. That means you have to *Think* about the situation at hand. An important first step, then, is to check that the type of analysis you plan is appropriate. The Categorical Data Condition is just the first of many such checks.

Contingency Tables: Children and First-Class Ticket Holders First?

AS **Activity: Children at Risk.**
This activity looks at the fates of children aboard the *Titanic*; the subsequent activity shows how to make such tables on a computer.

We know how many tickets of each class were sold on the *Titanic*, and we know that only about 32% of all those aboard the *Titanic* survived. After looking at the distribution of each variable by itself, it's natural and more interesting to ask how they relate. Was there a relationship between the kind of ticket a passenger held and the passenger's chances of making it into the lifeboat? To answer this question, we need to look at the two categorical variables *Class* and *Survival* together.

To look at two categorical variables together, we often arrange the counts in a two-way table. Here is a two-way table of those aboard the *Titanic*, classified according to the class of ticket and whether the ticket holder survived or didn't. Because the table shows how the individuals are distributed along each variable, contingent on the value of the other variable, such a table is called a **contingency table**.

Contingency table of ticket *Class* and *Survival*. The bottom line of "Totals" is the same as the previous frequency table.
Table 3.4

		Class				Total
		First	Second	Third	Crew	
Survival	Alive	203	118	178	212	711
	Dead	122	167	528	673	1490
	Total	325	285	706	885	2201

The margins of the table, both on the right and at the bottom, give totals. The bottom line of the table is just the frequency distribution of ticket *Class*. The right column of the table is the frequency distribution of the variable *Survival*. When presented like this, in the margins of a contingency table, the frequency distribution of one of the variables is called its **marginal distribution**.

Each **cell** of the table gives the count for a combination of values of the two variables. If you look down the column for second-class passengers to the first cell, you can see that 118 second-class passengers survived. Looking at the third-class passengers, you can see that more third-class passengers (178) survived. Were second-class passengers more likely to survive? Questions like this are easier to address by using percentages. The 118 survivors in second class were 41.4% of the total 285 second-class passengers, while the 178 surviving third-class passengers were only 25.2% of that class's total.

We know that 118 second-class passengers survived. We could display this number as a percentage—but as a percentage of what? The total number of passengers? (118 is 5.4% of the total: 2201.) The number of second-class passengers?



A bell-shaped artifact from the *Titanic*.

(118 is 41.4% of the 285 second-class passengers.) The number of survivors? (118 is 16.6% of the 711 survivors.) All of these are possibilities, and all are potentially useful or interesting. You'll probably wind up calculating (or letting your technology calculate) lots of percentages. Most statistics programs offer a choice of total percent, row percent, or column percent for contingency tables. Unfortunately, they often put them all together with several numbers in each cell of the table. The resulting table holds lots of information, but it can be hard to understand:

Another contingency table of ticket Class. This time we see not only the counts for each combination of *Class* and *Survival* (in bold) but the percentages these counts represent. For each count, there are three choices for the percentage: by row, by column, and by table total. There's probably too much information here for this table to be useful.

Table 3.5

		Class					
		First	Second	Third	Crew	Total	
Survival	Alive	Count	203	118	178	212	711
		% of Row	28.6%	16.6%	25.0%	29.8%	100%
		% of Column	62.5%	41.4%	25.2%	24.0%	32.3%
		% of Table	9.2%	5.4%	8.1%	9.6%	32.3%
	Dead	Count	122	167	528	673	1490
		% of Row	8.2%	11.2%	35.4%	45.2%	100%
		% of Column	37.5%	58.6%	74.8%	76.0%	67.7%
		% of Table	5.6%	7.6%	24.0%	30.6%	67.7%
	Total	Count	325	285	706	885	2201
		%of Row	14.8%	12.9%	32.1%	40.2%	100%
		% of Column	100%	100%	100%	100%	100%
		% of Table	14.8%	12.9%	32.1%	40.2%	100%

To simplify the table, let's first pull out the percent of table values:

A contingency table of Class by Survival with only the table percentages

Table 3.6

		Class				
		First	Second	Third	Crew	Total
Survival	Alive	9.2%	5.4%	8.1%	9.6%	32.3%
	Dead	5.6%	7.6%	24.0%	30.6%	67.7%
	Total	14.8%	12.9%	32.1%	40.2%	100%

These percentages tell us what percent of *all* passengers belong to each combination of column and row category. For example, we see that although 8.1% of the people aboard the *Titanic* were surviving third-class ticket holders, only 5.4% were surviving second-class ticket holders. Is this fact useful? Comparing these percentages, you might think that the chances of surviving were better in third class than in second. But be careful. There were many more third-class than second-class passengers on the *Titanic*, so there were more third-class survivors. That group is a larger percentage of the passengers, but is that really what we want to know?

Percent of what? The English language can be tricky when we talk about percentages. If you're asked "What percent of *the survivors* were in second class?" it's pretty clear that we're interested only in survivors. It's as if we're restricting the *Who* in the question to the survivors, so we should look at the number of second-class passengers among all the survivors—in other words, the row percent.

But if you're asked "What percent were second-class passengers who survived?" you have a different question. Be careful; here, the *Who* is everyone on board, so 2201 should be the denominator, and the answer is the table percent.

And if you're asked "What percent of the second-class passengers survived?" you have a third question. Now the *Who* is the second-class passengers, so the denominator is the 285 second-class passengers, and the answer is the column percent. Always be sure to ask "percent of what?" That will help you to know the *Who* and whether we want *row*, *column*, or *table* percentages.

FOR EXAMPLE

Finding marginal distributions

In January 2007, a Gallup poll asked 1008 Americans age 18 and over whether they planned to watch the upcoming Super Bowl. The pollster also asked those who planned to watch whether they were looking forward more to seeing the football game or the commercials. The results are summarized in the table:

Question: What's the marginal distribution of the responses?

To determine the percentages for the three responses, divide the count for each response by the total number of people polled:

$$\frac{479}{1008} = 47.5\% \quad \frac{237}{1008} = 23.5\% \quad \frac{292}{1008} = 29.0\%$$

According to the poll, 47.5% of American adults were looking forward to watching the Super Bowl game, 23.5% were looking forward to watching the commercials, and 29% didn't plan to watch at all.

Response	Sex			
	Male	Female	Total	
	Game	279	200	479
	Commercials	81	156	237
	Won't watch	132	160	292
	Total	492	516	1008

Conditional Distributions

The more interesting questions are *contingent*. We'd like to know, for example, what percentage of *second-class passengers* survived and how that compares with the survival rate for third-class passengers.

It's more interesting to ask whether the chance of surviving the *Titanic* sinking *depended* on ticket class. We can look at this question in two ways. First, we could ask how the distribution of ticket *Class* changes between survivors and non-survivors. To do that, we look at the *row percentages*:

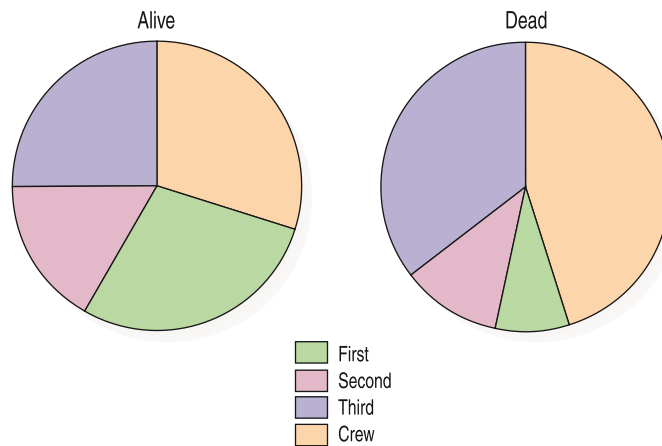
The conditional distribution of ticket *Class* conditioned on each value of *Survival*: *Alive* and *Dead*.
Table 3.7

		Class				
		First	Second	Third	Crew	Total
Survival	Alive	203	118	178	212	711
		28.6%	16.6%	25.0%	29.8%	100%
	Dead	122	167	528	673	1490
		8.2%	11.2%	35.4%	45.2%	100%

By focusing on each row separately, we see the distribution of class under the *condition* of surviving or not. The sum of the percentages in each row is 100%, and we divide that up by ticket class. In effect, we temporarily restrict the *Who* first to survivors and make a pie chart for them. Then we refocus the *Who* on the nonsurvivors and make their pie chart. These pie charts show the distribution of ticket classes *for each row* of the table: survivors and nonsurvivors. The distributions we create this way are called **conditional distributions**, because they show the distribution of one variable for just those cases that satisfy a condition on another variable.

FIGURE 3.6

Pie charts of the conditional distributions of ticket Class for the survivors and nonsurvivors, separately. Do the distributions appear to be the same? We're primarily concerned with percentages here, so pie charts are a reasonable choice.



FOR EXAMPLE

Finding conditional distributions

Recap: The table shows results of a poll asking adults whether they were looking forward to the Super Bowl game, looking forward to the commercials, or didn't plan to watch.

Question: How do the conditional distributions of interest in the commercials differ for men and women?

	Sex		
	Male	Female	Total
Game	279	200	479
Commercials	81	156	237
Won't watch	132	160	292
Total	492	516	1008

Look at the group of people who responded "Commercials" and determine what percent of them were male and female:

$$\frac{81}{237} = 34.2\% \quad \frac{156}{237} = 65.8\%$$

Women make up a sizable majority of the adult Americans who look forward to seeing Super Bowl commercials more than the game itself. Nearly 66% of people who voiced a preference for the commercials were women, and only 34% were men.

But we can also turn the question around. We can look at the distribution of *Survival* for each category of ticket *Class*. To do this, we look at the *column percentages*. Those show us whether the chance of surviving was roughly the same for each of the four classes. Now the percentages in each column add to 100%, because we've restricted the *Who*, in turn, to each of the four ticket classes:

A contingency table of *Class* by *Survival* with only counts and column percentages. Each column represents the conditional distribution of *Survival* for a given category of ticket *Class*.

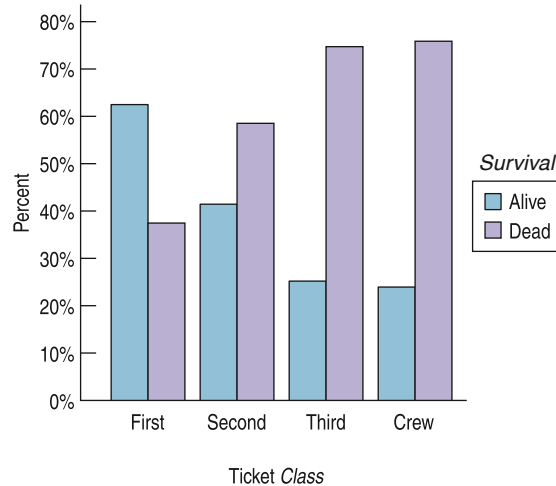
Table 3.8

		Class				
		First	Second	Third	Crew	Total
Survival	Alive	Count 203 62.5%	Count 118 41.4%	Count 178 25.2%	Count 212 24.0%	Count 711 32.3%
	Dead	Count 122 37.5%	Count 167 58.6%	Count 528 74.8%	Count 673 76.0%	Count 1490 67.7%
	Total	Count 325 100%	Count 285 100%	Count 706 100%	Count 885 100%	Count 2201 100%

Looking at how the percentages change across each row, it sure looks like ticket class mattered in whether a passenger survived. To make it more vivid, we could show the distribution of *Survival* for each ticket class in a display. Here's a side-by-side bar chart showing percentages of surviving and not for each category:

FIGURE 3.7

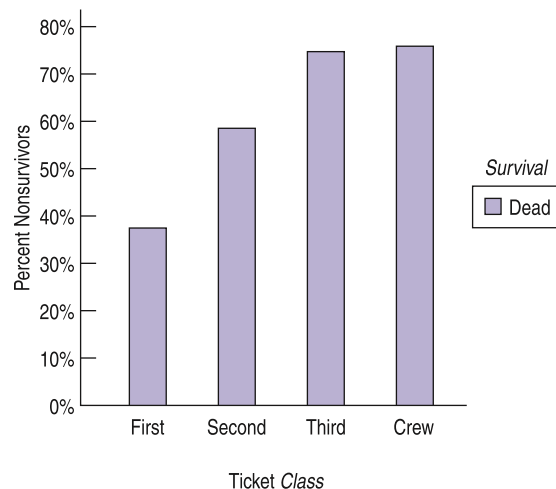
Side-by-side bar chart showing the conditional distribution of *Survival* for each category of *Ticket Class*. The corresponding pie charts would have only two categories in each of four pies, so bar charts seem the better alternative.



These bar charts are simple because, for the variable *Survival*, we have only two alternatives: Alive and Dead. When we have only two categories, we really need to know only the percentage of one of them. Knowing the percentage that survived tells us the percentage that died. We can use this fact to simplify the display even more by dropping one category. Here are the percentages of dying across the classes displayed in one chart:

FIGURE 3.8

Bar chart showing just nonsurvivor percentages for each value of *Ticket Class*. Because we have only two values, the second bar doesn't add any information. Compare this chart to the side-by-side bar chart shown earlier.



TI-*inspire*

Conditional distributions and association. Explore the *Titanic* data to see which passengers were most likely to survive.

Now it's easy to compare the risks. Among first-class passengers, 37.5% perished, compared to 58.6% for second-class ticket holders, 74.8% for those in third class, and 76.0% for crew members.

If the risk had been about the same across the ticket classes, we would have said that survival was *independent* of class. But it's not. The differences we see among these conditional distributions suggest that survival may have depended on ticket class. You may find it useful to consider conditioning on each variable in a contingency table in order to explore the dependence between them.

It is interesting to know that *Class* and *Survival* are associated. That's an important part of the *Titanic* story. And we know how important this is because the margins show us the actual numbers of people involved.

Variables can be associated in many ways and to different degrees. The best way to tell whether two variables are associated is to ask whether they are *not*.¹ In a contingency table, when the distribution of *one* variable is the same for all categories of another, we say that the variables are **independent**. That tells us there's no association between these variables. We'll see a way to check for independence formally later in the book. For now, we'll just compare the distributions.

FOR EXAMPLE

Looking for associations between variables

Recap: The table shows results of a poll asking adults whether they were looking forward to the Super Bowl game, looking forward to the commercials, or didn't plan to watch.

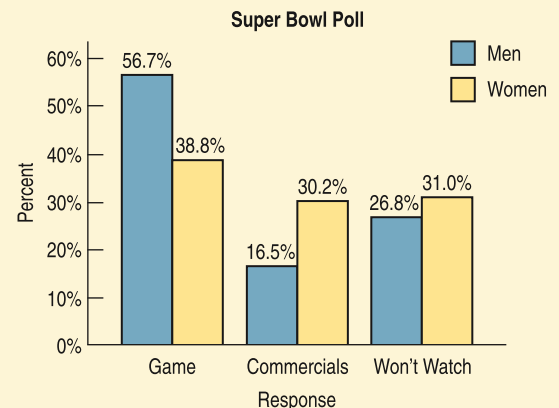
Question: Does it seem that there's an association between interest in Super Bowl TV coverage and a person's sex?

	Sex		
	Male	Female	Total
Game	279	200	479
Commercials	81	156	237
Won't watch	132	160	292
Total	492	516	1008

First find the distribution of the three responses for the men (the column percentages):

$$\frac{279}{492} = 56.7\% \quad \frac{81}{492} = 16.5\% \quad \frac{132}{492} = 26.8\%$$

Then do the same for the women who were polled, and display the two distributions with a side-by-side bar chart:



Based on this poll it appears that women were only slightly less interested than men in watching the Super Bowl telecast: 31% of the women said they didn't plan to watch, compared to just under 27% of men. Among those who planned to watch, however, there appears to be an association between the viewer's sex and what the viewer is most looking forward to. While more women are interested in the game (39%) than the commercials (30%), the margin among men is much wider: 57% of men said they were looking forward to seeing the game, compared to only 16.5% who cited the commercials.

¹This kind of "backwards" reasoning shows up surprisingly often in science—and in Statistics. We'll see it again.



JUST CHECKING

A Statistics class reports the following data on Sex and Eye Color for students in the class:

	Eye Color			
	Blue	Brown	Green/Hazel/Other	Total
Sex				
Males	6	20	6	32
Females	4	16	12	32
Total	10	36	18	64

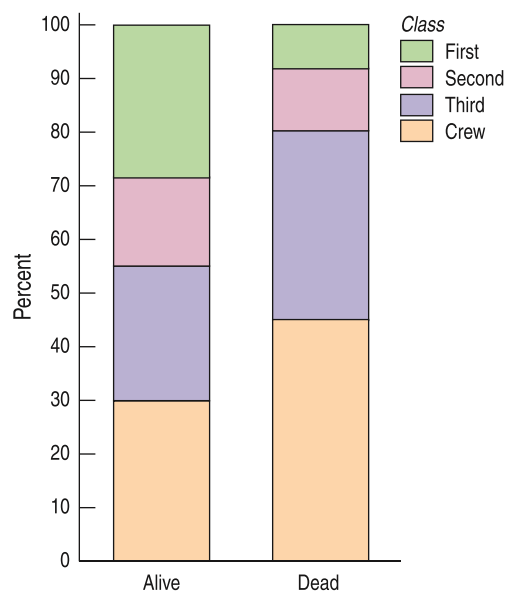
1. What percent of females are brown-eyed?
2. What percent of brown-eyed students are female?
3. What percent of students are brown-eyed females?
4. What's the distribution of Eye Color?
5. What's the conditional distribution of Eye Color for the males?
6. Compare the percent who are female among the blue-eyed students to the percent of all students who are female.
7. Does it seem that Eye Color and Sex are independent? Explain.

Segmented Bar Charts

We could display the *Titanic* information by dividing up bars rather than circles. The resulting **segmented bar chart** treats each bar as the “whole” and divides it proportionally into segments corresponding to the percentage in each group. We can clearly see that the distributions of ticket *Class* are different, indicating again that survival was not independent of ticket *Class*.

FIGURE 3.9 A segmented bar chart for Class by Survival

Notice that although the totals for survivors and nonsurvivors are quite different, the bars are the same height because we have converted the numbers to percentages. Compare this display with the side-by-side pie charts of the same data in Figure 3.6.



STEP-BY-STEP EXAMPLE

Examining Contingency Tables

Medical researchers followed 6272 Swedish men for 30 years to see if there was any association between the amount of fish in their diet and prostate cancer (“Fatty Fish Consumption and Risk of Prostate Cancer,” *Lancet*, June 2001). Their results are summarized in this table:



We asked for a picture of a man eating fish. This is what we got.

		Prostate Cancer	
		No	Yes
Fish Consumption	Never/seldom	110	14
	Small part of diet	2420	201
	Moderate part	2769	209
	Large part	507	42

Table 3.9

Question: Is there an association between fish consumption and prostate cancer?



Plan Be sure to state what the problem is about.

Variables Identify the variables and report the W's.

Be sure to check the appropriate condition.

I want to know if there is an association between fish consumption and prostate cancer.

The individuals are 6272 Swedish men followed by medical researchers for 30 years. The variables record their fish consumption and whether or not they were diagnosed with prostate cancer.

✓ **Categorical Data Condition:** I have counts for both fish consumption and cancer diagnosis. The categories of diet do not overlap, and the diagnoses do not overlap. It's okay to draw pie charts or bar charts.



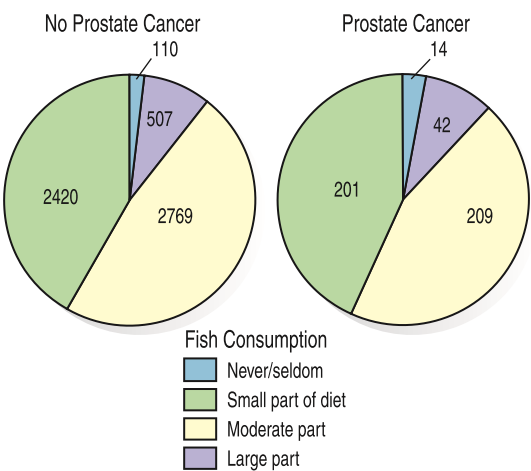
Mechanics It's a good idea to check the marginal distributions first before looking at the two variables together.

		Prostate Cancer		
		No	Yes	Total
Fish Consumption	Never/seldom	110	14	124 (2.0%)
	Small part of diet	2420	201	2621 (41.8%)
	Moderate part	2769	209	2978 (47.5%)
	Large part	507	42	549 (8.8%)
	Total	5806 (92.6%)	466 (7.4%)	6272 (100%)

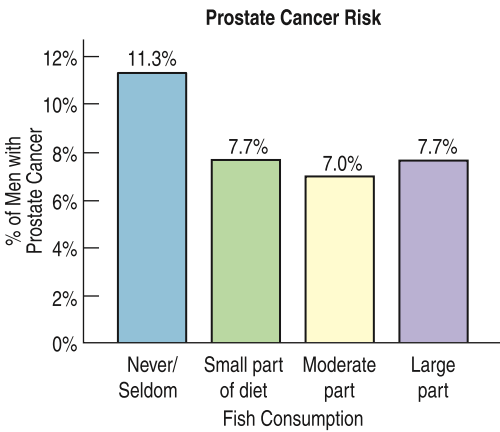
Two categories of the diet are quite small, with only 2.0% Never/Seldom eating fish and 8.8% in the “Large part” category. Overall, 7.4% of the men in this study had prostate cancer.

Then, make appropriate displays to see whether there is a difference in the relative proportions. These pie charts compare fish consumption for men who have prostate cancer to fish consumption for men who don't.

Both pie charts and bar charts can be used to compare conditional distributions. Here we compare prostate cancer rates based on differences in fish consumption.



It's hard to see much difference in the pie charts. So, I made a display of the row percentages. Because there are only two alternatives, I chose to display the risk of prostate cancer for each group:



Conclusion Interpret the patterns in the table and displays in context. If you can, discuss possible real-world consequences. Be careful not to overstate what you see. The results may not generalize to other situations.

Overall, there is a 7.4% rate of prostate cancer among men in this study. Most of the men (89.3%) ate fish either as a moderate or small part of their diet. From the pie charts, it's hard to see a difference in cancer rates among the groups. But in the bar chart, it looks like the cancer rate for those who never/seldom ate fish may be somewhat higher.

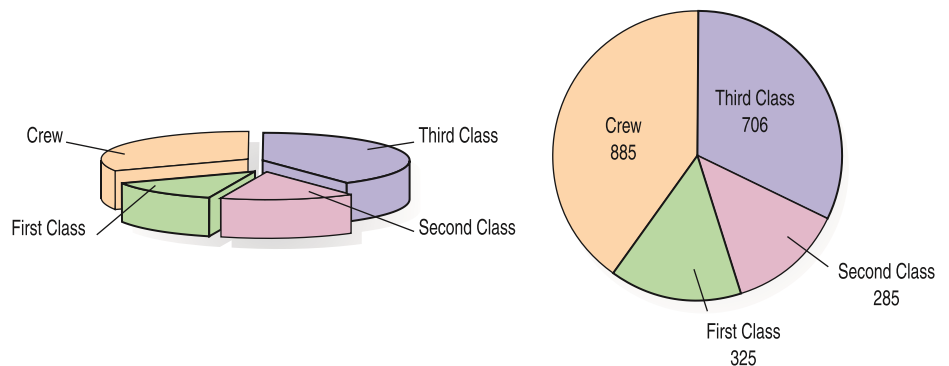
However, only 124 of the 6272 men in the study fell into this category, and only 14 of them developed prostate cancer. More study would probably be needed before we would recommend that men change their diets.²

² The original study actually used pairs of twins, which enabled the researchers to discern that the risk of cancer for those who never ate fish actually *was* substantially greater. Using pairs is a special way of gathering data. We'll discuss such study design issues and how to analyze the data in the later chapters.

This study is an example of looking at a sample of data to learn something about a larger population. We care about more than these particular 6272 Swedish men. We hope that learning about their experiences will tell us something about the value of eating fish in general. That raises the interesting question of what population we think this sample might represent. Do we hope to learn about all Swedish men? About all men? About the value of eating fish for all adult humans? ³ Often, it can be hard to decide just which population our findings may tell us about, but that also is how researchers decide what to look into in future studies.

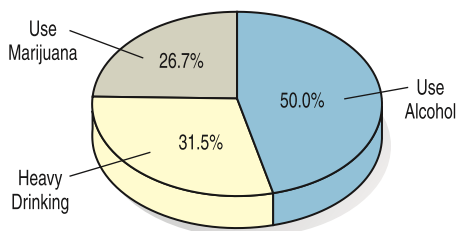
WHAT CAN GO WRONG?

- **Don't violate the area principle.** This is probably the most common mistake in a graphical display. It is often made in the cause of artistic presentation. Here, for example, are two displays of the pie chart of the *Titanic* passengers by class:



The one on the left looks pretty, doesn't it? But showing the pie on a slant violates the area principle and makes it much more difficult to compare fractions of the whole made up of each class—the principal feature that a pie chart ought to show.

- **Keep it honest.** Here's a pie chart that displays data on the percentage of high school students who engage in specified dangerous behaviors as reported by the Centers for Disease Control. What's wrong with this plot?

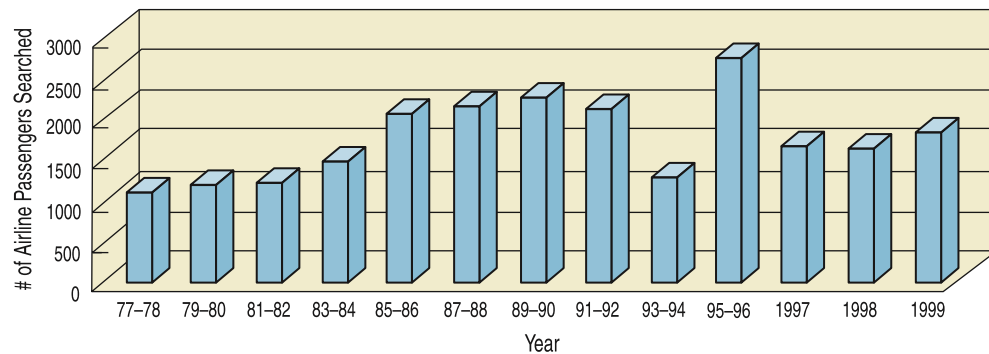


Try adding up the percentages. Or look at the 50% slice. Does it look right? Then think: What are these percentages of? Is there a "whole" that has been sliced up? In a pie chart, the proportions shown by each slice of the pie must add up to 100% and each individual must fall into only one category. Of course, showing the pie on a slant makes it even harder to detect the error.

(continued)

³ Probably not, since we're looking only at prostate cancer risk.

Here's another. This bar chart shows the number of airline passengers searched in security screening, by year:



Looks like things didn't change much in the final years of the 20th century—until you read the bar labels and see that the last three bars represent single years while all the others are for *pairs* of years. Of course, the false depth makes it harder to see the problem.

► **Don't confuse similar-sounding percentages.** These percentages sound similar but are different:

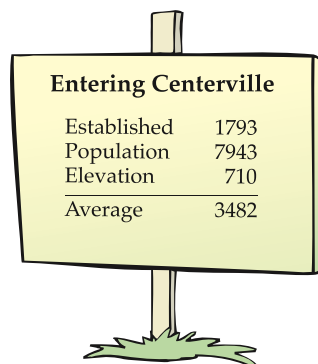
- The percentage of the passengers who were both in first class and survived: This would be $203/2201$, or 9.4%.
- The percentage of the first-class passengers who survived: This is $203/325$, or 62.5%.
- The percentage of the survivors who were in first class: This is $203/711$, or 28.6%.

In each instance, pay attention to the *Who* implicitly defined by the phrase. Often there is a restriction to a smaller group (all aboard the *Titanic*, those in first class, and those who survived, respectively) before a percentage is found. Your discussion of results must make these differences clear.

- **Don't forget to look at the variables separately, too.** When you make a contingency table or display a conditional distribution, be sure you also examine the marginal distributions. It's important to know how many cases are in each category.
- **Be sure to use enough individuals.** When you consider percentages, take care that they are based on a large enough number of individuals. Take care not to make a report such as this one:

We found that 66.67% of the rats improved their performance with training. The other rat died.

- **Don't overstate your case.** Independence is an important concept, but it is rare for two variables to be *entirely* independent. We can't conclude that one variable has no effect whatsoever on another. Usually, all we know is that little effect was observed in our study. Other studies of other groups under other circumstances could find different results.



SIMPSON'S PARADOX

- **Don't use unfair or silly averages.** Sometimes averages can be misleading. Sometimes they just don't make sense at all. Be careful when averaging different variables that the quantities you're averaging are comparable. The Centerville sign says it all.

When using averages of proportions across several different groups, it's important to make sure that the groups really are comparable.

It's easy to make up an example showing that averaging across very different values or groups can give absurd results. Here's how that might work: Suppose there are two pilots, Moe and Jill. Moe argues that he's the better pilot of the two, since he managed to land 83% of his last 120 flights on time compared with Jill's 78%. But let's look at the data a little more closely. Here are the results for each of their last 120 flights, broken down by the time of day they flew:

Table 3.10

On-time flights by *Time of Day* and *Pilot*. Look at the percentages within each *Time of Day* category. Who has a better on-time record during the day? At night? Who is better overall?

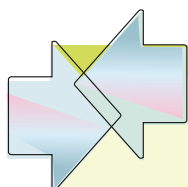
		Time of Day		
Pilot		Day	Night	Overall
	Moe	90 out of 100 90%	10 out of 20 50%	100 out of 120 83%
	Jill	19 out of 20 95%	75 out of 100 75%	94 out of 120 78%

One famous example of Simpson's paradox arose during an investigation of admission rates for men and women at the University of California at Berkeley's graduate schools. As reported in an article in *Science*, about 45% of male applicants were admitted, but only about 30% of female applicants got in. It looked like a clear case of discrimination. However, when the data were broken down by school (Engineering, Law, Medicine, etc.), it turned out that, within each school, the women were admitted at nearly the same or, in some cases, much *higher* rates than the men. How could this be? Women applied in large numbers to schools with very low admission rates (Law and Medicine, for example, admitted fewer than 10%). Men tended to apply to Engineering and Science. Those schools have admission rates above 50%. When the *average* was taken, the women had a much lower *overall* rate, but the average didn't really make sense.

Look at the daytime and nighttime flights separately. For day flights, Jill had a 95% on-time rate and Moe only a 90% rate. At night, Jill was on time 75% of the time and Moe only 50%. So Moe is better "overall," but Jill is better both during the day and at night. How can this be?

What's going on here is a problem known as **Simpson's paradox**, named for the statistician who discovered it in the 1960s. It comes up rarely in real life, but there have been several well-publicized cases. As we can see from the pilot example, the problem is *unfair averaging* over different groups. Jill has mostly night flights, which are more difficult, so her *overall average* is heavily influenced by her nighttime average. Moe, on the other hand, benefits from flying mostly during the day, with its higher on-time percentage. With their very different patterns of flying conditions, taking an overall average is misleading. It's not a fair comparison.

The moral of Simpson's paradox is to be careful when you average across different levels of a second variable. It's always better to compare percentages or other averages *within* each level of the other variable. The overall average may be misleading.



CONNECTIONS

All of the methods of this chapter work with *categorical variables*. You must know the *Who* of the data to know who is counted in each category and the *What* of the variable to know where the categories come from.



WHAT HAVE WE LEARNED?

We've learned that we can summarize categorical data by counting the number of cases in each category, sometimes expressing the resulting distribution as percents. We can display the distribution in a bar chart or a pie chart. When we want to see how two categorical variables are related, we put the counts (and/or percentages) in a two-way table called a contingency table.

- ▶ We look at the marginal distribution of each variable (found in the margins of the table).
- ▶ We also look at the conditional distribution of a variable within each category of the other variable.
- ▶ We can display these conditional and marginal distributions by using bar charts or pie charts.
- ▶ If the conditional distributions of one variable are (roughly) the same for every category of the other, the variables are independent.

Terms

Frequency table
(Relative frequency table)

Distribution

Area principle

Bar chart
(Relative frequency bar chart)

Pie chart

Categorical data condition

Contingency table

Marginal distribution

Conditional distribution

Independence

Segmented bar chart

Simpson's paradox

21. A frequency table lists the categories in a categorical variable and gives the count (or percentage) of observations for each category.

22. The distribution of a variable gives

- ▶ the possible values of the variable and
- ▶ the relative frequency of each value.

22. In a statistical display, each data value should be represented by the same amount of area.

22. Bar charts show a bar whose area represents the count (or percentage) of observations for each category of a categorical variable.

23. Pie charts show how a "whole" divides into categories by showing a wedge of a circle whose area corresponds to the proportion in each category.

24. The methods in this chapter are appropriate for displaying and describing categorical data. Be careful not to use them with quantitative data.

24. A contingency table displays counts and, sometimes, percentages of individuals falling into named categories on two or more variables. The table categorizes the individuals on all variables at once to reveal possible patterns in one variable that may be contingent on the category of the other.

24. In a contingency table, the distribution of either variable alone is called the marginal distribution. The counts or percentages are the totals found in the margins (last row or column) of the table.

26. The distribution of a variable restricting the *Who* to consider only a smaller group of individuals is called a conditional distribution.

29. Variables are said to be independent if the conditional distribution of one variable is the same for each category of the other. We'll show how to check for independence in a later chapter.

30. A segmented bar chart displays the conditional distribution of a categorical variable within each category of another variable.

34. When averages are taken across different groups, they can appear to contradict the overall averages. This is known as "Simpson's paradox."

Skills

THINK

- ▶ Be able to recognize when a variable is categorical and choose an appropriate display for it.
- ▶ Understand how to examine the association between categorical variables by comparing conditional and marginal percentages.

SHOW

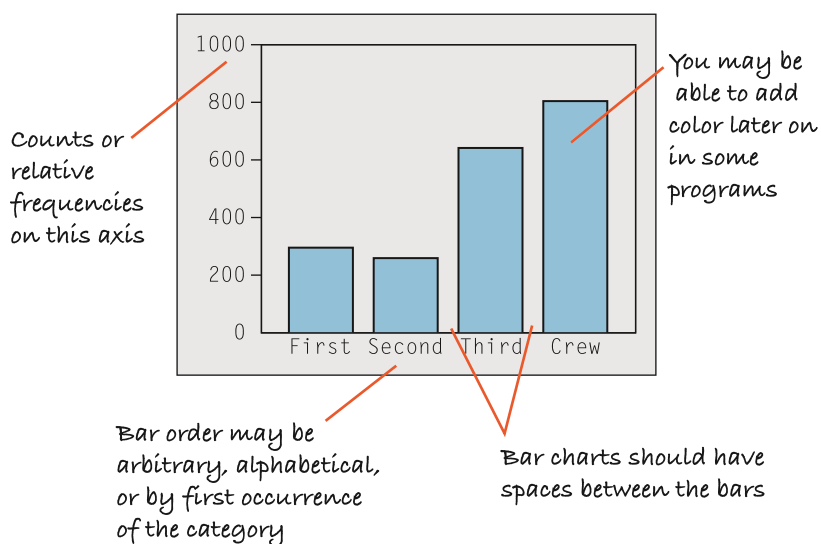
- ▶ Be able to summarize the distribution of a categorical variable with a frequency table.
- ▶ Be able to display the distribution of a categorical variable with a bar chart or pie chart.
- ▶ Know how to make and examine a contingency table.



- ▶ Know how to make and examine displays of the conditional distributions of one variable for two or more groups.
- ▶ Be able to describe the distribution of a categorical variable in terms of its possible values and relative frequencies.
- ▶ Know how to describe any anomalies or extraordinary features revealed by the display of a variable.
- ▶ Be able to describe and discuss patterns found in a contingency table and associated displays of conditional distributions.

DISPLAYING CATEGORICAL DATA ON THE COMPUTER

Although every package makes a slightly different bar chart, they all have similar features:



Sometimes the count or a percentage is printed above or on top of each bar to give some additional information. You may find that your statistics package sorts category names in annoying orders by default. For example, many packages sort categories alphabetically or by the order the categories are seen in the data set. Often, neither of these is the best choice.

EXERCISES

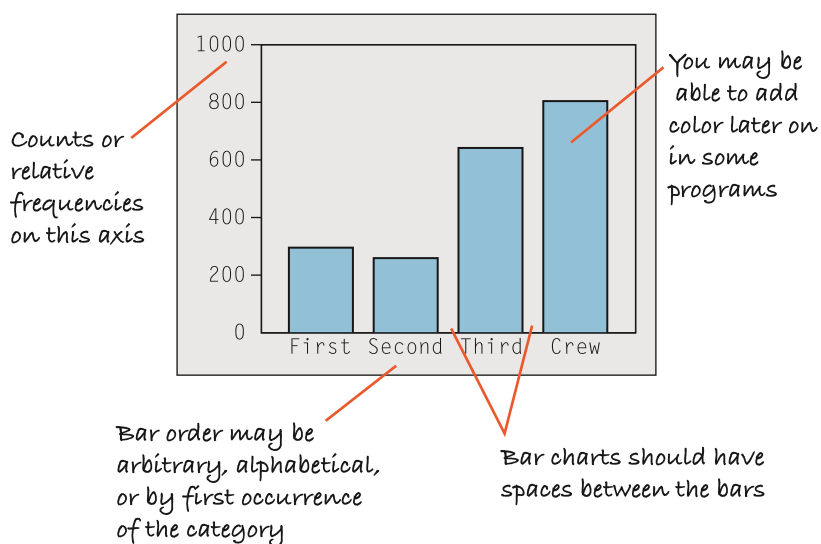
1. **Graphs in the news.** Find a bar graph of categorical data from a newspaper, a magazine, or the Internet.
 - a) Is the graph clearly labeled?
 - b) Does it violate the area principle?
 - c) Does the accompanying article tell the W's of the variable?
 - d) Do you think the article correctly interprets the data? Explain.
2. **Graphs in the news II.** Find a pie chart of categorical data from a newspaper, a magazine, or the Internet.
 - a) Is the graph clearly labeled?
 - b) Does it violate the area principle?
 - c) Does the accompanying article tell the W's of the variable?
 - d) Do you think the article correctly interprets the data? Explain.



- ▶ Know how to make and examine displays of the conditional distributions of one variable for two or more groups.
- ▶ Be able to describe the distribution of a categorical variable in terms of its possible values and relative frequencies.
- ▶ Know how to describe any anomalies or extraordinary features revealed by the display of a variable.
- ▶ Be able to describe and discuss patterns found in a contingency table and associated displays of conditional distributions.

DISPLAYING CATEGORICAL DATA ON THE COMPUTER

Although every package makes a slightly different bar chart, they all have similar features:



Sometimes the count or a percentage is printed above or on top of each bar to give some additional information. You may find that your statistics package sorts category names in annoying orders by default. For example, many packages sort categories alphabetically or by the order the categories are seen in the data set. Often, neither of these is the best choice.

EXERCISES

1. **Graphs in the news.** Find a bar graph of categorical data from a newspaper, a magazine, or the Internet.
 - a) Is the graph clearly labeled?
 - b) Does it violate the area principle?
 - c) Does the accompanying article tell the W's of the variable?
 - d) Do you think the article correctly interprets the data? Explain.
2. **Graphs in the news II.** Find a pie chart of categorical data from a newspaper, a magazine, or the Internet.
 - a) Is the graph clearly labeled?
 - b) Does it violate the area principle?
 - c) Does the accompanying article tell the W's of the variable?
 - d) Do you think the article correctly interprets the data? Explain.

3. **Tables in the news.** Find a frequency table of categorical data from a newspaper, a magazine, or the Internet.

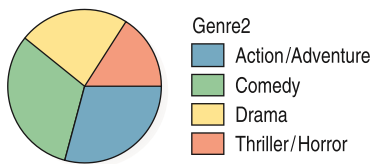
- Is it clearly labeled?
- Does it display percentages or counts?
- Does the accompanying article tell the W's of the variable?
- Do you think the article correctly interprets the data? Explain.

4. **Tables in the news II.** Find a contingency table of categorical data from a newspaper, a magazine, or the Internet.

- Is it clearly labeled?
- Does it display percentages or counts?
- Does the accompanying article tell the W's of the variables?
- Do you think the article correctly interprets the data? Explain.

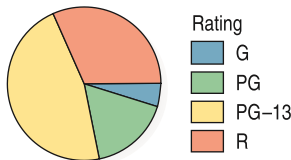
- T** 5. **Movie genres.** The pie chart summarizes the genres of 120 first-run movies released in 2005.

- Is this an appropriate display for the genres? Why/why not?
- Which genre was least common?



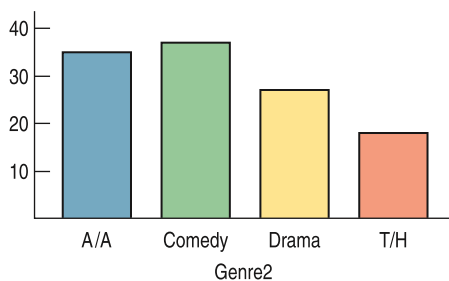
- T** 6. **Movie ratings.** The pie chart shows the ratings assigned to 120 first-run movies released in 2005.

- Is this an appropriate display for these data? Explain.
- Which was the most common rating?



- T** 7. **Genres again.** Here is a bar chart summarizing the 2005 movie genres, as seen in the pie chart in Exercise 5.

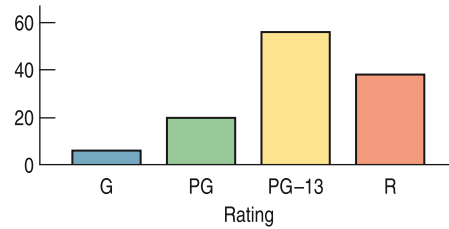
- Which genre was most common?
- Is it easier to see that in the pie chart or the bar chart? Explain.



- T** 8. **Ratings again.** Here is a bar chart summarizing the 2005 movie ratings, as seen in the pie chart in Exercise 6.

- Which was the least common rating?
- An editorial claimed that there's been a growth in PG-13 rated films that, according to the writer, "have too much sex and violence," at the expense of G-rated

films that offer "good, clean fun." The writer offered the bar chart below as evidence to support his claim. Does the bar chart support his claim? Explain.



9. **Magnet schools.** An article in the Winter 2003 issue of *Chance* magazine reported on the Houston Independent School District's magnet schools programs. Of the 1755 qualified applicants, 931 were accepted, 298 were wait-listed, and 526 were turned away for lack of space. Find the relative frequency distribution of the decisions made, and write a sentence describing it.

10. **Magnet schools again.** The *Chance* article about the Houston magnet schools program described in Exercise 9 also indicated that 517 applicants were black or Hispanic, 292 Asian, and 946 white. Summarize the relative frequency distribution of ethnicity with a sentence or two (in the proper context, of course).

11. **Causes of death 2004.** The Centers for Disease Control and Prevention (www.cdc.gov) lists causes of death in the United States during 2004:

Cause of Death	Percent
Heart disease	27.2
Cancer	23.1
Circulatory diseases and stroke	6.3
Respiratory diseases	5.1
Accidents	4.7

- Is it reasonable to conclude that heart or respiratory diseases were the cause of approximately 33% of U.S. deaths in 2003?
- What percent of deaths were from causes not listed here?
- Create an appropriate display for these data.

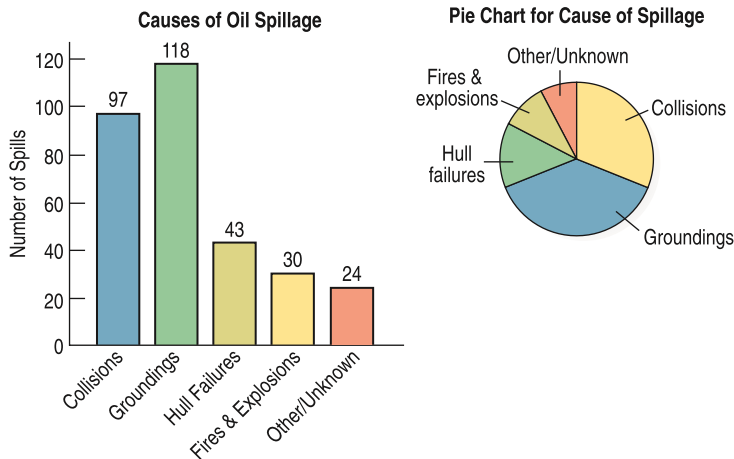
12. **Plane crashes.** An investigation compiled information about recent nonmilitary plane crashes (www.planecrashinfo.com). The causes, to the extent that they could be determined, are summarized in the table.

Cause	Percent
Pilot error	40
Other human error	5
Weather	6
Mechanical failure	14
Sabotage	6

- Is it reasonable to conclude that the weather or mechanical failures caused only about 20% of recent plane crashes?
- In what percent of crashes were the causes not determined?
- Create an appropriate display for these data.

13. **Oil spills 2006.** Data from the International Tanker Owners Pollution Federation Limited (www.itopf.com) give the cause of spillage for 312 large oil tanker accidents from 1974–2006. Here are displays.

- Write a brief report interpreting what the displays show.
- Is a pie chart an appropriate display for these data? Why or why not?

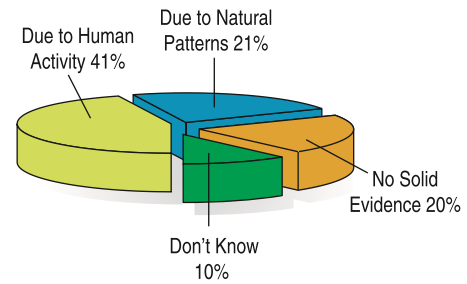


14. **Winter Olympics 2006.** Twenty-six countries won medals in the 2006 Winter Olympics. The table lists them, along with the total number of medals each won:

Country	Medals	Country	Medals
Germany	29	Finland	9
United States	25	Czech Republic	4
Canada	24	Estonia	3
Austria	23	Croatia	3
Russia	22	Australia	2
Norway	19	Poland	2
Sweden	14	Ukraine	2
Switzerland	14	Japan	1
South Korea	11	Belarus	1
Italy	11	Bulgaria	1
China	11	Great Britain	1
France	9	Slovakia	1
Netherlands	9	Latvia	1

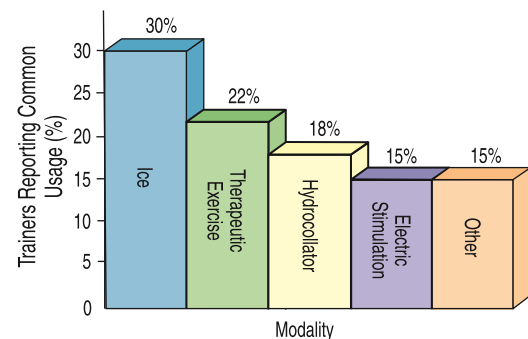
- Try to make a display of these data. What problems do you encounter?
- Can you find a way to organize the data so that the graph is more successful?

15. **Global Warming.** The Pew Research Center for the People and the Press (<http://people-press.org>) has asked a representative sample of U.S. adults about global warming, repeating the question over time. In January 2007, the responses reflected an increased belief that global warming is real and due to human activity. Here's a display of the percentages of respondents choosing each of the major alternatives offered:



List the errors in this display.

16. **Modalities.** A survey of athletic trainers (Scott F. Nadler, Michael Prybicien, Gerard A. Malanga, and Dan Sicher. "Complications from Therapeutic Modalities: Results of a National Survey of Athletic Trainers." *Archives of Physical Medical Rehabilitation* 84 [June 2003]) asked what modalities (treatment methods such as ice, whirlpool, ultrasound, or exercise) they commonly use to treat injuries. Respondents were each asked to list three modalities. The article included the following figure reporting the modalities used:



- What problems do you see with the graph?
- Consider the percentages for the named modalities. Do you see anything odd about them?

17. **Teen smokers.** The organization Monitoring the Future (www.monitoringthefuture.org) asked 2048 eighth graders who said they smoked cigarettes what brands they preferred. The table below shows brand preferences for two regions of the country. Write a few sentences describing the similarities and differences in brand preferences among eighth graders in the two regions listed.

Brand preference	South	West
Marlboro	58.4%	58.0%
Newport	22.5%	10.1%
Camel	3.3%	9.5%
Other (over 20 brands)	9.1%	9.5%
No usual brand	6.7%	12.9%

18. **Handguns.** In an effort to reduce the number of gun-related homicides, some cities have run buyback programs in which the police offer cash (often \$50) to anyone who turns in an operating handgun. *Chance* magazine looked at results from a four-year period in Milwaukee. The table on the next page shows what types of guns were turned in and what types were used in homicides during a four-year period. Write a few sentences comparing the two distributions.

Caliber of gun	Buyback	Homicide
Small (.22, .25, .32)	76.4%	20.3%
Medium (.357, .38, 9 mm)	19.3%	54.7%
Large (.40, .44, .45)	2.1%	10.8%
Other	2.2%	14.2%

- T 19. Movies by Genre and Rating.** Here's a table that classifies movies released in 2005 by genre and MPAA rating:

	G	PG	PG-13	R	Total
Action/Adventure	66.7	25	30.4	23.7	29.2
Comedy	33.3	60.0	35.7	10.5	31.7
Drama	0	15.0	14.3	44.7	23.3
Thriller/Horror	0	0	19.6	21.1	15.8
Total	100%	100%	100%	100%	100%

- The table gives column percents. How could you tell that from the table itself?
- What percentage of these movies were comedies?
- What percentage of the PG-rated movies were comedies?
- Which of the following can you learn from this table? Give the answer if you can find it from the table.
 - The percentage of PG-13 movies that were comedies
 - The percentage of dramas that were R-rated
 - The percentage of dramas that were G-rated
 - The percentage of 2005 movies that were PG-rated comedies

- T 20. The Last Picture Show.** Here's another table showing information about 120 movies released in 2005. This table gives percentages of the table total:

	G	PG	PG-13	R	Total
Action/Adventure	3.33%	4.17	14.2	7.50	29.2
Comedy	1.67	10	16.7	3.33	31.7
Drama	0	2.50	6.67	14.2	23.3
Thriller/Horror	0	0	9.17	6.67	15.8
Total	5	16.7	46.7	31.7	100%

- How can you tell that this table holds table percentages (rather than row or column percentages)?
- What was the most common genre/rating combination in 2005 movies?
- How many of these movies were PG-rated comedies?
- How many were G-rated?
- An editorial about the movies noted, "More than three-quarters of the movies made today can be seen only by patrons 13 years old or older." Does this table support that assertion? Explain.

- 21. Seniors.** Prior to graduation, a high school class was surveyed about its plans. The following table displays the results for white and minority students (the "Minority"

group included African-American, Asian, Hispanic, and Native American students):

Seniors		White	Minority
Plans	4-year college	198	44
	2-year college	36	6
	Military	4	1
	Employment	14	3
	Other	16	3

- What percent of the seniors are white?
 - What percent of the seniors are planning to attend a 2-year college?
 - What percent of the seniors are white and planning to attend a 2-year college?
 - What percent of the white seniors are planning to attend a 2-year college?
 - What percent of the seniors planning to attend a 2-year college are white?
- 22. Politics.** Students in an Intro Stats course were asked to describe their politics as "Liberal," "Moderate," or "Conservative." Here are the results:

Politics		L	M	C	Total
Sex	Female	35	36	6	77
	Male	50	44	21	115
	Total	85	80	27	192

- What percent of the class is male?
 - What percent of the class considers themselves to be "Conservative"?
 - What percent of the males in the class consider themselves to be "Conservative"?
 - What percent of all students in the class are males who consider themselves to be "Conservative"?
- 23. More about seniors.** Look again at the table of post-graduation plans for the senior class in Exercise 21.
- Find the conditional distributions (percentages) of plans for the white students.
 - Find the conditional distributions (percentages) of plans for the minority students.
 - Create a graph comparing the plans of white and minority students.
 - Do you see any important differences in the post-graduation plans of white and minority students? Write a brief summary of what these data show, including comparisons of conditional distributions.
- 24. Politics revisited.** Look again at the table of political views for the Intro Stats students in Exercise 22.
- Find the conditional distributions (percentages) of political views for the females.
 - Find the conditional distributions (percentages) of political views for the males.
 - Make a graphical display that compares the two distributions.
 - Do the variables *Politics* and *Sex* appear to be independent? Explain.

25. **Magnet schools revisited.** The *Chance* magazine article described in Exercise 9 further examined the impact of an applicant's ethnicity on the likelihood of admission to the Houston Independent School District's magnet schools programs. Those data are summarized in the table below:

		Admission Decision			
		Accepted	Wait-listed	Turned away	Total
Ethnicity	Black/Hispanic	485	0	32	517
	Asian	110	49	133	292
	White	336	251	359	946
	Total	931	300	524	1755

- a) What percent of all applicants were Asian?
 b) What percent of the students accepted were Asian?
 c) What percent of Asians were accepted?
 d) What percent of all students were accepted?
26. **More politics.** Look once more at the table summarizing the political views of Intro Stats students in Exercise 22.
- a) Produce a graphical display comparing the conditional distributions of males and females among the three categories of politics.
 b) Comment briefly on what you see from the display in a.
27. **Back to school.** Examine the table about ethnicity and acceptance for the Houston Independent School District's magnet schools program, shown in Exercise 25. Does it appear that the admissions decisions are made independent of the applicant's ethnicity? Explain.
28. **Cars.** A survey of autos parked in student and staff lots at a large university classified the brands by country of origin, as seen in the table.

		Driver	
		Student	Staff
Origin	American	107	105
	European	33	12
	Asian	55	47

- a) What percent of all the cars surveyed were foreign?
 b) What percent of the American cars were owned by students?
 c) What percent of the students owned American cars?
 d) What is the marginal distribution of origin?
 e) What are the conditional distributions of origin by driver classification?
 f) Do you think that the origin of the car is independent of the type of driver? Explain.
29. **Weather forecasts.** Just how accurate are the weather forecasts we hear every day? The following table compares the daily forecast with a city's actual weather for a year:

		Actual Weather	
		Rain	No rain
Forecast	Rain	27	63
	No rain	7	268

- a) On what percent of days did it actually rain?
 b) On what percent of days was rain predicted?
 c) What percent of the time was the forecast correct?
 d) Do you see evidence of an association between the type of weather and the ability of forecasters to make an accurate prediction? Write a brief explanation, including an appropriate graph.

30. **Twins.** In 2000, the *Journal of the American Medical Association (JAMA)* published a study that examined pregnancies that resulted in the birth of twins. Births were classified as preterm with intervention (induced labor or cesarean), preterm without procedures, or term/post-term. Researchers also classified the pregnancies by the level of prenatal medical care the mother received (inadequate, adequate, or intensive). The data, from the years 1995–1997, are summarized in the table below. Figures are in thousands of births. (*JAMA* 284 [2000]:335–341)

TWIN BIRTHS 1995–1997 (IN THOUSANDS)					
Level of Prenatal Care		Preterm (induced or cesarean)	Preterm (without procedures)	Term or post-term	Total
	Intensive	18	15	28	61
	Adequate	46	43	65	154
	Inadequate	12	13	38	63
	Total	76	71	131	278

- a) What percent of these mothers received inadequate medical care during their pregnancies?
 b) What percent of all twin births were preterm?
 c) Among the mothers who received inadequate medical care, what percent of the twin births were preterm?
 d) Create an appropriate graph comparing the outcomes of these pregnancies by the level of medical care the mother received.
 e) Write a few sentences describing the association between these two variables.
31. **Blood pressure.** A company held a blood pressure screening clinic for its employees. The results are summarized in the table below by age group and blood pressure level:

		Age		
		Under 30	30–49	Over 50
Blood Pressure	Low	27	37	31
	Normal	48	91	93
	High	23	51	73

- Find the marginal distribution of blood pressure level.
- Find the conditional distribution of blood pressure level within each age group.
- Compare these distributions with a segmented bar graph.
- Write a brief description of the association between age and blood pressure among these employees.
- Does this prove that people's blood pressure increases as they age? Explain.

32. **Obesity and exercise.** The Centers for Disease Control and Prevention (CDC) has estimated that 19.8% of Americans over 15 years old are obese. The CDC conducts a survey on obesity and various behaviors. Here is a table on self-reported exercise classified by body mass index (BMI):

		Body Mass Index		
		Normal (%)	Overweight (%)	Obese (%)
Physical Activity	Inactive	23.8	26.0	35.6
	Irregularly active	27.8	28.7	28.1
	Regular, not intense	31.6	31.1	27.2
	Regular, intense	16.8	14.2	9.1

- Are these percentages column percentages, row percentages, or table percentages?
 - Use graphical displays to show different percentages of physical activities for the three BMI groups.
 - Do these data prove that lack of exercise causes obesity? Explain.
33. **Anorexia.** Hearing anecdotal reports that some patients undergoing treatment for the eating disorder anorexia seemed to be responding positively to the antidepressant Prozac, medical researchers conducted an experiment to investigate. They found 93 women being treated for anorexia who volunteered to participate. For one year, 49 randomly selected patients were treated with Prozac and the other 44 were given an inert substance called a placebo. At the end of the year, patients were diagnosed as healthy or relapsed, as summarized in the table:

	Prozac	Placebo	Total
Healthy	35	32	67
Relapse	14	12	26
Total	49	44	93

Do these results provide evidence that Prozac might be helpful in treating anorexia? Explain.

34. **Antidepressants and bone fractures.** For a period of five years, physicians at McGill University Health Center followed more than 5000 adults over the age of 50. The

researchers were investigating whether people taking a certain class of antidepressants (SSRIs) might be at greater risk of bone fractures. Their observations are summarized in the table:

	Taking SSRI	No SSRI	Total
Experienced fractures	14	244	258
No fractures	123	4627	4750
Total	137	4871	5008

Do these results suggest there's an association between taking SSRI antidepressants and experiencing bone fractures? Explain.

35. **Drivers' licenses 2005.** The following table shows the number of licensed U.S. drivers by age and by sex (www.dot.gov):

Age	Male Drivers (number)	Female Drivers (number)	Total
19 and under	4,777,694	4,553,946	9,331,640
20–24	8,611,161	8,398,879	17,010,040
25–29	8,879,476	8,666,701	17,546,177
30–34	9,262,713	8,997,662	18,260,375
35–39	9,848,050	9,576,301	19,424,351
40–44	10,617,456	10,484,149	21,101,605
45–49	10,492,876	10,482,479	20,975,355
50–54	9,420,619	9,475,882	18,896,501
55–59	8,218,264	8,265,775	16,484,039
60–64	6,103,732	6,147,569	12,251,361
65–69	4,571,157	4,643,913	9,215,070
70–74	3,617,908	3,761,039	7,378,947
75–79	2,890,155	3,192,408	6,082,563
80–84	1,907,743	2,222,412	4,130,155
85 and over	1,170,817	1,406,271	2,577,088
Total	100,389,881	100,275,386	200,665,267

- What percent of total drivers are under 20?
 - What percent of total drivers are male?
 - Write a few sentences comparing the number of male and female licensed drivers in each age group.
 - Do a driver's age and sex appear to be independent? Explain?
36. **Tattoos.** A study by the University of Texas Southwestern Medical Center examined 626 people to see if an increased risk of contracting hepatitis C was associated with having a tattoo. If the subject had a tattoo, researchers asked whether it had been done in a commercial tattoo parlor or elsewhere. Write a brief description of the association between tattooing and hepatitis C, including an appropriate graphical display.

	Tattoo done in commercial parlor	Tattoo done elsewhere	No tattoo
Has hepatitis C	17	8	18
No hepatitis C	35	53	495

37. **Hospitals.** Most patients who undergo surgery make routine recoveries and are discharged as planned. Others suffer excessive bleeding, infection, or other postsurgical complications and have their discharges from the hospital delayed. Suppose your city has a large hospital and a small hospital, each performing major and minor surgeries. You collect data to see how many surgical patients have their discharges delayed by postsurgical complications, and you find the results shown in the following table.

	Discharge Delayed	
	Large hospital	Small hospital
Major surgery	120 of 800	10 of 50
Minor surgery	10 of 200	20 of 250

- Overall, for what percent of patients was discharge delayed?
 - Were the percentages different for major and minor surgery?
 - Overall, what were the discharge delay rates at each hospital?
 - What were the delay rates at each hospital for each kind of surgery?
 - The small hospital advertises that it has a lower rate of postsurgical complications. Do you agree?
 - Explain, in your own words, why this confusion occurs.
38. **Delivery service.** A company must decide which of two delivery services it will contract with. During a recent trial period, the company shipped numerous packages with each service and kept track of how often deliveries did not arrive on time. Here are the data:

Delivery Service	Type of Service	Number of Deliveries	Number of Late Packages
Pack Rats	Regular	400	12
	Overnight	100	16
Boxes R Us	Regular	100	2
	Overnight	400	28

- Compare the two services' overall percentage of late deliveries.
- On the basis of the results in part a, the company has decided to hire Pack Rats. Do you agree that Pack Rats delivers on time more often? Explain.
- The results here are an instance of what phenomenon?

39. **Graduate admissions.** A 1975 article in the magazine *Science* examined the graduate admissions process at Berkeley for evidence of sex discrimination. The table below shows the number of applicants accepted to each of four graduate programs:

		Males accepted (of applicants)	Females accepted (of applicants)
Program	1	511 of 825	89 of 108
	2	352 of 560	17 of 25
	3	137 of 407	132 of 375
	4	22 of 373	24 of 341
	Total	1022 of 2165	262 of 849

- What percent of total applicants were admitted?
 - Overall, was a higher percentage of males or females admitted?
 - Compare the percentage of males and females admitted in each program.
 - Which of the comparisons you made do you consider to be the most valid? Why?
40. **Be a Simpson!** Can you design a Simpson's paradox? Two companies are vying for a city's "Best Local Employer" award, to be given to the company most committed to hiring local residents. Although both employers hired 300 new people in the past year, Company A brags that it deserves the award because 70% of its new jobs went to local residents, compared to only 60% for Company B. Company B concedes that those percentages are correct, but points out that most of its new jobs were full-time, while most of Company A's were part-time. Not only that, says Company B, but a higher percentage of its full-time jobs went to local residents than did Company A's, and the same was true for part-time jobs. Thus, Company B argues, it's a better local employer than Company A.
- Show how it's possible for Company B to fill a higher percentage of both full-time and part-time jobs with local residents, even though Company A hired more local residents overall.



JUST CHECKING Answers

- 50.0%
- 44.4%
- 25.0%
- 15.6% Blue, 56.3% Brown, 28.1% Green/Hazel/Other
- 18.8% Blue, 62.5% Brown, 18.8% Green/Hazel/Other
- 40% of the blue-eyed students are female, while 50% of all students are female.
- Since blue-eyed students appear less likely to be female, it seems that *Sex* and *Eye Color* may not be independent. (But the numbers are small.)