

# OpenIntro Statistics

## Fourth Edition

David Diez  
*Data Scientist*  
*OpenIntro*

Mine Çetinkaya-Rundel  
*Associate Professor of the Practice, Duke University*  
*Professional Educator, RStudio*

Christopher D Barr  
*Investment Analyst*  
*Varadero Capital*

Editions 1, 2, and 3 can be found in the book's extra files,  
which also include tablet-friendly versions of some editions.

# Chapter 3

---

## Probability

---

**3.1 Defining probability**

**3.2 Conditional probability**

**3.3 Sampling from a small population**

**3.4 Random variables**

**3.5 Continuous distributions**

---

Probability forms the foundation of statistics, and you're probably already aware of many of the ideas presented in this chapter. However, formalization of probability concepts is likely new for most readers.

While this chapter provides a theoretical foundation for the ideas in later chapters and provides a path to a deeper understanding, mastery of the concepts introduced in this chapter is not required for applying the methods introduced in the rest of this book.

---



---

For videos, slides, and other resources, please visit  
[www.openintro.org/os](http://www.openintro.org/os)

## 3.4 Random variables

It's often useful to model a process using what's called a **random variable**. Such a model allows us to apply a mathematical framework and statistical principles for better understanding and predicting outcomes in the real world.

### EXAMPLE 3.54

Two books are assigned for a statistics class: a textbook and its corresponding study guide. The university bookstore determined 20% of enrolled students do not buy either book, 55% buy the textbook only, and 25% buy both books, and these percentages are relatively constant from one term to another. If there are 100 students enrolled, how many books should the bookstore expect to sell to this class?

E

Around 20 students will not buy either book (0 books total), about 55 will buy one book (55 books total), and approximately 25 will buy two books (totaling 50 books for these 25 students). The bookstore should expect to sell about 105 books for this class.

### GUIDED PRACTICE 3.55

Would you be surprised if the bookstore sold slightly more or less than 105 books?<sup>49</sup>

G

### EXAMPLE 3.56

The textbook costs \$137 and the study guide \$33. How much revenue should the bookstore expect from this class of 100 students?

About 55 students will just buy a textbook, providing revenue of

$$\$137 \times 55 = \$7,535$$

E

The roughly 25 students who buy both the textbook and the study guide would pay a total of

$$(\$137 + \$33) \times 25 = \$170 \times 25 = \$4,250$$

Thus, the bookstore should expect to generate about  $\$7,535 + \$4,250 = \$11,785$  from these 100 students for this one class. However, there might be some *sampling variability* so the actual amount may differ by a little bit.

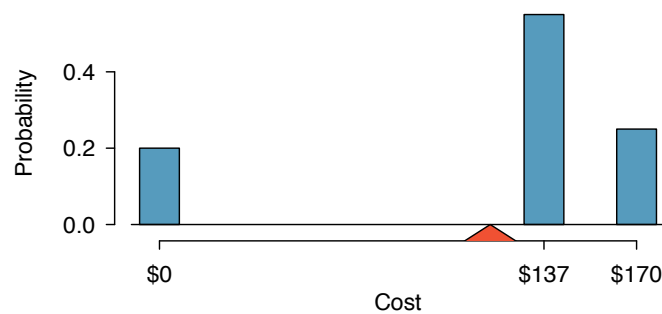


Figure 3.18: Probability distribution for the bookstore's revenue from one student. The triangle represents the average revenue per student.

<sup>49</sup>If they sell a little more or a little less, this should not be a surprise. Hopefully Chapter 1 helped make clear that there is natural variability in observed data. For example, if we would flip a coin 100 times, it will not usually come up heads exactly half the time, but it will probably be close.

**EXAMPLE 3.57**

What is the average revenue per student for this course?

E

The expected total revenue is \$11,785, and there are 100 students. Therefore the expected revenue per student is  $\$11,785/100 = \$117.85$ .

**3.4.1 Expectation**

We call a variable or process with a numerical outcome a **random variable**, and we usually represent this random variable with a capital letter such as  $X$ ,  $Y$ , or  $Z$ . The amount of money a single student will spend on her statistics books is a random variable, and we represent it by  $X$ .

**RANDOM VARIABLE**

A random process or variable with a numerical outcome.

The possible outcomes of  $X$  are labeled with a corresponding lower case letter  $x$  and subscripts. For example, we write  $x_1 = \$0$ ,  $x_2 = \$137$ , and  $x_3 = \$170$ , which occur with probabilities 0.20, 0.55, and 0.25. The distribution of  $X$  is summarized in Figure 3.18 and Figure 3.19.

$i$	1	2	3	Total
$x_i$	\$0	\$137	\$170	–
$P(X = x_i)$	0.20	0.55	0.25	1.00

Figure 3.19: The probability distribution for the random variable  $X$ , representing the bookstore's revenue from a single student.

We computed the average outcome of  $X$  as \$117.85 in Example 3.57. We call this average the **expected value** of  $X$ , denoted by  $E(X)$ . The expected value of a random variable is computed by adding each outcome weighted by its probability:

$$\begin{aligned} E(X) &= 0 \times P(X = 0) + 137 \times P(X = 137) + 170 \times P(X = 170) \\ &= 0 \times 0.20 + 137 \times 0.55 + 170 \times 0.25 = 117.85 \end{aligned}$$

**EXPECTED VALUE OF A DISCRETE RANDOM VARIABLE**

If  $X$  takes outcomes  $x_1, \dots, x_k$  with probabilities  $P(X = x_1), \dots, P(X = x_k)$ , the expected value of  $X$  is the sum of each outcome multiplied by its corresponding probability:

$$\begin{aligned} E(X) &= x_1 \times P(X = x_1) + \dots + x_k \times P(X = x_k) \\ &= \sum_{i=1}^k x_i P(X = x_i) \end{aligned}$$

The Greek letter  $\mu$  may be used in place of the notation  $E(X)$ .

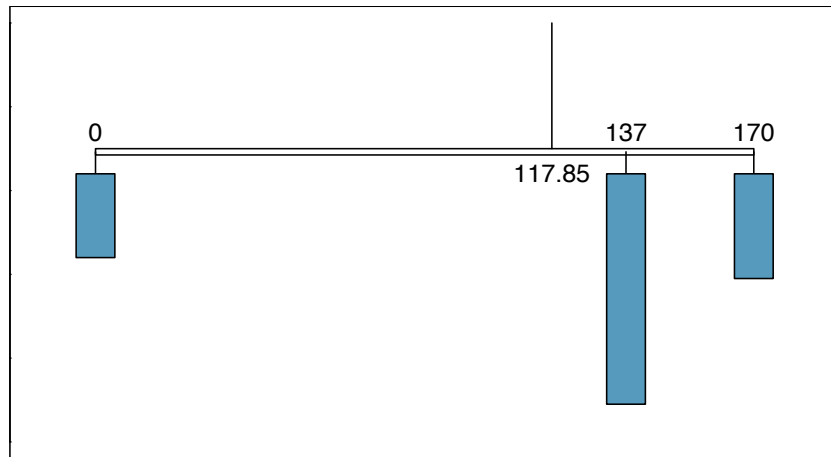


Figure 3.20: A weight system representing the probability distribution for  $X$ . The string holds the distribution at the mean to keep the system balanced.

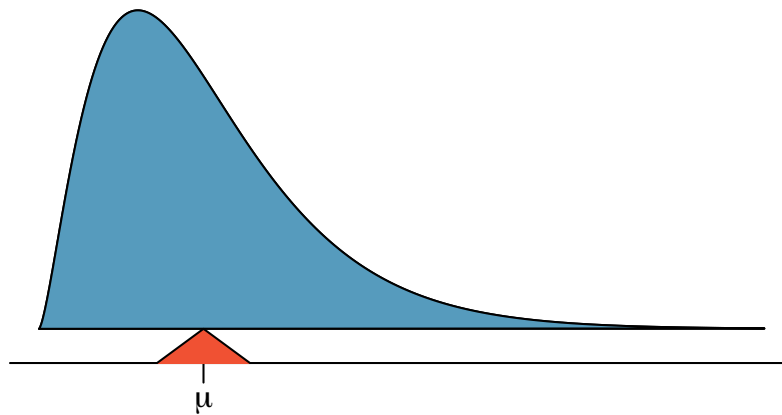


Figure 3.21: A continuous distribution can also be balanced at its mean.

The expected value for a random variable represents the average outcome. For example,  $E(X) = 117.85$  represents the average amount the bookstore expects to make from a single student, which we could also write as  $\mu = 117.85$ .

It is also possible to compute the expected value of a continuous random variable (see Section 3.5). However, it requires a little calculus and we save it for a later class.<sup>50</sup>

In physics, the expectation holds the same meaning as the center of gravity. The distribution can be represented by a series of weights at each outcome, and the mean represents the balancing point. This is represented in Figures 3.18 and 3.20. The idea of a center of gravity also expands to continuous probability distributions. Figure 3.21 shows a continuous probability distribution balanced atop a wedge placed at the mean.

<sup>50</sup> $\mu = \int xf(x)dx$  where  $f(x)$  represents a function for the density curve.

### 3.4.2 Variability in random variables

Suppose you ran the university bookstore. Besides how much revenue you expect to generate, you might also want to know the volatility (variability) in your revenue.

The variance and standard deviation can be used to describe the variability of a random variable. Section 2.1.4 introduced a method for finding the variance and standard deviation for a data set. We first computed deviations from the mean ( $x_i - \mu$ ), squared those deviations, and took an average to get the variance. In the case of a random variable, we again compute squared deviations. However, we take their sum weighted by their corresponding probabilities, just like we did for the expectation. This weighted sum of squared deviations equals the variance, and we calculate the standard deviation by taking the square root of the variance, just as we did in Section 2.1.4.

#### GENERAL VARIANCE FORMULA

If  $X$  takes outcomes  $x_1, \dots, x_k$  with probabilities  $P(X = x_1), \dots, P(X = x_k)$  and expected value  $\mu = E(X)$ , then the variance of  $X$ , denoted by  $Var(X)$  or the symbol  $\sigma^2$ , is

$$\begin{aligned}\sigma^2 &= (x_1 - \mu)^2 \times P(X = x_1) + \cdots \\ &\quad \cdots + (x_k - \mu)^2 \times P(X = x_k) \\ &= \sum_{j=1}^k (x_j - \mu)^2 P(X = x_j)\end{aligned}$$

The standard deviation of  $X$ , labeled  $\sigma$ , is the square root of the variance.

#### EXAMPLE 3.58

Compute the expected value, variance, and standard deviation of  $X$ , the revenue of a single statistics student for the bookstore.

It is useful to construct a table that holds computations for each outcome separately, then add up the results.

$i$	1	2	3	Total
$x_i$	\$0	\$137	\$170	
$P(X = x_i)$	0.20	0.55	0.25	
$x_i \times P(X = x_i)$	0	75.35	42.50	117.85

E

Thus, the expected value is  $\mu = 117.85$ , which we computed earlier. The variance can be constructed by extending this table:

$i$	1	2	3	Total
$x_i$	\$0	\$137	\$170	
$P(X = x_i)$	0.20	0.55	0.25	
$x_i \times P(X = x_i)$	0	75.35	42.50	117.85
$x_i - \mu$	-117.85	19.15	52.15	
$(x_i - \mu)^2$	13888.62	366.72	2719.62	
$(x_i - \mu)^2 \times P(X = x_i)$	2777.7	201.7	679.9	3659.3

The variance of  $X$  is  $\sigma^2 = 3659.3$ , which means the standard deviation is  $\sigma = \sqrt{3659.3} = \$60.49$ .

**GUIDED PRACTICE 3.59**

The bookstore also offers a chemistry textbook for \$159 and a book supplement for \$41. From past experience, they know about 25% of chemistry students just buy the textbook while 60% buy both the textbook and supplement.<sup>51</sup>

- (a) What proportion of students don't buy either book? Assume no students buy the supplement without the textbook.
- (b) Let  $Y$  represent the revenue from a single student. Write out the probability distribution of  $Y$ , i.e. a table for each outcome and its associated probability.
- (c) Compute the expected revenue from a single chemistry student.
- (d) Find the standard deviation to describe the variability associated with the revenue from a single student.

**3.4.3 Linear combinations of random variables**

So far, we have thought of each variable as being a complete story in and of itself. Sometimes it is more appropriate to use a combination of variables. For instance, the amount of time a person spends commuting to work each week can be broken down into several daily commutes. Similarly, the total gain or loss in a stock portfolio is the sum of the gains and losses in its components.

**EXAMPLE 3.60**

John travels to work five days a week. We will use  $X_1$  to represent his travel time on Monday,  $X_2$  to represent his travel time on Tuesday, and so on. Write an equation using  $X_1, \dots, X_5$  that represents his travel time for the week, denoted by  $W$ .

His total weekly travel time is the sum of the five daily values:

$$W = X_1 + X_2 + X_3 + X_4 + X_5$$

Breaking the weekly travel time  $W$  into pieces provides a framework for understanding each source of randomness and is useful for modeling  $W$ .

<sup>51</sup>(a)  $100\% - 25\% - 60\% = 15\%$  of students do not buy any books for the class. Part (b) is represented by the first two lines in the table below. The expectation for part (c) is given as the total on the line  $y_i \times P(Y = y_i)$ . The result of part (d) is the square-root of the variance listed on in the total on the last line:  $\sigma = \sqrt{Var(Y)} = \$69.28$ .

$i$ (scenario)	1 (noBook)	2 (textbook)	3 (both)	Total
$y_i$	0.00	159.00	200.00	
$P(Y = y_i)$	0.15	0.25	0.60	
$y_i \times P(Y = y_i)$	0.00	39.75	120.00	$E(Y) = 159.75$
$y_i - E(Y)$	-159.75	-0.75	40.25	
$(y_i - E(Y))^2$	25520.06	0.56	1620.06	
$(y_i - E(Y))^2 \times P(Y)$	3828.0	0.1	972.0	$Var(Y) \approx 4800$



**EXAMPLE 3.61**

It takes John an average of 18 minutes each day to commute to work. What would you expect his average commute time to be for the week?

We were told that the average (i.e. expected value) of the commute time is 18 minutes per day:  $E(X_i) = 18$ . To get the expected time for the sum of the five days, we can add up the expected time for each individual day:

E

$$\begin{aligned} E(W) &= E(X_1 + X_2 + X_3 + X_4 + X_5) \\ &= E(X_1) + E(X_2) + E(X_3) + E(X_4) + E(X_5) \\ &= 18 + 18 + 18 + 18 + 18 = 90 \text{ minutes} \end{aligned}$$

The expectation of the total time is equal to the sum of the expected individual times. More generally, the expectation of a sum of random variables is always the sum of the expectation for each random variable.

**GUIDED PRACTICE 3.62**

G

Elena is selling a TV at a cash auction and also intends to buy a toaster oven in the auction. If  $X$  represents the profit for selling the TV and  $Y$  represents the cost of the toaster oven, write an equation that represents the net change in Elena's cash.<sup>52</sup>

**GUIDED PRACTICE 3.63**

G

Based on past auctions, Elena figures she should expect to make about \$175 on the TV and pay about \$23 for the toaster oven. In total, how much should she expect to make or spend?<sup>53</sup>

**GUIDED PRACTICE 3.64**

G

Would you be surprised if John's weekly commute wasn't exactly 90 minutes or if Elena didn't make exactly \$152? Explain.<sup>54</sup>

Two important concepts concerning combinations of random variables have so far been introduced. First, a final value can sometimes be described as the sum of its parts in an equation. Second, intuition suggests that putting the individual average values into this equation gives the average value we would expect in total. This second point needs clarification – it is guaranteed to be true in what are called *linear combinations of random variables*.

A **linear combination** of two random variables  $X$  and  $Y$  is a fancy phrase to describe a combination

$$aX + bY$$

where  $a$  and  $b$  are some fixed and known numbers. For John's commute time, there were five random variables – one for each work day – and each random variable could be written as having a fixed coefficient of 1:

$$1X_1 + 1X_2 + 1X_3 + 1X_4 + 1X_5$$

For Elena's net gain or loss, the  $X$  random variable had a coefficient of +1 and the  $Y$  random variable had a coefficient of -1.

<sup>52</sup>She will make  $X$  dollars on the TV but spend  $Y$  dollars on the toaster oven:  $X - Y$ .

<sup>53</sup> $E(X - Y) = E(X) - E(Y) = 175 - 23 = \$152$ . She should expect to make about \$152.

<sup>54</sup>No, since there is probably some variability. For example, the traffic will vary from one day to next, and auction prices will vary depending on the quality of the merchandise and the interest of the attendees.

When considering the average of a linear combination of random variables, it is safe to plug in the mean of each random variable and then compute the final result. For a few examples of nonlinear combinations of random variables – cases where we cannot simply plug in the means – see the footnote.<sup>55</sup>

#### LINEAR COMBINATIONS OF RANDOM VARIABLES AND THE AVERAGE RESULT

If  $X$  and  $Y$  are random variables, then a linear combination of the random variables is given by

$$aX + bY$$

where  $a$  and  $b$  are some fixed numbers. To compute the average value of a linear combination of random variables, plug in the average of each individual random variable and compute the result:

$$a \times E(X) + b \times E(Y)$$

Recall that the expected value is the same as the mean, e.g.  $E(X) = \mu_X$ .

#### EXAMPLE 3.65

Leonard has invested \$6000 in Caterpillar Inc (stock ticker: CAT) and \$2000 in Exxon Mobil Corp (XOM). If  $X$  represents the change in Caterpillar's stock next month and  $Y$  represents the change in Exxon Mobil's stock next month, write an equation that describes how much money will be made or lost in Leonard's stocks for the month.

E

For simplicity, we will suppose  $X$  and  $Y$  are not in percents but are in decimal form (e.g. if Caterpillar's stock increases 1%, then  $X = 0.01$ ; or if it loses 1%, then  $X = -0.01$ ). Then we can write an equation for Leonard's gain as

$$\$6000 \times X + \$2000 \times Y$$

If we plug in the change in the stock value for  $X$  and  $Y$ , this equation gives the change in value of Leonard's stock portfolio for the month. A positive value represents a gain, and a negative value represents a loss.

#### GUIDED PRACTICE 3.66

G

Caterpillar stock has recently been rising at 2.0% and Exxon Mobil's at 0.2% per month, respectively. Compute the expected change in Leonard's stock portfolio for next month.<sup>56</sup>

#### GUIDED PRACTICE 3.67

G

You should have found that Leonard expects a positive gain in Guided Practice 3.66. However, would you be surprised if he actually had a loss this month?<sup>57</sup>

<sup>55</sup>If  $X$  and  $Y$  are random variables, consider the following combinations:  $X^{1+Y}$ ,  $X \times Y$ ,  $X/Y$ . In such cases, plugging in the average value for each random variable and computing the result will not generally lead to an accurate average value for the end result.

<sup>56</sup> $E(\$6000 \times X + \$2000 \times Y) = \$6000 \times 0.020 + \$2000 \times 0.002 = \$124$ .

<sup>57</sup>No. While stocks tend to rise over time, they are often volatile in the short term.

### 3.4.4 Variability in linear combinations of random variables

Quantifying the average outcome from a linear combination of random variables is helpful, but it is also important to have some sense of the uncertainty associated with the total outcome of that combination of random variables. The expected net gain or loss of Leonard's stock portfolio was considered in Guided Practice 3.66. However, there was no quantitative discussion of the volatility of this portfolio. For instance, while the average monthly gain might be about \$124 according to the data, that gain is not guaranteed. Figure 3.22 shows the monthly changes in a portfolio like Leonard's during a three year period. The gains and losses vary widely, and quantifying these fluctuations is important when investing in stocks.

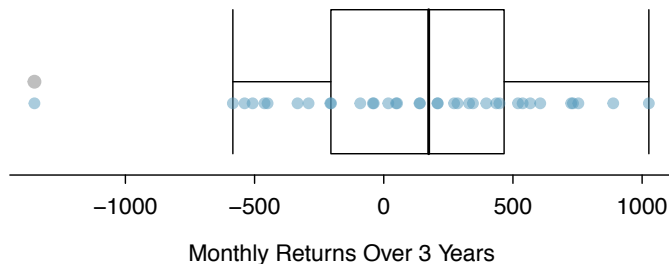


Figure 3.22: The change in a portfolio like Leonard's for 36 months, where \$6000 is in Caterpillar's stock and \$2000 is in Exxon Mobil's.

Just as we have done in many previous cases, we use the variance and standard deviation to describe the uncertainty associated with Leonard's monthly returns. To do so, the variances of each stock's monthly return will be useful, and these are shown in Figure 3.23. The stocks' returns are nearly independent.

Here we use an equation from probability theory to describe the uncertainty of Leonard's monthly returns; we leave the proof of this method to a dedicated probability course. The variance of a linear combination of random variables can be computed by plugging in the variances of the individual random variables and squaring the coefficients of the random variables:

$$\text{Var}(aX + bY) = a^2 \times \text{Var}(X) + b^2 \times \text{Var}(Y)$$

It is important to note that this equality assumes the random variables are independent; if independence doesn't hold, then a modification to this equation would be required that we leave as a topic for a future course to cover. This equation can be used to compute the variance of Leonard's monthly return:

$$\begin{aligned} \text{Var}(6000 \times X + 2000 \times Y) &= 6000^2 \times \text{Var}(X) + 2000^2 \times \text{Var}(Y) \\ &= 36,000,000 \times 0.0057 + 4,000,000 \times 0.0021 \\ &\approx 213,600 \end{aligned}$$

The standard deviation is computed as the square root of the variance:  $\sqrt{213,600} = \$463$ . While an average monthly return of \$124 on an \$8000 investment is nothing to scoff at, the monthly returns are so volatile that Leonard should not expect this income to be very stable.

	Mean ( $\bar{x}$ )	Standard deviation ( $s$ )	Variance ( $s^2$ )
CAT	0.0204	0.0757	0.0057
XOM	0.0025	0.0455	0.0021

Figure 3.23: The mean, standard deviation, and variance of the CAT and XOM stocks. These statistics were estimated from historical stock data, so notation used for sample statistics has been used.

**VARIABILITY OF LINEAR COMBINATIONS OF RANDOM VARIABLES**

The variance of a linear combination of random variables may be computed by squaring the constants, substituting in the variances for the random variables, and computing the result:

$$\text{Var}(aX + bY) = a^2 \times \text{Var}(X) + b^2 \times \text{Var}(Y)$$

This equation is valid as long as the random variables are independent of each other. The standard deviation of the linear combination may be found by taking the square root of the variance.

**EXAMPLE 3.68**

Suppose John's daily commute has a standard deviation of 4 minutes. What is the uncertainty in his total commute time for the week?

The expression for John's commute time was

$$X_1 + X_2 + X_3 + X_4 + X_5$$

E

Each coefficient is 1, and the variance of each day's time is  $4^2 = 16$ . Thus, the variance of the total weekly commute time is

$$\text{variance} = 1^2 \times 16 + 1^2 \times 16 + 1^2 \times 16 + 1^2 \times 16 + 1^2 \times 16 = 5 \times 16 = 80$$

$$\text{standard deviation} = \sqrt{\text{variance}} = \sqrt{80} = 8.94$$

The standard deviation for John's weekly work commute time is about 9 minutes.

**GUIDED PRACTICE 3.69**

G

The computation in Example 3.68 relied on an important assumption: the commute time for each day is independent of the time on other days of that week. Do you think this is valid? Explain.<sup>58</sup>

**GUIDED PRACTICE 3.70**

G

Consider Elena's two auctions from Guided Practice 3.62 on page 120. Suppose these auctions are approximately independent and the variability in auction prices associated with the TV and toaster oven can be described using standard deviations of \$25 and \$8. Compute the standard deviation of Elena's net gain.<sup>59</sup>

Consider again Guided Practice 3.70. The negative coefficient for  $Y$  in the linear combination was eliminated when we squared the coefficients. This generally holds true: negatives in a linear combination will have no impact on the variability computed for a linear combination, but they do impact the expected value computations.

<sup>58</sup>One concern is whether traffic patterns tend to have a weekly cycle (e.g. Fridays may be worse than other days). If that is the case, and John drives, then the assumption is probably not reasonable. However, if John walks to work, then his commute is probably not affected by any weekly traffic cycle.

<sup>59</sup>The equation for Elena can be written as

$$(1) \times X + (-1) \times Y$$

The variances of  $X$  and  $Y$  are 625 and 64. We square the coefficients and plug in the variances:

$$(1)^2 \times \text{Var}(X) + (-1)^2 \times \text{Var}(Y) = 1 \times 625 + 1 \times 64 = 689$$

The variance of the linear combination is 689, and the standard deviation is the square root of 689: about \$26.25.

---

## Exercises

**3.29 College smokers.** At a university, 13% of students smoke.

- Calculate the expected number of smokers in a random sample of 100 students from this university.
- The university gym opens at 9 am on Saturday mornings. One Saturday morning at 8:55 am there are 27 students outside the gym waiting for it to open. Should you use the same approach from part (a) to calculate the expected number of smokers among these 27 students?

**3.30 Ace of clubs wins.** Consider the following card game with a well-shuffled deck of cards. If you draw a red card, you win nothing. If you get a spade, you win \$5. For any club, you win \$10 plus an extra \$20 for the ace of clubs.

- Create a probability model for the amount you win at this game. Also, find the expected winnings for a single game and the standard deviation of the winnings.
- What is the maximum amount you would be willing to pay to play this game? Explain your reasoning.

**3.31 Hearts win.** In a new card game, you start with a well-shuffled full deck and draw 3 cards without replacement. If you draw 3 hearts, you win \$50. If you draw 3 black cards, you win \$25. For any other draws, you win nothing.

- Create a probability model for the amount you win at this game, and find the expected winnings. Also compute the standard deviation of this distribution.
- If the game costs \$5 to play, what would be the expected value and standard deviation of the net profit (or loss)? (*Hint: profit = winnings - cost;  $X - 5$* )
- If the game costs \$5 to play, should you play this game? Explain.

**3.32 Is it worth it?** Andy is always looking for ways to make money fast. Lately, he has been trying to make money by gambling. Here is the game he is considering playing: The game costs \$2 to play. He draws a card from a deck. If he gets a number card (2-10), he wins nothing. For any face card (jack, queen or king), he wins \$3. For any ace, he wins \$5, and he wins an *extra* \$20 if he draws the ace of clubs.

- Create a probability model and find Andy's expected profit per game.
- Would you recommend this game to Andy as a good way to make money? Explain.

**3.33 Portfolio return.** A portfolio's value increases by 18% during a financial boom and by 9% during normal times. It decreases by 12% during a recession. What is the expected return on this portfolio if each scenario is equally likely?

**3.34 Baggage fees.** An airline charges the following baggage fees: \$25 for the first bag and \$35 for the second. Suppose 54% of passengers have no checked luggage, 34% have one piece of checked luggage and 12% have two pieces. We suppose a negligible portion of people check more than two bags.

- Build a probability model, compute the average revenue per passenger, and compute the corresponding standard deviation.
- About how much revenue should the airline expect for a flight of 120 passengers? With what standard deviation? Note any assumptions you make and if you think they are justified.

**3.35 American roulette.** The game of American roulette involves spinning a wheel with 38 slots: 18 red, 18 black, and 2 green. A ball is spun onto the wheel and will eventually land in a slot, where each slot has an equal chance of capturing the ball. Gamblers can place bets on red or black. If the ball lands on their color, they double their money. If it lands on another color, they lose their money. Suppose you bet \$1 on red. What's the expected value and standard deviation of your winnings?

**3.36 European roulette.** The game of European roulette involves spinning a wheel with 37 slots: 18 red, 18 black, and 1 green. A ball is spun onto the wheel and will eventually land in a slot, where each slot has an equal chance of capturing the ball. Gamblers can place bets on red or black. If the ball lands on their color, they double their money. If it lands on another color, they lose their money.

- Suppose you play roulette and bet \$3 on a single round. What is the expected value and standard deviation of your total winnings?
- Suppose you bet \$1 in three different rounds. What is the expected value and standard deviation of your total winnings?
- How do your answers to parts (a) and (b) compare? What does this say about the riskiness of the two games?

## 3.5 Continuous distributions

So far in this chapter we've discussed cases where the outcome of a variable is discrete. In this section, we consider a context where the outcome is a continuous numerical variable.

### EXAMPLE 3.71

Figure 3.24 shows a few different hollow histograms for the heights of US adults. How does changing the number of bins allow you to make different interpretations of the data?

E

Adding more bins provides greater detail. This sample is extremely large, which is why much smaller bins still work well. Usually we do not use so many bins with smaller sample sizes since small counts per bin mean the bin heights are very volatile.

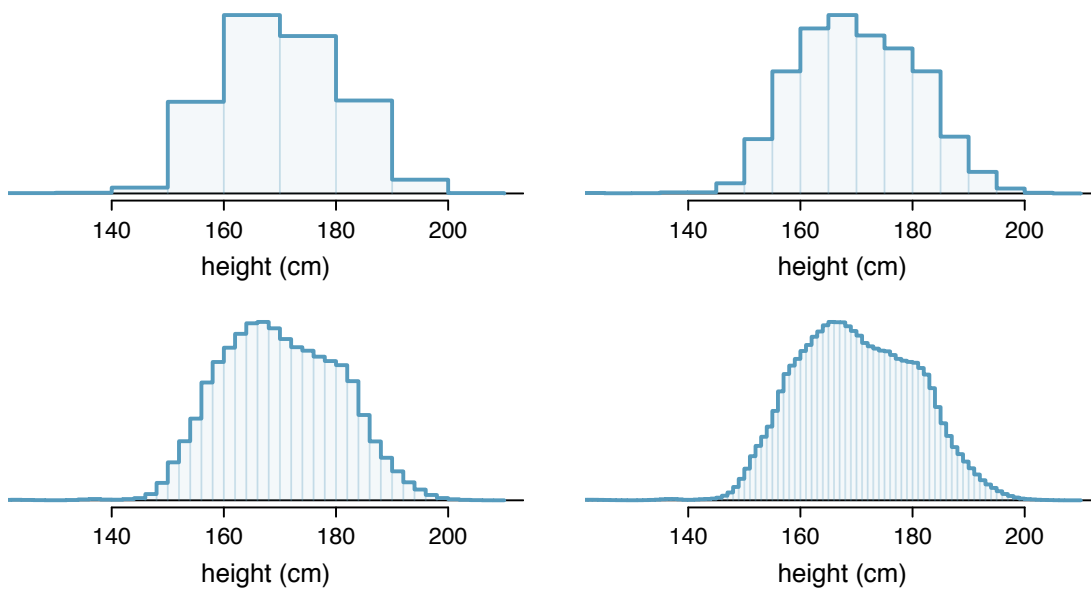


Figure 3.24: Four hollow histograms of US adults heights with varying bin widths.

### EXAMPLE 3.72

What proportion of the sample is between 180 cm and 185 cm tall (about 5'11" to 6'1")?

We can add up the heights of the bins in the range 180 cm and 185 and divide by the sample size. For instance, this can be done with the two shaded bins shown in Figure 3.25. The two bins in this region have counts of 195,307 and 156,239 people, resulting in the following estimate of the probability:

E

$$\frac{195307 + 156239}{3,000,000} = 0.1172$$

This fraction is the same as the proportion of the histogram's area that falls in the range 180 to 185 cm.

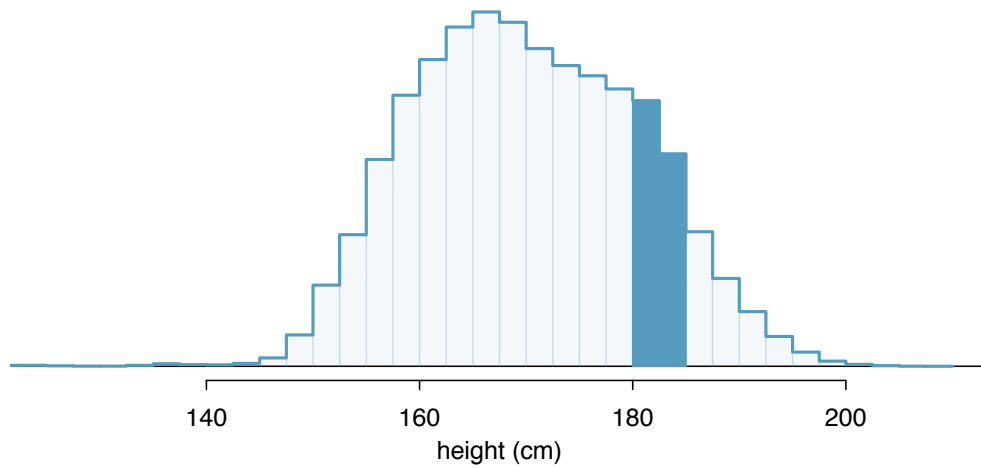


Figure 3.25: A histogram with bin sizes of 2.5 cm. The shaded region represents individuals with heights between 180 and 185 cm.

### 3.5.1 From histograms to continuous distributions

Examine the transition from a boxy hollow histogram in the top-left of Figure 3.24 to the much smoother plot in the lower-right. In this last plot, the bins are so slim that the hollow histogram is starting to resemble a smooth curve. This suggests the population height as a *continuous* numerical variable might best be explained by a curve that represents the outline of extremely slim bins.

This smooth curve represents a **probability density function** (also called a **density** or **distribution**), and such a curve is shown in Figure 3.26 overlaid on a histogram of the sample. A density has a special property: the total area under the density's curve is 1.

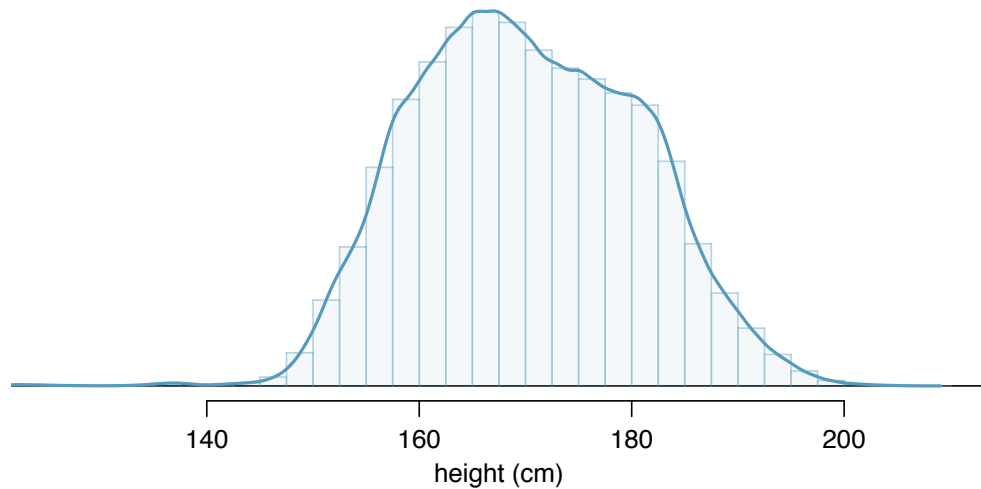


Figure 3.26: The continuous probability distribution of heights for US adults.

### 3.5.2 Probabilities from continuous distributions

We computed the proportion of individuals with heights 180 to 185 cm in Example 3.72 as a fraction:

$$\frac{\text{number of people between 180 and 185}}{\text{total sample size}}$$

We found the number of people with heights between 180 and 185 cm by determining the fraction of the histogram's area in this region. Similarly, we can use the area in the shaded region under the curve to find a probability (with the help of a computer):

$$P(\text{height between 180 and 185}) = \text{area between 180 and 185} = 0.1157$$

The probability that a randomly selected person is between 180 and 185 cm is 0.1157. This is very close to the estimate from Example 3.72: 0.1172.

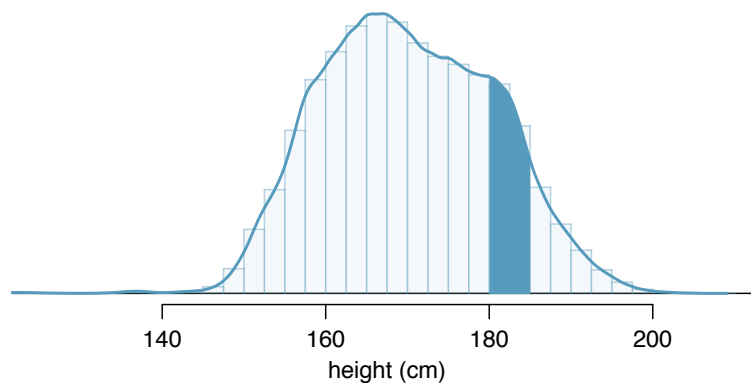


Figure 3.27: Density for heights in the US adult population with the area between 180 and 185 cm shaded. Compare this plot with Figure 3.25.

#### GUIDED PRACTICE 3.73

Three US adults are randomly selected. The probability a single adult is between 180 and 185 cm is 0.1157.<sup>60</sup>

- What is the probability that all three are between 180 and 185 cm tall?
- What is the probability that none are between 180 and 185 cm?

#### EXAMPLE 3.74

What is the probability that a randomly selected person is **exactly** 180 cm? Assume you can measure perfectly.

This probability is zero. A person might be close to 180 cm, but not exactly 180 cm tall. This also makes sense with the definition of probability as area; there is no area captured between 180 cm and 180 cm.

#### GUIDED PRACTICE 3.75

Suppose a person's height is rounded to the nearest centimeter. Is there a chance that a random person's **measured** height will be 180 cm?<sup>61</sup>

<sup>60</sup>Brief answers: (a)  $0.1157 \times 0.1157 \times 0.1157 = 0.0015$ . (b)  $(1 - 0.1157)^3 = 0.692$

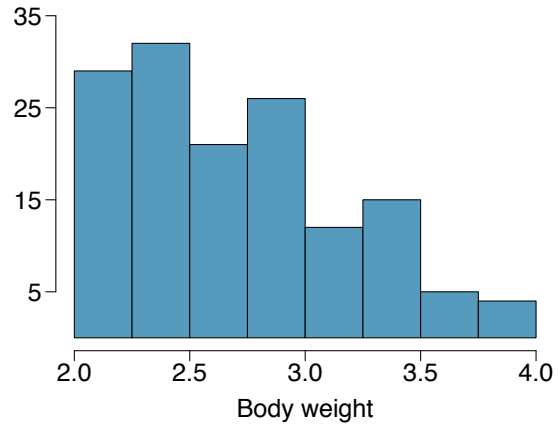
<sup>61</sup>This has positive probability. Anyone between 179.5 cm and 180.5 cm will have a *measured* height of 180 cm. This is probably a more realistic scenario to encounter in practice versus Example 3.74.



## Exercises

**3.37 Cat weights.** The histogram shown below represents the weights (in kg) of 47 female and 97 male cats.<sup>62</sup>

- What fraction of these cats weigh less than 2.5 kg?
- What fraction of these cats weigh between 2.5 and 2.75 kg?
- What fraction of these cats weigh between 2.75 and 3.5 kg?



**3.38 Income and gender.** The relative frequency table below displays the distribution of annual total personal income (in 2009 inflation-adjusted dollars) for a representative sample of 96,420,486 Americans. These data come from the American Community Survey for 2005-2009. This sample is comprised of 59% males and 41% females.<sup>63</sup>

- Describe the distribution of total personal income.
- What is the probability that a randomly chosen US resident makes less than \$50,000 per year?
- What is the probability that a randomly chosen US resident makes less than \$50,000 per year and is female? Note any assumptions you make.
- The same data source indicates that 71.8% of females make less than \$50,000 per year. Use this value to determine whether or not the assumption you made in part (c) is valid.

<i>Income</i>	<i>Total</i>
\$1 to \$9,999 or loss	2.2%
\$10,000 to \$14,999	4.7%
\$15,000 to \$24,999	15.8%
\$25,000 to \$34,999	18.3%
\$35,000 to \$49,999	21.2%
\$50,000 to \$64,999	13.9%
\$65,000 to \$74,999	5.8%
\$75,000 to \$99,999	8.4%
\$100,000 or more	9.7%

<sup>62</sup>W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Fourth Edition. www.stats.ox.ac.uk/pub/MASS4. New York: Springer, 2002.

<sup>63</sup>U.S. Census Bureau, 2005-2009 American Community Survey.