

## Chapter 4

# Foundations for inference

Statistical inference is concerned primarily with understanding the quality of parameter estimates. For example, a classic inferential question is, “How sure are we that the estimated mean,  $\bar{x}$ , is near the true population mean,  $\mu$ ?” While the equations and details change depending on the setting, the foundations for inference are the same throughout all of statistics. We introduce these common themes in Sections 4.1-4.4 by discussing inference about the population mean,  $\mu$ , and set the stage for other parameters and scenarios in Section 4.5. Some advanced considerations are discussed in Section 4.6. Understanding this chapter will make the rest of this book, and indeed the rest of statistics, seem much more familiar.

Throughout the next few sections we consider a data set called `run10`, which represents all 16,924 runners who finished the 2012 Cherry Blossom 10 mile run in Washington, DC.<sup>1</sup> Part of this data set is shown in Table 4.1, and the variables are described in Table 4.2.

| ID       | time     | age      | gender   | state    |
|----------|----------|----------|----------|----------|
| 1        | 92.25    | 38.00    | M        | MD       |
| 2        | 106.35   | 33.00    | M        | DC       |
| 3        | 89.33    | 55.00    | F        | VA       |
| 4        | 113.50   | 24.00    | F        | VA       |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| 16923    | 122.87   | 37.00    | F        | VA       |
| 16924    | 93.30    | 27.00    | F        | DC       |

Table 4.1: Six observations from the `run10` data set.

| variable            | description                                |
|---------------------|--|
| <code>time</code>   | Ten mile run time, in minutes              |
| <code>age</code>    | Age, in years                              |
| <code>gender</code> | Gender (M for male, F for female)          |
| <code>state</code>  | Home state (or country if not from the US) |

Table 4.2: Variables and their descriptions for the `run10` data set.

---

<sup>1</sup><http://www.cherryblossom.org>

| ID       | time     | age      | gender   | state    |
|----------|----------|----------|----------|----------|
| 1983     | 88.31    | 59       | M        | MD       |
| 8192     | 100.67   | 32       | M        | VA       |
| 11020    | 109.52   | 33       | F        | VA       |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| 1287     | 89.49    | 26       | M        | DC       |

Table 4.3: Four observations for the `run10Samp` data set, which represents a simple random sample of 100 runners from the 2012 Cherry Blossom Run.

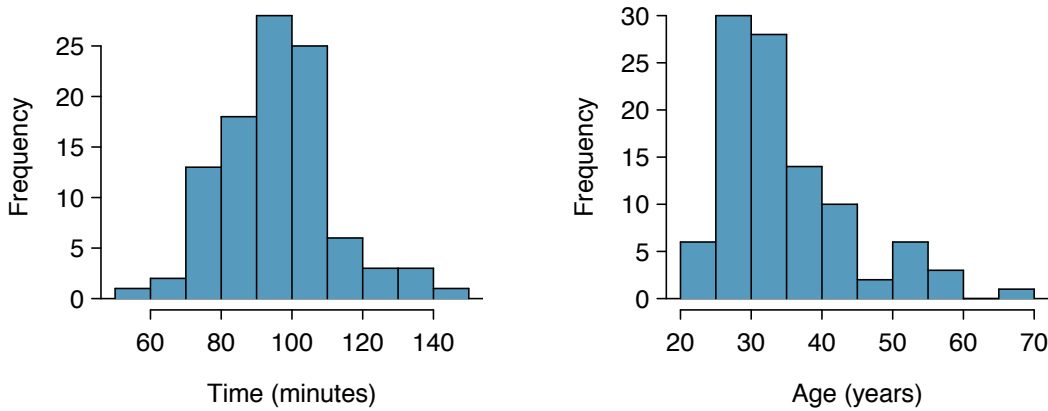


Figure 4.4: Histograms of `time` and `age` for the sample Cherry Blossom Run data. The average time is in the mid-90s, and the average age is in the mid-to-upper 30s. The age distribution is moderately skewed to the right.

These data are special because they include the results for the entire population of runners who finished the 2012 Cherry Blossom Run. We took a simple random sample of this population, which is represented in Table 4.3. We will use this sample, which we refer to as the `run10Samp` data set, to draw conclusions about the entire population. This is the practice of statistical inference in the broadest sense. Two histograms summarizing the time and age variables in the `run10Samp` data set are shown in Figure 4.4.

## 4.1 Variability in estimates

We would like to estimate two features of the Cherry Blossom runners using the sample.

- (1) How long does it take a runner, on average, to complete the 10 miles?
- (2) What is the average age of the runners?

These questions may be informative for planning the Cherry Blossom Run in future years.<sup>2</sup> We will use  $x_1, \dots, x_{100}$  to represent the 10 mile time for each runner in our sample, and  $y_1, \dots, y_{100}$  will represent the age of each of these participants.

<sup>2</sup>While we focus on the mean in this chapter, questions regarding variation are often just as important in practice. For instance, we would plan an event very differently if the standard deviation of runner age was 2 versus if it was 20.

### 4.1.1 Point estimates

We want to estimate the **population mean** based on the sample. The most intuitive way to go about doing this is to simply take the **sample mean**. That is, to estimate the average 10 mile run time of all participants, take the average time for the sample:

$$\bar{x} = \frac{88.22 + 100.58 + \cdots + 89.40}{100} = 95.61$$

The sample mean  $\bar{x} = 95.61$  minutes is called a **point estimate** of the population mean: if we can only choose one value to estimate the population mean, this is our best guess. Suppose we take a new sample of 100 people and recompute the mean; we will probably not get the exact same answer that we got using the `run10Samp` data set. Estimates generally vary from one sample to another, and this **sampling variation** suggests our estimate may be close, but it will not be exactly equal to the parameter.

We can also estimate the average age of participants by examining the sample mean of **age**:

$$\bar{y} = \frac{59 + 32 + \cdots + 26}{100} = 35.05$$

What about generating point estimates of other **population parameters**, such as the population median or population standard deviation? Once again we might estimate parameters based on sample statistics, as shown in Table 4.5. For example, we estimate the population standard deviation for the running time using the sample standard deviation, 15.78 minutes.

| time     | estimate | parameter |
|----------|----------|-----------|
| mean     | 95.61    | 94.52     |
| median   | 95.37    | 94.03     |
| st. dev. | 15.78    | 15.93     |

Table 4.5: Point estimates and parameter values for the `time` variable.

- ⊙ **Exercise 4.1** Suppose we want to estimate the difference in run times for men and women. If  $\bar{x}_{men} = 87.65$  and  $\bar{x}_{women} = 102.13$ , then what would be a good point estimate for the population difference?<sup>3</sup>
- ⊙ **Exercise 4.2** If you had to provide a point estimate of the population IQR for the run time of participants, how might you make such an estimate using a sample?<sup>4</sup>

### 4.1.2 Point estimates are not exact

Estimates are usually not exactly equal to the truth, but they get better as more data become available. We can see this by plotting a running mean from our `run10Samp` sample. A **running mean** is a sequence of means, where each mean uses one more observation in its calculation than the mean directly before it in the sequence. For example, the second mean in the sequence is the average of the first two observations and the third in the

<sup>3</sup>We could take the difference of the two sample means:  $102.13 - 87.65 = 14.48$ . Men ran about 14.48 minutes faster on average in the 2012 Cherry Blossom Run.

<sup>4</sup>To obtain a point estimate of the IQR for the population, we could take the IQR of the sample.

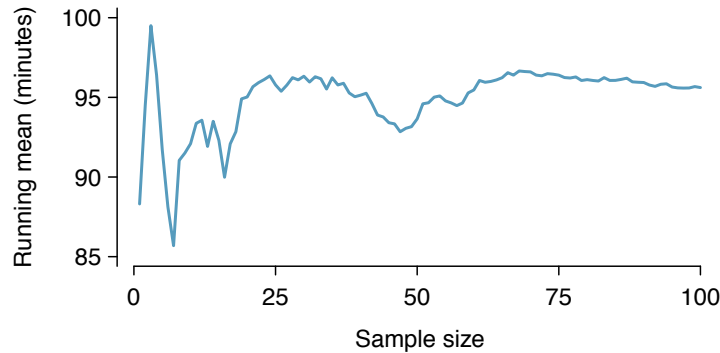


Figure 4.6: The mean computed after adding each individual to the sample. The mean tends to approach the true population average as more data become available.

sequence is the average of the first three. The running mean for the 10 mile run time in the `run10Samp` data set is shown in Figure 4.6, and it approaches the true population average, 94.52 minutes, as more data become available.

Sample point estimates only approximate the population parameter, and they vary from one sample to another. If we took another simple random sample of the Cherry Blossom runners, we would find that the sample mean for the run time would be a little different. It will be useful to quantify how variable an estimate is from one sample to another. If this variability is small (i.e. the sample mean doesn't change much from one sample to another) then that estimate is probably very accurate. If it varies widely from one sample to another, then we should not expect our estimate to be very good.

### 4.1.3 Standard error of the mean

From the random sample represented in `run10Samp`, we guessed the average time it takes to run 10 miles is 95.61 minutes. Suppose we take another random sample of 100 individuals and take its mean: 95.30 minutes. Suppose we took another (93.43 minutes) and another (94.16 minutes), and so on. If we do this many many times – which we can do only because we have the entire population data set – we can build up a **sampling distribution** for the sample mean when the sample size is 100, shown in Figure 4.7.

#### Sampling distribution

The sampling distribution represents the distribution of the point estimates based on samples of a fixed size from a certain population. It is useful to think of a particular point estimate as being drawn from such a distribution. Understanding the concept of a sampling distribution is central to understanding statistical inference.

The sampling distribution shown in Figure 4.7 is unimodal and approximately symmetric. It is also centered exactly at the true population mean:  $\mu = 94.52$ . Intuitively, this makes sense. The sample means should tend to “fall around” the population mean.

We can see that the sample mean has some variability around the population mean, which can be quantified using the standard deviation of this distribution of sample means:  $\sigma_{\bar{x}} = 1.59$ . The standard deviation of the sample mean tells us how far the typical estimate

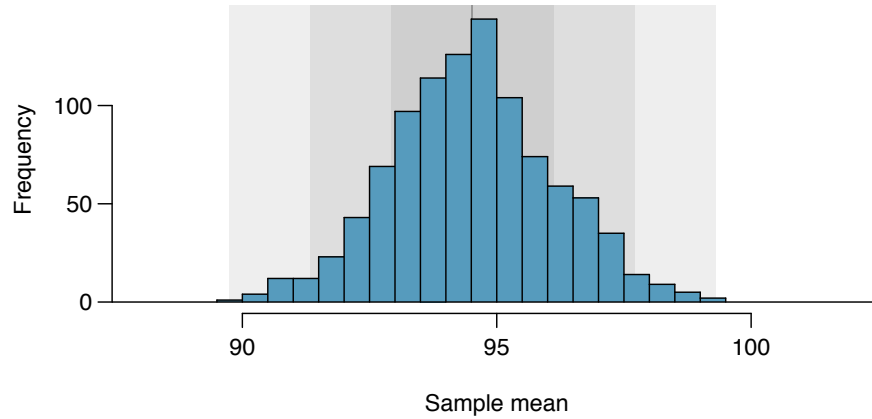


Figure 4.7: A histogram of 1000 sample means for run time, where the samples are of size  $n = 100$ .

is away from the actual population mean, 94.52 minutes. It also describes the typical **error** of the point estimate, and for this reason we usually call this standard deviation the **standard error (SE)** of the estimate.

*SE*  
standard  
error

#### Standard error of an estimate

The standard deviation associated with an estimate is called the *standard error*. It describes the typical error or uncertainty associated with the estimate.

When considering the case of the point estimate  $\bar{x}$ , there is one problem: there is no obvious way to estimate its standard error from a single sample. However, statistical theory provides a helpful tool to address this issue.

- ⊙ **Exercise 4.3** (a) Would you rather use a small sample or a large sample when estimating a parameter? Why? (b) Using your reasoning from (a), would you expect a point estimate based on a small sample to have smaller or larger standard error than a point estimate based on a larger sample?<sup>5</sup>

In the sample of 100 runners, the standard error of the sample mean is equal to one-tenth of the population standard deviation:  $1.59 = 15.93/10$ . In other words, the standard error of the sample mean based on 100 observations is equal to

$$SE_{\bar{x}} = \sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} = \frac{15.93}{\sqrt{100}} = 1.59$$

where  $\sigma_x$  is the standard deviation of the individual observations. This is no coincidence. We can show mathematically that this equation is correct when the observations are independent using the probability tools of Section 2.4.

<sup>5</sup>(a) Consider two random samples: one of size 10 and one of size 1000. Individual observations in the small sample are highly influential on the estimate while in larger samples these individual observations would more often average each other out. The larger sample would tend to provide a more accurate estimate. (b) If we think an estimate is better, we probably mean it typically has less error. Based on (a), our intuition suggests that a larger sample size corresponds to a smaller standard error.

**Computing SE for the sample mean**

Given  $n$  independent observations from a population with standard deviation  $\sigma$ , the standard error of the sample mean is equal to

$$SE = \frac{\sigma}{\sqrt{n}} \quad (4.4)$$

A reliable method to ensure sample observations are independent is to conduct a simple random sample consisting of less than 10% of the population.

There is one subtle issue of Equation (4.4): the population standard deviation is typically unknown. You might have already guessed how to resolve this problem: we can use the point estimate of the standard deviation from the sample. This estimate tends to be sufficiently good when the sample size is at least 30 and the population distribution is not strongly skewed. Thus, we often just use the sample standard deviation  $s$  instead of  $\sigma$ . When the sample size is smaller than 30, we will need to use a method to account for extra uncertainty in the standard error. If the skew condition is not met, a larger sample is needed to compensate for the extra skew. These topics are further discussed in Section 4.4.

- ⊙ **Exercise 4.5** In the sample of 100 runners, the standard deviation of the runners' ages is  $s_y = 8.97$ . Because the sample is simple random and consists of less than 10% of the population, the observations are independent. (a) What is the standard error of the sample mean,  $\bar{y} = 35.05$  years? (b) Would you be surprised if someone told you the average age of all the runners was actually 36 years?<sup>6</sup>
- ⊙ **Exercise 4.6** (a) Would you be more trusting of a sample that has 100 observations or 400 observations? (b) We want to show mathematically that our estimate tends to be better when the sample size is larger. If the standard deviation of the individual observations is 10, what is our estimate of the standard error when the sample size is 100? What about when it is 400? (c) Explain how your answer to (b) mathematically justifies your intuition in part (a).<sup>7</sup>

### 4.1.4 Basic properties of point estimates

We achieved three goals in this section. First, we determined that point estimates from a sample may be used to estimate population parameters. We also determined that these point estimates are not exact: they vary from one sample to another. Lastly, we quantified the uncertainty of the sample mean using what we call the standard error, mathematically represented in Equation (4.4). While we could also quantify the standard error for other estimates – such as the median, standard deviation, or any other number of statistics – we will postpone these extensions until later chapters or courses.

<sup>6</sup>(a) Use Equation (4.4) with the sample standard deviation to compute the standard error:  $SE_{\bar{y}} = 8.97/\sqrt{100} = 0.90$  years. (b) It would not be surprising. Our sample is about 1 standard error from 36 years. In other words, 36 years old does not seem to be implausible given that our sample was relatively close to it. (We use the standard error to identify what is close.)

<sup>7</sup>(a) Extra observations are usually helpful in understanding the population, so a point estimate with 400 observations seems more trustworthy. (b) The standard error when the sample size is 100 is given by  $SE_{100} = 10/\sqrt{100} = 1$ . For 400:  $SE_{400} = 10/\sqrt{400} = 0.5$ . The larger sample has a smaller standard error. (c) The standard error of the sample with 400 observations is lower than that of the sample with 100 observations. The standard error describes the typical error, and since it is lower for the larger sample, this mathematically shows the estimate from the larger sample tends to be better – though it does not guarantee that every large sample will provide a better estimate than a particular small sample.

## 4.2 Confidence intervals

A point estimate provides a single plausible value for a parameter. However, a point estimate is rarely perfect; usually there is some error in the estimate. Instead of supplying just a point estimate of a parameter, a next logical step would be to provide a plausible *range of values* for the parameter.

In this section and in Section 4.3, we will emphasize the special case where the point estimate is a sample mean and the parameter is the population mean. In Section 4.5, we generalize these methods for a variety of point estimates and population parameters that we will encounter in Chapter 5 and beyond.

### 4.2.1 Capturing the population parameter

A plausible range of values for the population parameter is called a **confidence interval**.

Using only a point estimate is like fishing in a murky lake with a spear, and using a confidence interval is like fishing with a net. We can throw a spear where we saw a fish, but we will probably miss. On the other hand, if we toss a net in that area, we have a good chance of catching the fish.

If we report a point estimate, we probably will not hit the exact population parameter. On the other hand, if we report a range of plausible values – a confidence interval – we have a good shot at capturing the parameter.

- ⊙ **Exercise 4.7** If we want to be very certain we capture the population parameter, should we use a wider interval or a smaller interval?<sup>8</sup>

### 4.2.2 An approximate 95% confidence interval

Our point estimate is the most plausible value of the parameter, so it makes sense to build the confidence interval around the point estimate. The standard error, which is a measure of the uncertainty associated with the point estimate, provides a guide for how large we should make the confidence interval.

The standard error represents the standard deviation associated with the estimate, and roughly 95% of the time the estimate will be within 2 standard errors of the parameter. If the interval spreads out 2 standard errors from the point estimate, we can be roughly 95% **confident** that we have captured the true parameter:

$$\text{point estimate} \pm 2 \times SE \quad (4.8)$$

But what does “95% confident” mean? Suppose we took many samples and built a confidence interval from each sample using Equation (4.8). Then about 95% of those intervals would contain the actual mean,  $\mu$ . Figure 4.8 shows this process with 25 samples, where 24 of the resulting confidence intervals contain the average time for all the runners,  $\mu = 94.52$  minutes, and one does not.

- ⊙ **Exercise 4.9** In Figure 4.8, one interval does not contain 94.52 minutes. Does this imply that the mean cannot be 94.52? <sup>9</sup>

<sup>8</sup>If we want to be more certain we will capture the fish, we might use a wider net. Likewise, we use a wider confidence interval if we want to be more certain that we capture the parameter.

<sup>9</sup>Just as some observations occur more than 2 standard deviations from the mean, some point estimates will be more than 2 standard errors from the parameter. A confidence interval only provides a plausible range of values for a parameter. While we might say other values are implausible based on the data, this does not mean they are impossible.

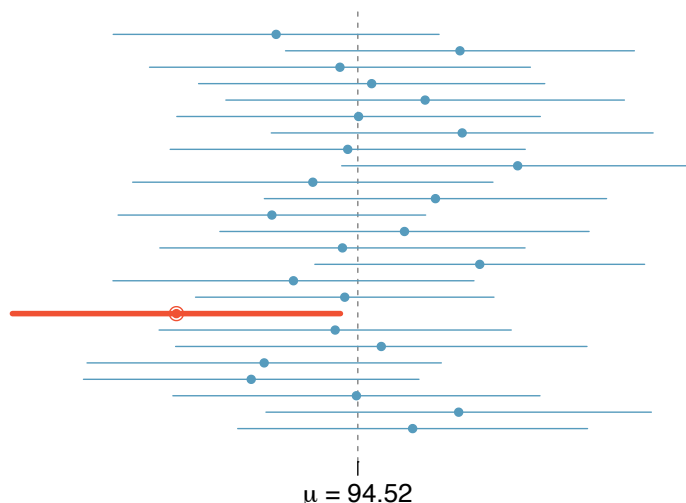


Figure 4.8: Twenty-five samples of size  $n = 100$  were taken from the `run10` data set. For each sample, a confidence interval was created to try to capture the average 10 mile time for the population. Only 1 of these 25 intervals did not capture the true mean,  $\mu = 94.52$  minutes.

The rule where about 95% of observations are within 2 standard deviations of the mean is only approximately true. However, it holds very well for the normal distribution. As we will soon see, the mean tends to be normally distributed when the sample size is sufficiently large.

- **Example 4.10** If the sample mean of times from `run10Samp` is 95.61 minutes and the standard error, as estimated using the sample standard deviation, is 1.58 minutes, what would be an approximate 95% confidence interval for the average 10 mile time of all runners in the race? Apply the standard error calculated using the sample standard deviation ( $SE = \frac{15.78}{\sqrt{100}} = 1.58$ ), which is how we usually proceed since the population standard deviation is generally unknown.

We apply Equation (4.8):

$$95.61 \pm 2 \times 1.58 \rightarrow (92.45, 98.77)$$

Based on these data, we are about 95% confident that the average 10 mile time for all runners in the race was larger than 92.45 but less than 98.77 minutes. Our interval extends out 2 standard errors from the point estimate,  $\bar{x}$ .

- ⊙ **Exercise 4.11** The sample data suggest the average runner's age is about 35.05 years with a standard error of 0.90 years (estimated using the sample standard deviation, 8.97). What is an approximate 95% confidence interval for the average age of all of the runners?<sup>10</sup>

<sup>10</sup>Again apply Equation (4.8):  $35.05 \pm 2 \times 0.90 \rightarrow (33.25, 36.85)$ . We interpret this interval as follows: We are about 95% confident the average age of all participants in the 2012 Cherry Blossom Run was between 33.25 and 36.85 years.



### 4.2.3 A sampling distribution for the mean

In Section 4.1.3, we introduced a sampling distribution for  $\bar{x}$ , the average run time for samples of size 100. We examined this distribution earlier in Figure 4.7. Now we'll take 100,000 samples, calculate the mean of each, and plot them in a histogram to get an especially accurate depiction of the sampling distribution. This histogram is shown in the left panel of Figure 4.9.

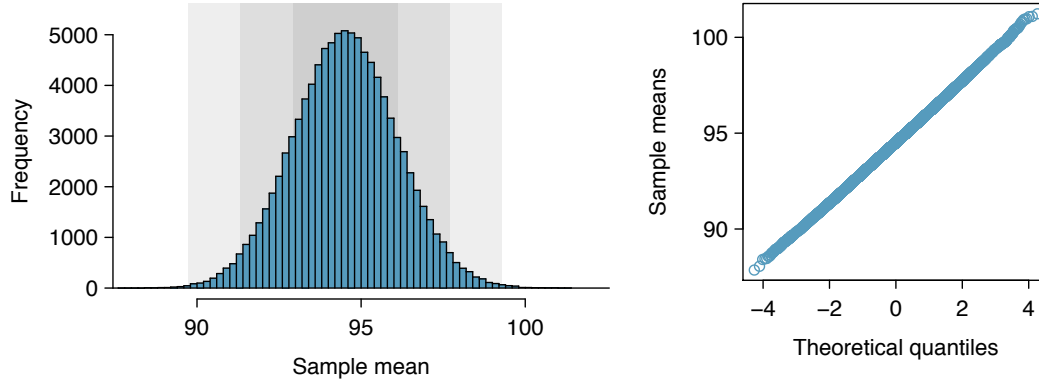


Figure 4.9: The left panel shows a histogram of the sample means for 100,000 different random samples. The right panel shows a normal probability plot of those sample means.

Does this distribution look familiar? Hopefully so! The distribution of sample means closely resembles the normal distribution (see Section 3.1). A normal probability plot of these sample means is shown in the right panel of Figure 4.9. Because all of the points closely fall around a straight line, we can conclude the distribution of sample means is nearly normal. This result can be explained by the Central Limit Theorem.

#### Central Limit Theorem, informal description

If a sample consists of at least 30 independent observations and the data are not strongly skewed, then the distribution of the sample mean is well approximated by a normal model.

We will apply this informal version of the Central Limit Theorem for now, and discuss its details further in Section 4.4.

The choice of using 2 standard errors in Equation (4.8) was based on our general guideline that roughly 95% of the time, observations are within two standard deviations of the mean. Under the normal model, we can make this more accurate by using 1.96 in place of 2.

$$\text{point estimate} \pm 1.96 \times SE \quad (4.12)$$

If a point estimate, such as  $\bar{x}$ , is associated with a normal model and standard error  $SE$ , then we use this more precise 95% confidence interval.

### 4.2.4 Changing the confidence level

Suppose we want to consider confidence intervals where the confidence level is somewhat higher than 95%: perhaps we would like a confidence level of 99%. Think back to the analogy about trying to catch a fish: if we want to be more sure that we will catch the fish, we should use a wider net. To create a 99% confidence level, we must also widen our 95% interval. On the other hand, if we want an interval with lower confidence, such as 90%, we could make our original 95% interval slightly slimmer.

The 95% confidence interval structure provides guidance in how to make intervals with new confidence levels. Below is a general 95% confidence interval for a point estimate that comes from a nearly normal distribution:

$$\text{point estimate} \pm 1.96 \times SE \quad (4.13)$$

There are three components to this interval: the point estimate, “1.96”, and the standard error. The choice of  $1.96 \times SE$  was based on capturing 95% of the data since the estimate is within 1.96 standard deviations of the parameter about 95% of the time. The choice of 1.96 corresponds to a 95% confidence level.

⊙ **Exercise 4.14** If  $X$  is a normally distributed random variable, how often will  $X$  be within 2.58 standard deviations of the mean?<sup>11</sup>

To create a 99% confidence interval, change 1.96 in the 95% confidence interval formula to be 2.58. Exercise 4.14 highlights that 99% of the time a normal random variable will be within 2.58 standard deviations of the mean. This approach – using the  $Z$  scores in the normal model to compute confidence levels – is appropriate when  $\bar{x}$  is associated with a normal distribution with mean  $\mu$  and standard deviation  $SE_{\bar{x}}$ . Thus, the formula for a 99% confidence interval is

$$\bar{x} \pm 2.58 \times SE_{\bar{x}} \quad (4.15)$$

The normal approximation is crucial to the precision of these confidence intervals. Section 4.4 provides a more detailed discussion about when the normal model can safely be applied. When the normal model is not a good fit, we will use alternative distributions that better characterize the sampling distribution.

#### Conditions for $\bar{x}$ being nearly normal and $SE$ being accurate

Important conditions to help ensure the sampling distribution of  $\bar{x}$  is nearly normal and the estimate of  $SE$  sufficiently accurate:

- The sample observations are independent.
- The sample size is large:  $n \geq 30$  is a good rule of thumb.
- The distribution of sample observations is not strongly skewed.

Additionally, the larger the sample size, the more lenient we can be with the sample's skew.

<sup>11</sup>This is equivalent to asking how often the  $Z$  score will be larger than -2.58 but less than 2.58. (For a picture, see Figure 4.10.) To determine this probability, look up -2.58 and 2.58 in the normal probability table (0.0049 and 0.9951). Thus, there is a  $0.9951 - 0.0049 \approx 0.99$  probability that the unobserved random variable  $X$  will be within 2.58 standard deviations of  $\mu$ .

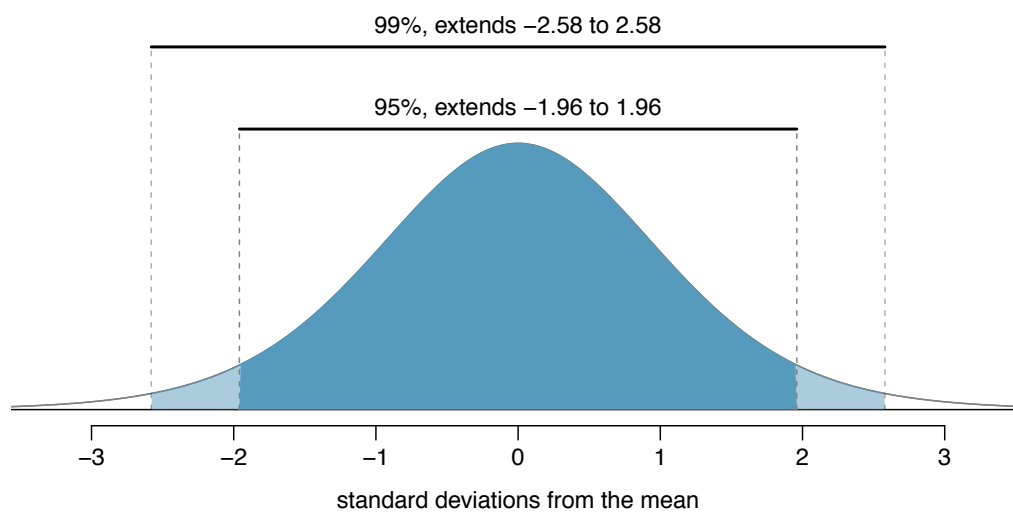


Figure 4.10: The area between  $-z^*$  and  $z^*$  increases as  $|z^*|$  becomes larger. If the confidence level is 99%, we choose  $z^*$  such that 99% of the normal curve is between  $-z^*$  and  $z^*$ , which corresponds to 0.5% in the lower tail and 0.5% in the upper tail:  $z^* = 2.58$ .

Verifying independence is often the most difficult of the conditions to check, and the way to check for independence varies from one situation to another. However, we can provide simple rules for the most common scenarios.

**TIP: How to verify sample observations are independent**

Observations in a simple random sample consisting of less than 10% of the population are independent.

**Caution: Independence for random processes and experiments**

If a sample is from a random process or experiment, it is important to verify the observations from the process or subjects in the experiment are nearly independent and maintain their independence throughout the process or experiment. Usually subjects are considered independent if they undergo random assignment in an experiment.

- ⊙ **Exercise 4.16** Create a 99% confidence interval for the average age of all runners in the 2012 Cherry Blossom Run. The point estimate is  $\bar{y} = 35.05$  and the standard error is  $SE_{\bar{y}} = 0.90$ .<sup>12</sup>

<sup>12</sup>The observations are independent (simple random sample, < 10% of the population), the sample size is at least 30 ( $n = 100$ ), and the distribution is only slightly skewed (Figure 4.4); the normal approximation and estimate of SE should be reasonable. Apply the 99% confidence interval formula:  $\bar{y} \pm 2.58 \times SE_{\bar{y}} \rightarrow (32.7, 37.4)$ . We are 99% confident that the average age of all runners is between 32.7 and 37.4 years.

**Confidence interval for any confidence level**

If the point estimate follows the normal model with standard error  $SE$ , then a confidence interval for the population parameter is

$$\text{point estimate} \pm z^* SE$$

where  $z^*$  corresponds to the confidence level selected.

Figure 4.10 provides a picture of how to identify  $z^*$  based on a confidence level. We select  $z^*$  so that the area between  $-z^*$  and  $z^*$  in the normal model corresponds to the confidence level.

**Margin of error**

In a confidence interval,  $z^* \times SE$  is called the **margin of error**.

- ⊙ **Exercise 4.17** Use the data in Exercise 4.16 to create a 90% confidence interval for the average age of all runners in the 2012 Cherry Blossom Run.<sup>13</sup>

**4.2.5 Interpreting confidence intervals**

A careful eye might have observed the somewhat awkward language used to describe confidence intervals. Correct interpretation:

We are XX% confident that the population parameter is between...

*Incorrect* language might try to describe the confidence interval as capturing the population parameter with a certain probability. This is one of the most common errors: while it might be useful to think of it as a probability, the confidence level only quantifies how plausible it is that the parameter is in the interval.

Another especially important consideration of confidence intervals is that they *only try to capture the population parameter*. Our intervals say nothing about the confidence of capturing individual observations, a proportion of the observations, or about capturing point estimates. Confidence intervals only attempt to capture population parameters.

**4.2.6 Nearly normal population with known SD (special topic)**

In rare circumstances we know important characteristics of a population. For instance, we might know a population is nearly normal and we may also know its parameter values. Even so, we may still like to study characteristics of a random sample from the population. Consider the conditions required for modeling a sample mean using the normal distribution:

- (1) The observations are independent.
- (2) The sample size  $n$  is at least 30.
- (3) The data distribution is not strongly skewed.

<sup>13</sup>We first find  $z^*$  such that 90% of the distribution falls between  $-z^*$  and  $z^*$  in the standard normal model,  $N(\mu = 0, \sigma = 1)$ . We can look up  $-z^*$  in the normal probability table by looking for a lower tail of 5% (the other 5% is in the upper tail), thus  $z^* = 1.65$ . The 90% confidence interval can then be computed as  $\bar{y} \pm 1.65 \times SE_{\bar{y}} \rightarrow (33.6, 36.5)$ . (We had already verified conditions for normality and the standard error.) That is, we are 90% confident the average age is larger than 33.6 but less than 36.5 years.

These conditions are required so we can adequately estimate the standard deviation and so we can ensure the distribution of sample means is nearly normal. However, if the population is known to be nearly normal, the sample mean is always nearly normal (this is a special case of the Central Limit Theorem). If the standard deviation is also known, then conditions (2) and (3) are not necessary for those data.

- **Example 4.18** The heights of male seniors in high school closely follow a normal distribution  $N(\mu = 70.43, \sigma = 2.73)$ , where the units are inches.<sup>14</sup> If we randomly sampled the heights of five male seniors, what distribution should the sample mean follow?

The population is nearly normal, the population standard deviation is known, and the heights represent a random sample from a much larger population, satisfying the independence condition. Therefore the sample mean of the heights will follow a nearly normal distribution with mean  $\mu = 70.43$  inches and standard error  $SE = \sigma/\sqrt{n} = 2.73/\sqrt{5} = 1.22$  inches.

**Alternative conditions for applying the normal distribution to model the sample mean**

If the population of cases is known to be nearly normal and the population standard deviation  $\sigma$  is known, then the sample mean  $\bar{x}$  will follow a nearly normal distribution  $N(\mu, \sigma/\sqrt{n})$  if the sampled observations are also independent.

Sometimes the mean changes over time but the standard deviation remains the same. In such cases, a sample mean of small but nearly normal observations paired with a known standard deviation can be used to produce a confidence interval for the current population mean using the normal distribution.

- **Example 4.19** Is there a connection between height and popularity in high school? Many students may suspect as much, but what do the data say? Suppose the top 5 nominees for prom king at a high school have an average height of 71.8 inches. Does this provide strong evidence that these seniors' heights are not representative of all male seniors at their high school?

If these five seniors are height-representative, then their heights should be like a random sample from the distribution given in Example 4.18,  $N(\mu = 70.43, \sigma = 2.73)$ , and the sample mean should follow  $N(\mu = 70.43, \sigma/\sqrt{n} = 1.22)$ . Formally we are conducting what is called a *hypothesis test*, which we will discuss in greater detail during the next section. We are weighing two possibilities:

$H_0$ : The prom king nominee heights are representative;  $\bar{x}$  will follow a normal distribution with mean 70.43 inches and standard error 1.22 inches.

$H_A$ : The heights are not representative; we suspect the mean height is different from 70.43 inches.

If there is strong evidence that the sample mean is not from the normal distribution provided in  $H_0$ , then that suggests the heights of prom king nominees are not a simple random sample (i.e.  $H_A$  is true). We can look at the Z score of the sample mean to

<sup>14</sup>These values were computed using the USDA Food Commodity Intake Database.

tell us how unusual our sample is. If  $H_0$  is true:

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{71.8 - 70.43}{1.22} = 1.12$$

A Z score of just 1.12 is not very unusual (we typically use a threshold of  $\pm 2$  to decide what is unusual), so there is not strong evidence against the claim that the heights are representative. This does not mean the heights are actually representative, only that this very small sample does not necessarily show otherwise.

**TIP: Relaxing the nearly normal condition**

As the sample size becomes larger, it is reasonable to *slowly* relax the nearly normal assumption on the data when dealing with small samples. By the time the sample size reaches 30, the data must show strong skew for us to be concerned about the normality of the sampling distribution.

## 4.3 Hypothesis testing

Is the typical US runner getting faster or slower over time? We consider this question in the context of the Cherry Blossom Run, comparing runners in 2006 and 2012. Technological advances in shoes, training, and diet might suggest runners would be faster in 2012. An opposing viewpoint might say that with the average body mass index on the rise, people tend to run slower. In fact, all of these components might be influencing run time.

In addition to considering run times in this section, we consider a topic near and dear to most students: sleep. A recent study found that college students average about 7 hours of sleep per night.<sup>15</sup> However, researchers at a rural college are interested in showing that their students sleep longer than seven hours on average. We investigate this topic in Section 4.3.4.

### 4.3.1 Hypothesis testing framework

The average time for all runners who finished the Cherry Blossom Run in 2006 was 93.29 minutes (93 minutes and about 17 seconds). We want to determine if the `run10Samp` data set provides strong evidence that the participants in 2012 were faster or slower than those runners in 2006, versus the other possibility that there has been no change.<sup>16</sup> We simplify these three options into two competing **hypotheses**:

$H_0$ : The average 10 mile run time was the same for 2006 and 2012.

$H_A$ : The average 10 mile run time for 2012 was *different* than that of 2006.

We call  $H_0$  the null hypothesis and  $H_A$  the alternative hypothesis.

$H_0$   
null hypothesis

$H_A$   
alternative  
hypothesis

**Null and alternative hypotheses**

The **null hypothesis** ( $H_0$ ) often represents either a skeptical perspective or a claim to be tested. The **alternative hypothesis** ( $H_A$ ) represents an alternative claim under consideration and is often represented by a range of possible parameter values.

<sup>15</sup><http://theloquitur.com/?p=1161>

<sup>16</sup>While we could answer this question by examining the entire population data (`run10`), we only consider the sample data (`run10Samp`), which is more realistic since we rarely have access to population data.

The null hypothesis often represents a skeptical position or a perspective of no difference. The alternative hypothesis often represents a new perspective, such as the possibility that there has been a change.

**TIP: Hypothesis testing framework**

The skeptic will not reject the null hypothesis ( $H_0$ ), unless the evidence in favor of the alternative hypothesis ( $H_A$ ) is so strong that she rejects  $H_0$  in favor of  $H_A$ .

The hypothesis testing framework is a very general tool, and we often use it without a second thought. If a person makes a somewhat unbelievable claim, we are initially skeptical. However, if there is sufficient evidence that supports the claim, we set aside our skepticism and reject the null hypothesis in favor of the alternative. The hallmarks of hypothesis testing are also found in the US court system.

- ⊙ **Exercise 4.20** A US court considers two possible claims about a defendant: she is either innocent or guilty. If we set these claims up in a hypothesis framework, which would be the null hypothesis and which the alternative?<sup>17</sup>

Jurors examine the evidence to see whether it convincingly shows a defendant is guilty. Even if the jurors leave unconvinced of guilt beyond a reasonable doubt, this does not mean they believe the defendant is innocent. This is also the case with hypothesis testing: *even if we fail to reject the null hypothesis, we typically do not accept the null hypothesis as true*. Failing to find strong evidence for the alternative hypothesis is not equivalent to accepting the null hypothesis.

In the example with the Cherry Blossom Run, the null hypothesis represents no difference in the average time from 2006 to 2012. The alternative hypothesis represents something new or more interesting: there was a difference, either an increase or a decrease. These hypotheses can be described in mathematical notation using  $\mu_{12}$  as the average run time for 2012:

$$H_0: \mu_{12} = 93.29$$

$$H_A: \mu_{12} \neq 93.29$$

where 93.29 minutes (93 minutes and about 17 seconds) is the average 10 mile time for all runners in the 2006 Cherry Blossom Run. Using this mathematical notation, the hypotheses can now be evaluated using statistical tools. We call 93.29 the **null value** since it represents the value of the parameter if the null hypothesis is true. We will use the `run10Samp` data set to evaluate the hypothesis test.

### 4.3.2 Testing hypotheses using confidence intervals

We can start the evaluation of the hypothesis setup by comparing 2006 and 2012 run times using a point estimate from the 2012 sample:  $\bar{x}_{12} = 95.61$  minutes. This estimate suggests the average time is actually longer than the 2006 time, 93.29 minutes. However, to evaluate whether this provides strong evidence that there has been a change, we must consider the uncertainty associated with  $\bar{x}_{12}$ .

<sup>17</sup>The jury considers whether the evidence is so convincing (strong) that there is no reasonable doubt regarding the person's guilt; in such a case, the jury rejects innocence (the null hypothesis) and concludes the defendant is guilty (alternative hypothesis).

We learned in Section 4.1 that there is fluctuation from one sample to another, and it is very unlikely that the sample mean will be exactly equal to our parameter; we should not expect  $\bar{x}_{12}$  to exactly equal  $\mu_{12}$ . Given that  $\bar{x}_{12} = 95.61$ , it might still be possible that the population average in 2012 has remained unchanged from 2006. The difference between  $\bar{x}_{12}$  and 93.29 could be due to *sampling variation*, i.e. the variability associated with the point estimate when we take a random sample.

In Section 4.2, confidence intervals were introduced as a way to find a range of plausible values for the population mean. Based on `run10Samp`, a 95% confidence interval for the 2012 population mean,  $\mu_{12}$ , was calculated as

$$(92.45, 98.77)$$

Because the 2006 mean, 93.29, falls in the range of plausible values, we cannot say the null hypothesis is implausible. That is, we failed to reject the null hypothesis,  $H_0$ .

**TIP: Double negatives can sometimes be used in statistics**

In many statistical explanations, we use double negatives. For instance, we might say that the null hypothesis is *not implausible* or we *failed to reject* the null hypothesis. Double negatives are used to communicate that while we are not rejecting a position, we are also not saying it is correct.

- **Example 4.21** Next consider whether there is strong evidence that the average age of runners has changed from 2006 to 2012 in the Cherry Blossom Run. In 2006, the average age was 36.13 years, and in the 2012 `run10Samp` data set, the average was 35.05 years with a standard deviation of 8.97 years for 100 runners.

First, set up the hypotheses:

$H_0$ : The average age of runners has not changed from 2006 to 2012,  $\mu_{age} = 36.13$ .

$H_A$ : The average age of runners has changed from 2006 to 2012,  $\mu_{age} \neq 36.13$ .

We have previously verified conditions for this data set. The normal model may be applied to  $\bar{y}$  and the estimate of  $SE$  should be very accurate. Using the sample mean and standard error, we can construct a 95% confidence interval for  $\mu_{age}$  to determine if there is sufficient evidence to reject  $H_0$ :

$$\bar{y} \pm 1.96 \times \frac{s}{\sqrt{100}} \rightarrow 35.05 \pm 1.96 \times 0.90 \rightarrow (33.29, 36.81)$$

This confidence interval contains the *null value*, 36.13. Because 36.13 is not implausible, we cannot reject the null hypothesis. We have not found strong evidence that the average age is different than 36.13 years.

- ⊙ **Exercise 4.22** Colleges frequently provide estimates of student expenses such as housing. A consultant hired by a community college claimed that the average student housing expense was \$650 per month. What are the null and alternative hypotheses to test whether this claim is accurate?<sup>18</sup>

<sup>18</sup> $H_0$ : The average cost is \$650 per month,  $\mu = \$650$ .

$H_A$ : The average cost is different than \$650 per month,  $\mu \neq \$650$ .



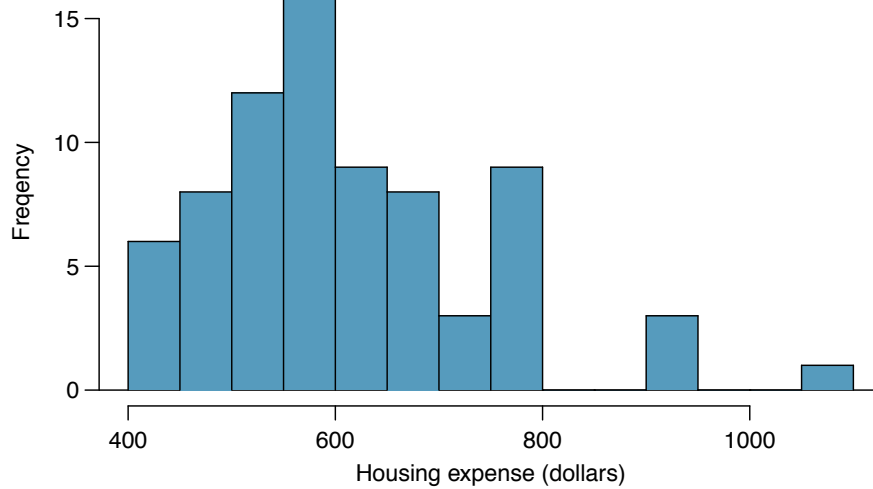


Figure 4.11: Sample distribution of student housing expense. These data are moderately skewed, roughly determined using the outliers on the right.

⊙ **Exercise 4.23** The community college decides to collect data to evaluate the \$650 per month claim. They take a random sample of 75 students at their school and obtain the data represented in Figure 4.11. Can we apply the normal model to the sample mean?<sup>19</sup>

● **Example 4.24** The sample mean for student housing is \$611.63 and the sample standard deviation is \$132.85. Construct a 95% confidence interval for the population mean and evaluate the hypotheses of Exercise 4.22.

The standard error associated with the mean may be estimated using the sample standard deviation divided by the square root of the sample size. Recall that  $n = 75$  students were sampled.

$$SE = \frac{s}{\sqrt{n}} = \frac{132.85}{\sqrt{75}} = 15.34$$

You showed in Exercise 4.23 that the normal model may be applied to the sample mean. This ensures a 95% confidence interval may be accurately constructed:

$$\bar{x} \pm z^* SE \rightarrow 611.63 \pm 1.96 \times 15.34 \rightarrow (581.56, 641.70)$$

Because the null value \$650 is not in the confidence interval, a true mean of \$650 is implausible and we reject the null hypothesis. The data provide statistically significant evidence that the actual average housing expense is less than \$650 per month.

<sup>19</sup>Applying the normal model requires that certain conditions are met. Because the data are a simple random sample and the sample (presumably) represents no more than 10% of all students at the college, the observations are independent. The sample size is also sufficiently large ( $n = 75$ ) and the data exhibit only moderate skew. Thus, the normal model may be applied to the sample mean.

### 4.3.3 Decision errors

Hypothesis tests are not flawless. Just think of the court system: innocent people are sometimes wrongly convicted and the guilty sometimes walk free. Similarly, we can make a wrong decision in statistical hypothesis tests. However, the difference is that we have the tools necessary to quantify how often we make such errors.

There are two competing hypotheses: the null and the alternative. In a hypothesis test, we make a statement about which one might be true, but we might choose incorrectly. There are four possible scenarios in a hypothesis test, which are summarized in Table 4.12.

| Truth | Test conclusion     |                                |
|-------|---------------------|--------------------------------|
|       | do not reject $H_0$ | reject $H_0$ in favor of $H_A$ |
|       | $H_0$ true          | $H_A$ true                     |
|       | okay                | Type 1 Error                   |
|       | Type 2 Error        | okay                           |

Table 4.12: Four different scenarios for hypothesis tests.

A **Type 1 Error** is rejecting the null hypothesis when  $H_0$  is actually true. A **Type 2 Error** is failing to reject the null hypothesis when the alternative is actually true.

- ⊙ **Exercise 4.25** In a US court, the defendant is either innocent ( $H_0$ ) or guilty ( $H_A$ ). What does a Type 1 Error represent in this context? What does a Type 2 Error represent? Table 4.12 may be useful.<sup>20</sup>
- ⊙ **Exercise 4.26** How could we reduce the Type 1 Error rate in US courts? What influence would this have on the Type 2 Error rate?<sup>21</sup>
- ⊙ **Exercise 4.27** How could we reduce the Type 2 Error rate in US courts? What influence would this have on the Type 1 Error rate?<sup>22</sup>

Exercises 4.25-4.27 provide an important lesson: if we reduce how often we make one type of error, we generally make more of the other type.

Hypothesis testing is built around rejecting or failing to reject the null hypothesis. That is, we do not reject  $H_0$  unless we have strong evidence. But what precisely does *strong evidence* mean? As a general rule of thumb, for those cases where the null hypothesis is actually true, we do not want to incorrectly reject  $H_0$  more than 5% of the time. This corresponds to a **significance level** of 0.05. We often write the significance level using  $\alpha$  (the Greek letter *alpha*):  $\alpha = 0.05$ . We discuss the appropriateness of different significance levels in Section 4.3.6.

If we use a 95% confidence interval to test a hypothesis where the null hypothesis is true, we will make an error whenever the point estimate is at least 1.96 standard errors

<sup>20</sup>If the court makes a Type 1 Error, this means the defendant is innocent ( $H_0$  true) but wrongly convicted. A Type 2 Error means the court failed to reject  $H_0$  (i.e. failed to convict the person) when she was in fact guilty ( $H_A$  true).

<sup>21</sup>To lower the Type 1 Error rate, we might raise our standard for conviction from “beyond a reasonable doubt” to “beyond a conceivable doubt” so fewer people would be wrongly convicted. However, this would also make it more difficult to convict the people who are actually guilty, so we would make more Type 2 Errors.

<sup>22</sup>To lower the Type 2 Error rate, we want to convict more guilty people. We could lower the standards for conviction from “beyond a reasonable doubt” to “beyond a little doubt”. Lowering the bar for guilt will also result in more wrongful convictions, raising the Type 1 Error rate.

away from the population parameter. This happens about 5% of the time (2.5% in each tail). Similarly, using a 99% confidence interval to evaluate a hypothesis is equivalent to a significance level of  $\alpha = 0.01$ .

A confidence interval is, in one sense, simplistic in the world of hypothesis tests. Consider the following two scenarios:

- The null value (the parameter value under the null hypothesis) is in the 95% confidence interval but just barely, so we would not reject  $H_0$ . However, we might like to somehow say, quantitatively, that it was a close decision.
- The null value is very far outside of the interval, so we reject  $H_0$ . However, we want to communicate that, not only did we reject the null hypothesis, but it wasn't even close. Such a case is depicted in Figure 4.13.

In Section 4.3.4, we introduce a tool called the *p-value* that will be helpful in these cases. The p-value method also extends to hypothesis tests where confidence intervals cannot be easily constructed or applied.

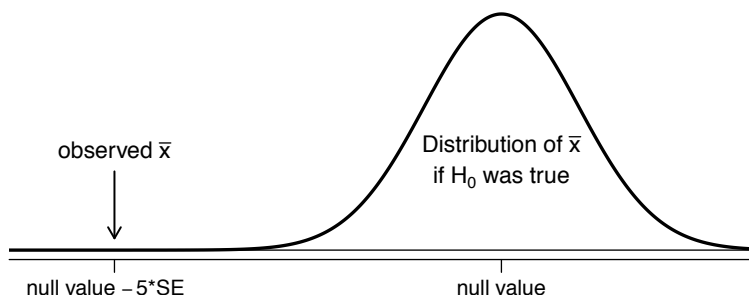


Figure 4.13: It would be helpful to quantify the strength of the evidence against the null hypothesis. In this case, the evidence is extremely strong.

#### 4.3.4 Formal testing using p-values

The p-value is a way of quantifying the strength of the evidence against the null hypothesis and in favor of the alternative. Formally the *p-value* is a conditional probability.

##### p-value

The **p-value** is the probability of observing data at least as favorable to the alternative hypothesis as our current data set, if the null hypothesis is true. We typically use a summary statistic of the data, in this chapter the sample mean, to help compute the p-value and evaluate the hypotheses.

- ⊙ **Exercise 4.28** A poll by the National Sleep Foundation found that college students average about 7 hours of sleep per night. Researchers at a rural school are interested in showing that students at their school sleep longer than seven hours on average, and they would like to demonstrate this using a sample of students. What would be an appropriate skeptical position for this research?<sup>23</sup>

<sup>23</sup>A skeptic would have no reason to believe that sleep patterns at this school are different than the sleep patterns at another school.

We can set up the null hypothesis for this test as a skeptical perspective: the students at this school average 7 hours of sleep per night. The alternative hypothesis takes a new form reflecting the interests of the research: the students average more than 7 hours of sleep. We can write these hypotheses as

$$H_0: \mu = 7.$$

$$H_A: \mu > 7.$$

Using  $\mu > 7$  as the alternative is an example of a **one-sided** hypothesis test. In this investigation, there is no apparent interest in learning whether the mean is less than 7 hours.<sup>24</sup> Earlier we encountered a **two-sided** hypothesis where we looked for any clear difference, greater than or less than the null value.

Always use a two-sided test unless it was made clear prior to data collection that the test should be one-sided. Switching a two-sided test to a one-sided test after observing the data is dangerous because it can inflate the Type 1 Error rate.

**TIP: One-sided and two-sided tests**

If the researchers are only interested in showing an increase or a decrease, but not both, use a one-sided test. If the researchers would be interested in any difference from the null value – an increase or decrease – then the test should be two-sided.

**TIP: Always write the null hypothesis as an equality**

We will find it most useful if we always list the null hypothesis as an equality (e.g.  $\mu = 7$ ) while the alternative always uses an inequality (e.g.  $\mu \neq 7$ ,  $\mu > 7$ , or  $\mu < 7$ ).

The researchers at the rural school conducted a simple random sample of  $n = 110$  students on campus. They found that these students averaged 7.42 hours of sleep and the standard deviation of the amount of sleep for the students was 1.75 hours. A histogram of the sample is shown in Figure 4.14.

Before we can use a normal model for the sample mean or compute the standard error of the sample mean, we must verify conditions. (1) Because this is a simple random sample from less than 10% of the student body, the observations are independent. (2) The sample size in the sleep study is sufficiently large since it is greater than 30. (3) The data show moderate skew in Figure 4.14 and the presence of a couple of outliers. This skew and the outliers (which are not too extreme) are acceptable for a sample size of  $n = 110$ . With these conditions verified, the normal model can be safely applied to  $\bar{x}$  and the estimated standard error will be very accurate.

⊙ **Exercise 4.29** What is the standard deviation associated with  $\bar{x}$ ? That is, estimate the standard error of  $\bar{x}$ .<sup>25</sup>

The hypothesis test will be evaluated using a significance level of  $\alpha = 0.05$ . We want to consider the data under the scenario that the null hypothesis is true. In this case, the sample mean is from a distribution that is nearly normal and has mean 7 and standard deviation of about 0.17. Such a distribution is shown in Figure 4.15.

<sup>24</sup>This is entirely based on the interests of the researchers. Had they been only interested in the opposite case – showing that their students were actually averaging fewer than seven hours of sleep but not interested in showing more than 7 hours – then our setup would have set the alternative as  $\mu < 7$ .

<sup>25</sup>The standard error can be estimated from the sample standard deviation and the sample size:  $SE_{\bar{x}} = \frac{s_x}{\sqrt{n}} = \frac{1.75}{\sqrt{110}} = 0.17$ .

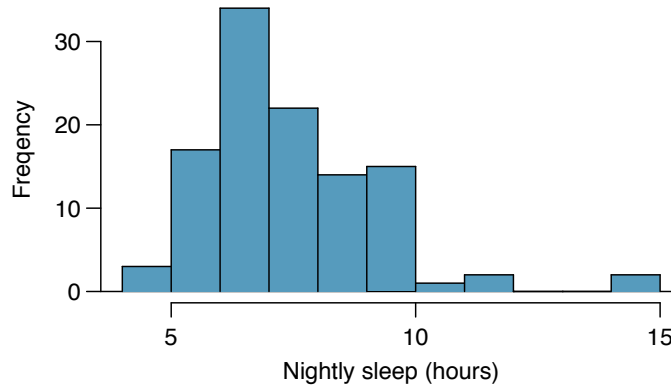


Figure 4.14: Distribution of a night of sleep for 110 college students. These data are moderately skewed.

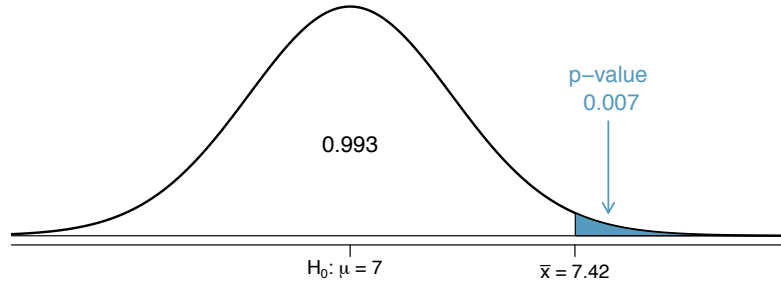


Figure 4.15: If the null hypothesis is true, then the sample mean  $\bar{x}$  came from this nearly normal distribution. The right tail describes the probability of observing such a large sample mean if the null hypothesis is true.

The shaded tail in Figure 4.15 represents the chance of observing such a large mean, conditional on the null hypothesis being true. That is, the shaded tail represents the p-value. We shade all means larger than our sample mean,  $\bar{x} = 7.42$ , because they are more favorable to the alternative hypothesis than the observed mean.

We compute the p-value by finding the tail area of this normal distribution, which we learned to do in Section 3.1. First compute the Z score of the sample mean,  $\bar{x} = 7.42$ :

$$Z = \frac{\bar{x} - \text{null value}}{SE_{\bar{x}}} = \frac{7.42 - 7}{0.17} = 2.47$$

Using the normal probability table, the lower unshaded area is found to be 0.993. Thus the shaded area is  $1 - 0.993 = 0.007$ . *If the null hypothesis is true, the probability of observing such a large sample mean for a sample of 110 students is only 0.007.* That is, if the null hypothesis is true, we would not often see such a large mean.

We evaluate the hypotheses by comparing the p-value to the significance level. Because the p-value is less than the significance level ( $\text{p-value} = 0.007 < 0.05 = \alpha$ ), we reject the null hypothesis. What we observed is so unusual with respect to the null hypothesis that it casts serious doubt on  $H_0$  and provides strong evidence favoring  $H_A$ .

**p-value as a tool in hypothesis testing**

The p-value quantifies how strongly the data favor  $H_A$  over  $H_0$ . A small p-value (usually  $< 0.05$ ) corresponds to sufficient evidence to reject  $H_0$  in favor of  $H_A$ .

**TIP: It is useful to first draw a picture to find the p-value**

It is useful to draw a picture of the distribution of  $\bar{x}$  as though  $H_0$  was true (i.e.  $\mu$  equals the null value), and shade the region (or regions) of sample means that are at least as favorable to the alternative hypothesis. These shaded regions represent the p-value.

The ideas below review the process of evaluating hypothesis tests with p-values:

- The null hypothesis represents a skeptic's position or a position of no difference. We reject this position only if the evidence strongly favors  $H_A$ .
- A small p-value means that if the null hypothesis is true, there is a low probability of seeing a point estimate at least as extreme as the one we saw. We interpret this as strong evidence in favor of the alternative.
- We reject the null hypothesis if the p-value is smaller than the significance level,  $\alpha$ , which is usually 0.05. Otherwise, we fail to reject  $H_0$ .
- We should always state the conclusion of the hypothesis test in plain language so non-statisticians can also understand the results.

The p-value is constructed in such a way that we can directly compare it to the significance level ( $\alpha$ ) to determine whether or not to reject  $H_0$ . This method ensures that the Type 1 Error rate does not exceed the significance level standard.

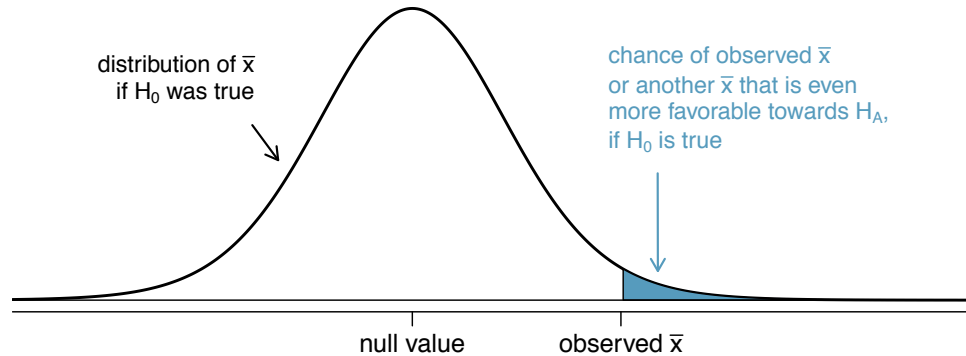


Figure 4.16: To identify the p-value, the distribution of the sample mean is considered as if the null hypothesis was true. Then the p-value is defined and computed as the probability of the observed  $\bar{x}$  or an  $\bar{x}$  even more favorable to  $H_A$  under this distribution.

⊙ **Exercise 4.30** If the null hypothesis is true, how often should the p-value be less than 0.05?<sup>26</sup>

<sup>26</sup>About 5% of the time. If the null hypothesis is true, then the data only has a 5% chance of being in the 5% of data most favorable to  $H_A$ .

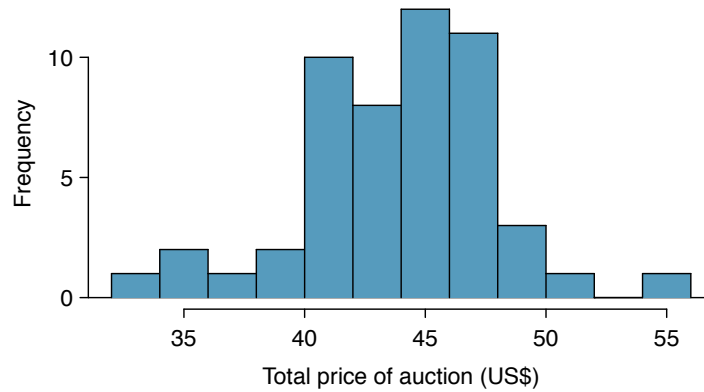


Figure 4.17: A histogram of the total auction prices for 52 Ebay auctions.

- ⊙ **Exercise 4.31** Suppose we had used a significance level of 0.01 in the sleep study. Would the evidence have been strong enough to reject the null hypothesis? (The  $p$ -value was 0.007.) What if the significance level was  $\alpha = 0.001$ ?<sup>27</sup>
- ⊙ **Exercise 4.32** Ebay might be interested in showing that buyers on its site tend to pay less than they would for the corresponding new item on Amazon. We'll research this topic for one particular product: a video game called *Mario Kart* for the Nintendo Wii. During early October 2009, Amazon sold this game for \$46.99. Set up an appropriate (one-sided!) hypothesis test to check the claim that Ebay buyers pay less during auctions at this same time.<sup>28</sup>
- ⊙ **Exercise 4.33** During early October, 2009, 52 Ebay auctions were recorded for *Mario Kart*.<sup>29</sup> The total prices for the auctions are presented using a histogram in Figure 4.17, and we may like to apply the normal model to the sample mean. Check the three conditions required for applying the normal model: (1) independence, (2) at least 30 observations, and (3) the data are not strongly skewed.<sup>30</sup>
- **Example 4.34** The average sale price of the 52 Ebay auctions for *Wii Mario Kart* was \$44.17 with a standard deviation of \$4.15. Does this provide sufficient evidence to reject the null hypothesis in Exercise 4.32? Use a significance level of  $\alpha = 0.01$ .

The hypotheses were set up and the conditions were checked in Exercises 4.32 and 4.33. The next step is to find the standard error of the sample mean and produce a sketch

<sup>27</sup>We reject the null hypothesis whenever  $p\text{-value} < \alpha$ . Thus, we would still reject the null hypothesis if  $\alpha = 0.01$  but not if the significance level had been  $\alpha = 0.001$ .

<sup>28</sup>The skeptic would say the average is the same on Ebay, and we are interested in showing the average price is lower.

$H_0$ : The average auction price on Ebay is equal to (or more than) the price on Amazon. We write only the equality in the statistical notation:  $\mu_{\text{ebay}} = 46.99$ .

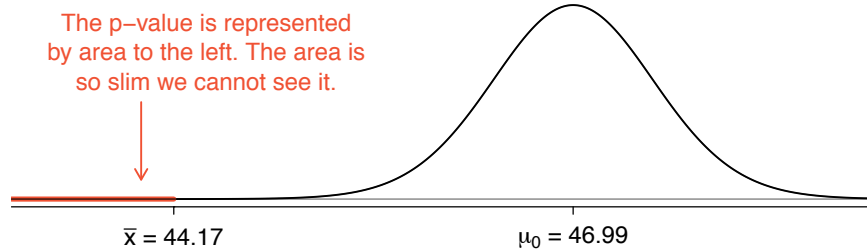
$H_A$ : The average price on Ebay is less than the price on Amazon,  $\mu_{\text{ebay}} < 46.99$ .

<sup>29</sup>These data were collected by OpenIntro staff.

<sup>30</sup>(1) The independence condition is unclear. We will make the assumption that the observations are independent, which we should report with any final results. (2) The sample size is sufficiently large:  $n = 52 \geq 30$ . (3) The data distribution is not strongly skewed; it is approximately symmetric.

to help find the p-value.

$$SE_{\bar{x}} = s/\sqrt{n} = 4.15/\sqrt{52} = 0.5755$$



Because the alternative hypothesis says we are looking for a smaller mean, we shade the lower tail. We find this shaded area by using the Z score and normal probability table:  $Z = \frac{44.17 - 46.99}{0.5755} = -4.90$ , which has area less than 0.0002. The area is so small we cannot really see it on the picture. This lower tail area corresponds to the p-value.

Because the p-value is so small – specifically, smaller than  $\alpha = 0.01$  – this provides sufficiently strong evidence to reject the null hypothesis in favor of the alternative. The data provide statistically significant evidence that the average price on Ebay is lower than Amazon’s asking price.

#### 4.3.5 Two-sided hypothesis testing with p-values

We now consider how to compute a p-value for a two-sided test. In one-sided tests, we shade the single tail in the direction of the alternative hypothesis. For example, when the alternative had the form  $\mu > 7$ , then the p-value was represented by the upper tail (Figure 4.16). When the alternative was  $\mu < 46.99$ , the p-value was the lower tail (Exercise 4.32). In a two-sided test, *we shade two tails* since evidence in either direction is favorable to  $H_A$ .

⊙ **Exercise 4.35** Earlier we talked about a research group investigating whether the students at their school slept longer than 7 hours each night. Let’s consider a second group of researchers who want to evaluate whether the students at their college differ from the norm of 7 hours. Write the null and alternative hypotheses for this investigation.<sup>31</sup>

● **Example 4.36** The second college randomly samples 72 students and finds a mean of  $\bar{x} = 6.83$  hours and a standard deviation of  $s = 1.8$  hours. Does this provide strong evidence against  $H_0$  in Exercise 4.35? Use a significance level of  $\alpha = 0.05$ .

First, we must verify assumptions. (1) A simple random sample of less than 10% of the student body means the observations are independent. (2) The sample size is 72, which is greater than 30. (3) Based on the earlier distribution and what we already know about college student sleep habits, the distribution is probably not strongly skewed.

Next we can compute the standard error ( $SE_{\bar{x}} = \frac{s}{\sqrt{n}} = 0.21$ ) of the estimate and create a picture to represent the p-value, shown in Figure 4.18. Both tails are shaded.

<sup>31</sup>Because the researchers are interested in any difference, they should use a two-sided setup:  $H_0 : \mu = 7$ ,  $H_A : \mu \neq 7$ .



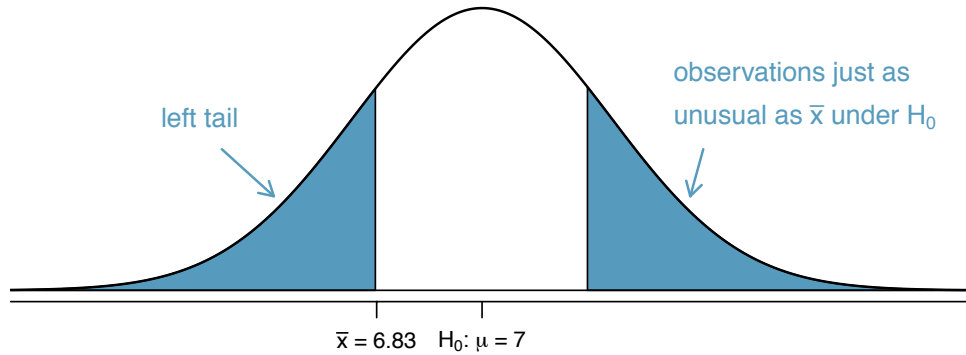


Figure 4.18:  $H_A$  is two-sided, so *both* tails must be counted for the p-value.

An estimate of 7.17 or more provides at least as strong of evidence against the null hypothesis and in favor of the alternative as the observed estimate,  $\bar{x} = 6.83$ .

We can calculate the tail areas by first finding the lower tail corresponding to  $\bar{x}$ :

$$Z = \frac{6.83 - 7.00}{0.21} = -0.81 \xrightarrow{\text{table}} \text{left tail} = 0.2090$$

Because the normal model is symmetric, the right tail will have the same area as the left tail. The p-value is found as the sum of the two shaded tails:

$$\text{p-value} = \text{left tail} + \text{right tail} = 2 \times (\text{left tail}) = 0.4180$$

This p-value is relatively large (larger than  $\alpha = 0.05$ ), so we should not reject  $H_0$ . That is, if  $H_0$  is true, it would not be very unusual to see a sample mean this far from 7 hours simply due to sampling variation. Thus, we do not have sufficient evidence to conclude that the mean is different than 7 hours.

- **Example 4.37** It is never okay to change two-sided tests to one-sided tests after observing the data. In this example we explore the consequences of ignoring this advice. Using  $\alpha = 0.05$ , we show that freely switching from two-sided tests to one-sided tests will cause us to make twice as many Type 1 Errors as intended.

Suppose the sample mean was larger than the null value,  $\mu_0$  (e.g.  $\mu_0$  would represent 7 if  $H_0: \mu = 7$ ). Then if we can flip to a one-sided test, we would use  $H_A: \mu > \mu_0$ . Now if we obtain any observation with a Z score greater than 1.65, we would reject  $H_0$ . If the null hypothesis is true, we incorrectly reject the null hypothesis about 5% of the time when the sample mean is above the null value, as shown in Figure 4.19.

Suppose the sample mean was smaller than the null value. Then if we change to a one-sided test, we would use  $H_A: \mu < \mu_0$ . If  $\bar{x}$  had a Z score smaller than -1.65, we would reject  $H_0$ . If the null hypothesis is true, then we would observe such a case about 5% of the time.

By examining these two scenarios, we can determine that we will make a Type 1 Error  $5\% + 5\% = 10\%$  of the time if we are allowed to swap to the “best” one-sided test for the data. This is twice the error rate we prescribed with our significance level:  $\alpha = 0.05$  (!).

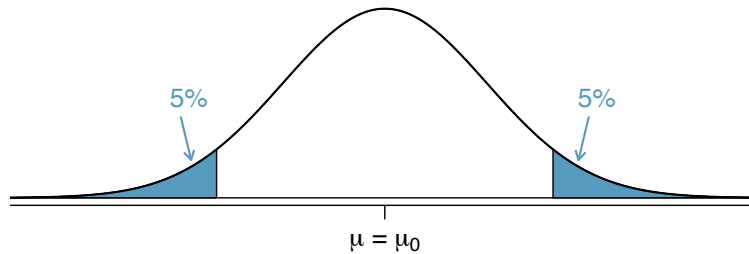


Figure 4.19: The shaded regions represent areas where we would reject  $H_0$  under the bad practices considered in Example 4.37 when  $\alpha = 0.05$ .

**Caution: One-sided hypotheses are allowed only *before* seeing data**

After observing data, it is tempting to turn a two-sided test into a one-sided test. Avoid this temptation. Hypotheses must be set up *before* observing the data. If they are not, the test must be two-sided.

### 4.3.6 Choosing a significance level

Choosing a significance level for a test is important in many contexts, and the traditional level is 0.05. However, it is often helpful to adjust the significance level based on the application. We may select a level that is smaller or larger than 0.05 depending on the consequences of any conclusions reached from the test.

If making a Type 1 Error is dangerous or especially costly, we should choose a small significance level (e.g. 0.01). Under this scenario we want to be very cautious about rejecting the null hypothesis, so we demand very strong evidence favoring  $H_A$  before we would reject  $H_0$ .

If a Type 2 Error is relatively more dangerous or much more costly than a Type 1 Error, then we should choose a higher significance level (e.g. 0.10). Here we want to be cautious about failing to reject  $H_0$  when the null is actually false. We will discuss this particular case in greater detail in Section 4.6.

**Significance levels should reflect consequences of errors**

The significance level selected for a test should reflect the consequences associated with Type 1 and Type 2 Errors.

- **Example 4.38** A car manufacturer is considering a higher quality but more expensive supplier for window parts in its vehicles. They sample a number of parts from their current supplier and also parts from the new supplier. They decide that if the high quality parts will last more than 12% longer, it makes financial sense to switch to this more expensive supplier. Is there good reason to modify the significance level in such a hypothesis test?

The null hypothesis is that the more expensive parts last no more than 12% longer while the alternative is that they do last more than 12% longer. This decision is just one of the many regular factors that have a marginal impact on the car and company. A significance level of 0.05 seems reasonable since neither a Type 1 or Type 2 error should be dangerous or (relatively) much more expensive.

- **Example 4.39** The same car manufacturer is considering a slightly more expensive supplier for parts related to safety, not windows. If the durability of these safety components is shown to be better than the current supplier, they will switch manufacturers. Is there good reason to modify the significance level in such an evaluation?

The null hypothesis would be that the suppliers' parts are equally reliable. Because safety is involved, the car company should be eager to switch to the slightly more expensive manufacturer (reject  $H_0$ ) even if the evidence of increased safety is only moderately strong. A slightly larger significance level, such as  $\alpha = 0.10$ , might be appropriate.

- **Exercise 4.40** A part inside of a machine is very expensive to replace. However, the machine usually functions properly even if this part is broken, so the part is replaced only if we are extremely certain it is broken based on a series of measurements. Identify appropriate hypotheses for this test (in plain language) and suggest an appropriate significance level.<sup>32</sup>

## 4.4 Examining the Central Limit Theorem

The normal model for the sample mean tends to be very good when the sample consists of at least 30 independent observations and the population data are not strongly skewed. The Central Limit Theorem provides the theory that allows us to make this assumption.

### Central Limit Theorem, informal definition

The distribution of  $\bar{x}$  is approximately normal. The approximation can be poor if the sample size is small, but it improves with larger sample sizes.

The Central Limit Theorem states that when the sample size is small, the normal approximation may not be very good. However, as the sample size becomes large, the normal approximation improves. We will investigate three cases to see roughly when the approximation is reasonable.

We consider three data sets: one from a *uniform* distribution, one from an *exponential* distribution, and the other from a *log-normal* distribution. These distributions are shown in the top panels of Figure 4.20. The uniform distribution is symmetric, the exponential distribution may be considered as having moderate skew since its right tail is relatively short (few outliers), and the log-normal distribution is strongly skewed and will tend to produce more apparent outliers.

The left panel in the  $n = 2$  row represents the sampling distribution of  $\bar{x}$  if it is the sample mean of two observations from the uniform distribution shown. The dashed line represents the closest approximation of the normal distribution. Similarly, the center and right panels of the  $n = 2$  row represent the respective distributions of  $\bar{x}$  for data from exponential and log-normal distributions.

<sup>32</sup>Here the null hypothesis is that the part is not broken, and the alternative is that it is broken. If we don't have sufficient evidence to reject  $H_0$ , we would not replace the part. It sounds like failing to fix the part if it is broken ( $H_0$  false,  $H_A$  true) is not very problematic, and replacing the part is expensive. Thus, we should require very strong evidence against  $H_0$  before we replace the part. Choose a small significance level, such as  $\alpha = 0.01$ .

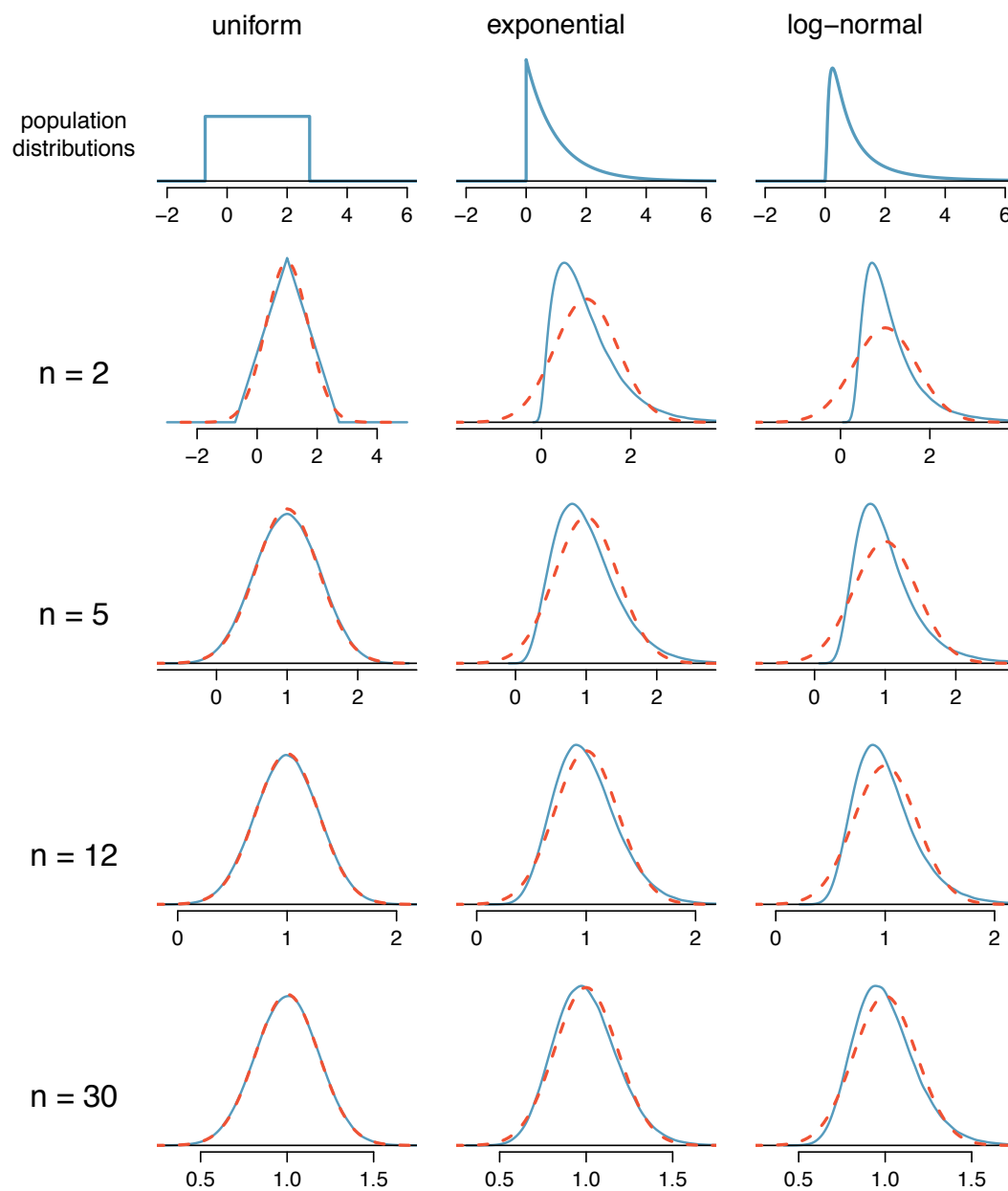


Figure 4.20: Sampling distributions for the mean at different sample sizes and for three different distributions. The dashed red lines show normal distributions.

⊙ **Exercise 4.41** Examine the distributions in each row of Figure 4.20. What do you notice about the normal approximation for each sampling distribution as the sample size becomes larger?<sup>33</sup>

● **Example 4.42** Would the normal approximation be good in all applications where the sample size is at least 30?

Not necessarily. For example, the normal approximation for the log-normal example is questionable for a sample size of 30. Generally, the more skewed a population distribution or the more common the frequency of outliers, the larger the sample required to guarantee the distribution of the sample mean is nearly normal.

**TIP: With larger  $n$ , the sampling distribution of  $\bar{x}$  becomes more normal**

As the sample size increases, the normal model for  $\bar{x}$  becomes more reasonable. We can also relax our condition on skew when the sample size is very large.

We discussed in Section 4.1.3 that the sample standard deviation,  $s$ , could be used as a substitute of the population standard deviation,  $\sigma$ , when computing the standard error. This estimate tends to be reasonable when  $n \geq 30$ . We will encounter alternative distributions for smaller sample sizes in Chapters 5 and 6.

● **Example 4.43** Figure 4.21 shows a histogram of 50 observations. These represent winnings and losses from 50 consecutive days of a professional poker player. Can the normal approximation be applied to the sample mean, 90.69?

We should consider each of the required conditions.

- (1) These are referred to as **time series data**, because the data arrived in a particular sequence. If the player wins on one day, it may influence how she plays the next. To make the assumption of independence we should perform careful checks on such data. While the supporting analysis is not shown, no evidence was found to indicate the observations are not independent.
- (2) The sample size is 50, satisfying the sample size condition.
- (3) There are two outliers, one very extreme, which suggests the data are very strongly skewed or very distant outliers may be common for this type of data. Outliers can play an important role and affect the distribution of the sample mean and the estimate of the standard error.

Since we should be skeptical of the independence of observations and the very extreme upper outlier poses a challenge, we should not use the normal model for the sample mean of these 50 observations. If we can obtain a much larger sample, perhaps several hundred observations, then the concerns about skew and outliers would no longer apply.

**Caution: Examine data structure when considering independence**

Some data sets are collected in such a way that they have a natural underlying structure between observations, e.g. when observations occur consecutively. Be especially cautious about independence assumptions regarding such data sets.

<sup>33</sup>The normal approximation becomes better as larger samples are used.

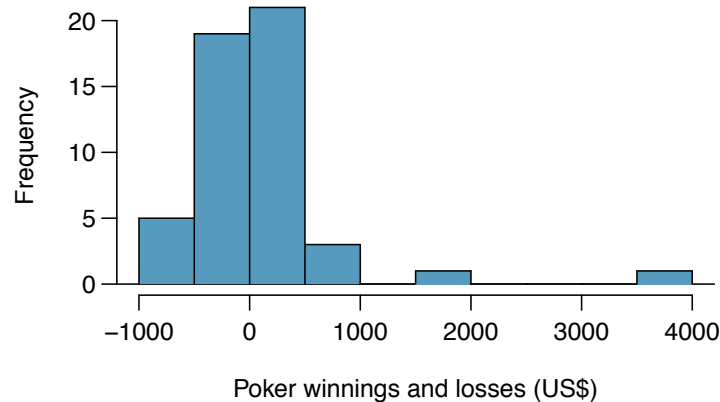


Figure 4.21: Sample distribution of poker winnings. These data include some very clear outliers. These are problematic when considering the normality of the sample mean. For example, outliers are often an indicator of very strong skew.

**Caution: Watch out for strong skew and outliers**

Strong skew is often identified by the presence of clear outliers. If a data set has prominent outliers, or such observations are somewhat common for the type of data under study, then it is useful to collect a sample with many more than 30 observations if the normal model will be used for  $\bar{x}$ . There are no simple guidelines for what sample size is big enough for all situations, so proceed with caution when working in the presence of strong skew or more extreme outliers.

## 4.5 Inference for other estimators

The sample mean is not the only point estimate for which the sampling distribution is nearly normal. For example, the sampling distribution of sample proportions closely resembles the normal distribution when the sample size is sufficiently large. In this section, we introduce a number of examples where the normal approximation is reasonable for the point estimate. Chapters 5 and 6 will revisit each of the point estimates you see in this section along with some other new statistics.

We make another important assumption about each point estimate encountered in this section: the estimate is unbiased. A point estimate is **unbiased** if the sampling distribution of the estimate is centered at the parameter it estimates. That is, an unbiased estimate does not naturally over or underestimate the parameter. Rather, it tends to provide a “good” estimate. The sample mean is an example of an unbiased point estimate, as are each of the examples we introduce in this section.

Finally, we will discuss the general case where a point estimate may follow some distribution other than the normal distribution. We also provide guidance about how to handle scenarios where the statistical techniques you are familiar with are insufficient for the problem at hand.

### 4.5.1 Confidence intervals for nearly normal point estimates

In Section 4.2, we used the point estimate  $\bar{x}$  with a standard error  $SE_{\bar{x}}$  to create a 95% confidence interval for the population mean:

$$\bar{x} \pm 1.96 \times SE_{\bar{x}} \quad (4.44)$$

We constructed this interval by noting that the sample mean is within 1.96 standard errors of the actual mean about 95% of the time. This same logic generalizes to any unbiased point estimate that is nearly normal. We may also generalize the confidence level by using a place-holder  $z^*$ .

#### General confidence interval for the normal sampling distribution case

A confidence interval based on an unbiased and nearly normal point estimate is

$$\text{point estimate} \pm z^* SE \quad (4.45)$$

where  $z^*$  is selected to correspond to the confidence level, and  $SE$  represents the standard error. The value  $z^* SE$  is called the *margin of error*.

Generally the standard error for a point estimate is estimated from the data and computed using a formula. For example, the standard error for the sample mean is

$$SE_{\bar{x}} = \frac{s}{\sqrt{n}}$$

In this section, we provide the computed standard error for each example and exercise without detailing where the values came from. In future chapters, you will learn to fill in these and other details for each situation.

- **Example 4.46** In Exercise 4.1 on page 161, we computed a point estimate for the average difference in run times between men and women:  $\bar{x}_{\text{women}} - \bar{x}_{\text{men}} = 14.48$  minutes. This point estimate is associated with a nearly normal distribution with standard error  $SE = 2.78$  minutes. What is a reasonable 95% confidence interval for the difference in average run times?

The normal approximation is said to be valid, so we apply Equation (4.45):

$$\text{point estimate} \pm z^* SE \rightarrow 14.48 \pm 1.96 \times 2.78 \rightarrow (9.03, 19.93)$$

Thus, we are 95% confident that the men were, on average, between 9.03 and 19.93 minutes faster than women in the 2012 Cherry Blossom Run. That is, the actual average difference is plausibly between 9.03 and 19.93 minutes with 95% confidence.

- **Example 4.47** Does Example 4.46 guarantee that if a husband and wife both ran in the race, the husband would run between 9.03 and 19.93 minutes faster than the wife?

Our confidence interval says absolutely nothing about individual observations. It only makes a statement about a plausible range of values for the *average* difference between all men and women who participated in the run.

- ⊙ **Exercise 4.48** What  $z^*$  would be appropriate for a 99% confidence level? For help, see Figure 4.10 on page 169.<sup>34</sup>
- ⊙ **Exercise 4.49** The proportion of men in the `run10Samp` sample is  $\hat{p} = 0.45$ . This sample meets certain conditions that ensure  $\hat{p}$  will be nearly normal, and the standard error of the estimate is  $SE_{\hat{p}} = 0.05$ . Create a 90% confidence interval for the proportion of participants in the 2012 Cherry Blossom Run who are men.<sup>35</sup>

### 4.5.2 Hypothesis testing for nearly normal point estimates

Just as the confidence interval method works with many other point estimates, we can generalize our hypothesis testing methods to new point estimates. Here we only consider the p-value approach, introduced in Section 4.3.4, since it is the most commonly used technique and also extends to non-normal cases.

#### Hypothesis testing using the normal model

1. First write the hypotheses in plain language, then set them up in mathematical notation.
2. Identify an appropriate point estimate of the parameter of interest.
3. Verify conditions to ensure the standard error estimate is reasonable and the point estimate is nearly normal and unbiased.
4. Compute the standard error. Draw a picture depicting the distribution of the estimate under the idea that  $H_0$  is true. Shade areas representing the p-value.
5. Using the picture and normal model, compute the *test statistic* (Z score) and identify the p-value to evaluate the hypotheses. Write a conclusion in plain language.

- ⊙ **Exercise 4.50** A drug called sulphinpyrazone was under consideration for use in reducing the death rate in heart attack patients. To determine whether the drug was effective, a set of 1,475 patients were recruited into an experiment and randomly split into two groups: a control group that received a placebo and a treatment group that received the new drug. What would be an appropriate null hypothesis? And the alternative?<sup>36</sup>

We can formalize the hypotheses from Exercise 4.50 by letting  $p_{\text{control}}$  and  $p_{\text{treatment}}$  represent the proportion of patients who died in the control and treatment groups, respec-

<sup>34</sup>We seek  $z^*$  such that 99% of the area under the normal curve will be between the Z scores  $-z^*$  and  $z^*$ . Because the remaining 1% is found in the tails, each tail has area 0.5%, and we can identify  $-z^*$  by looking up 0.0050 in the normal probability table:  $z^* = 2.58$ . See also Figure 4.10 on page 169.

<sup>35</sup>We use  $z^* = 1.65$  (see Exercise 4.17 on page 170), and apply the general confidence interval formula:

$$\hat{p} \pm z^* SE_{\hat{p}} \rightarrow 0.45 \pm 1.65 \times 0.05 \rightarrow (0.3675, 0.5325)$$

Thus, we are 90% confident that between 37% and 53% of the participants were men.

<sup>36</sup>The skeptic's perspective is that the drug does not work at reducing deaths in heart attack patients ( $H_0$ ), while the alternative is that the drug does work ( $H_A$ ).



tively. Then the hypotheses can be written as

$$\begin{aligned} H_0 : p_{\text{control}} &= p_{\text{treatment}} && \text{(the drug doesn't work)} \\ H_A : p_{\text{control}} &> p_{\text{treatment}} && \text{(the drug works)} \end{aligned}$$

or equivalently,

$$\begin{aligned} H_0 : p_{\text{control}} - p_{\text{treatment}} &= 0 && \text{(the drug doesn't work)} \\ H_A : p_{\text{control}} - p_{\text{treatment}} &> 0 && \text{(the drug works)} \end{aligned}$$

Strong evidence against the null hypothesis and in favor of the alternative would correspond to an observed difference in death rates,

$$\text{point estimate} = \hat{p}_{\text{control}} - \hat{p}_{\text{treatment}}$$

being larger than we would expect from chance alone. This difference in sample proportions represents a point estimate that is useful in evaluating the hypotheses.

- **Example 4.51** We want to evaluate the hypothesis setup from Exercise 4.50 using data from the actual study.<sup>37</sup> In the control group, 60 of 742 patients died. In the treatment group, 41 of 733 patients died. The sample difference in death rates can be summarized as

$$\text{point estimate} = \hat{p}_{\text{control}} - \hat{p}_{\text{treatment}} = \frac{60}{742} - \frac{41}{733} = 0.025$$

This point estimate is nearly normal and is an unbiased estimate of the actual difference in death rates. The standard error of this sample difference is  $SE = 0.013$ . Evaluate the hypothesis test at a 5% significance level:  $\alpha = 0.05$ .

We would like to identify the p-value to evaluate the hypotheses. If the null hypothesis is true, then the point estimate would have come from a nearly normal distribution, like the one shown in Figure 4.22. The distribution is centered at zero since  $p_{\text{control}} - p_{\text{treatment}} = 0$  under the null hypothesis. Because a large positive difference provides evidence against the null hypothesis and in favor of the alternative, the upper tail has been shaded to represent the p-value. We need not shade the lower tail since this is a one-sided test: an observation in the lower tail does not support the alternative hypothesis.

The p-value can be computed by using the Z score of the point estimate and the normal probability table.

$$Z = \frac{\text{point estimate} - \text{null value}}{SE_{\text{point estimate}}} = \frac{0.025 - 0}{0.013} = 1.92 \quad (4.52)$$

Examining  $Z$  in the normal probability table, we find that the lower unshaded tail is about 0.973. Thus, the upper shaded tail representing the p-value is

$$\text{p-value} = 1 - 0.973 = 0.027$$

Because the p-value is less than the significance level ( $\alpha = 0.05$ ), we say the null hypothesis is implausible. That is, we reject the null hypothesis in favor of the alternative and conclude that the drug is effective at reducing deaths in heart attack patients.

<sup>37</sup>Anturane Reinfarction Trial Research Group. 1980. Sulfapyrazone in the prevention of sudden death after myocardial infarction. *New England Journal of Medicine* 302(5):250-256.

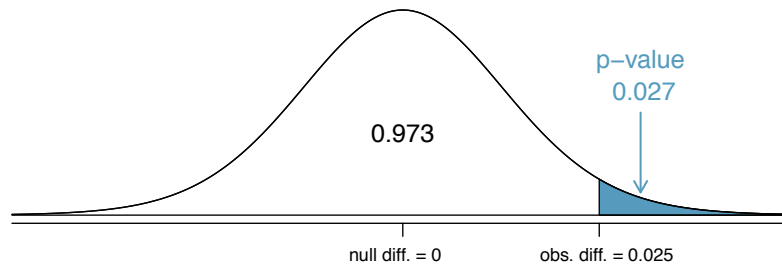


Figure 4.22: The distribution of the sample difference if the null hypothesis is true.

The Z score in Equation (4.52) is called a **test statistic**. In most hypothesis tests, a test statistic is a particular data summary that is especially useful for computing the p-value and evaluating the hypothesis test. In the case of point estimates that are nearly normal, the test statistic is the Z score.

#### Test statistic

A *test statistic* is a special summary statistic that is particularly useful for evaluating a hypothesis test or identifying the p-value. When a point estimate is nearly normal, we use the Z score of the point estimate as the test statistic. In later chapters we encounter situations where other test statistics are helpful.

### 4.5.3 Non-normal point estimates

We may apply the ideas of confidence intervals and hypothesis testing to cases where the point estimate or test statistic is not necessarily normal. There are many reasons why such a situation may arise:

- the sample size is too small for the normal approximation to be valid;
- the standard error estimate may be poor; or
- the point estimate tends towards some distribution that is not the normal distribution.

For each case where the normal approximation is not valid, our first task is always to understand and characterize the sampling distribution of the point estimate or test statistic. Next, we can apply the general frameworks for confidence intervals and hypothesis testing to these alternative distributions.

### 4.5.4 When to retreat

Statistical tools rely on conditions. When the conditions are not met, these tools are unreliable and drawing conclusions from them is treacherous. The conditions for these tools typically come in two forms.

- **The individual observations must be independent.** A random sample from less than 10% of the population ensures the observations are independent. In experiments, we generally require that subjects are randomized into groups. If independence fails,

then advanced techniques must be used, and in some such cases, inference may not be possible.

- **Other conditions focus on sample size and skew.** For example, if the sample size is too small, the skew too strong, or extreme outliers are present, then the normal model for the sample mean will fail.

Verification of conditions for statistical tools is always necessary. Whenever conditions are not satisfied for a statistical technique, there are three options. The first is to learn new methods that are appropriate for the data. The second route is to consult a statistician.<sup>38</sup> The third route is to ignore the failure of conditions. This last option effectively invalidates any analysis and may discredit novel and interesting findings.

Finally, we caution that there may be no inference tools helpful when considering data that include unknown biases, such as convenience samples. For this reason, there are books, courses, and researchers devoted to the techniques of sampling and experimental design. See Sections 1.3-1.5 for basic principles of data collection.

## 4.6 Sample size and power (special topic)

The Type 2 Error rate and the magnitude of the error for a point estimate are controlled by the sample size. Real differences from the null value, even large ones, may be difficult to detect with small samples. If we take a very large sample, we might find a statistically significant difference but the magnitude might be so small that it is of no practical value. In this section we describe techniques for selecting an appropriate sample size based on these considerations.

### 4.6.1 Finding a sample size for a certain margin of error

Many companies are concerned about rising healthcare costs. A company may estimate certain health characteristics of its employees, such as blood pressure, to project its future cost obligations. However, it might be too expensive to measure the blood pressure of every employee at a large company, and the company may choose to take a sample instead.

- **Example 4.53** Blood pressure oscillates with the beating of the heart, and the systolic pressure is defined as the peak pressure when a person is at rest. The average systolic blood pressure for people in the U.S. is about 130 mmHg with a standard deviation of about 25 mmHg. How large of a sample is necessary to estimate the average systolic blood pressure with a margin of error of 4 mmHg using a 95% confidence level?

First, we frame the problem carefully. Recall that the margin of error is the part we add and subtract from the point estimate when computing a confidence interval. The margin of error for a 95% confidence interval estimating a mean can be written as

$$ME_{95\%} = 1.96 \times SE = 1.96 \times \frac{\sigma_{employee}}{\sqrt{n}}$$

---

<sup>38</sup>If you work at a university, then there may be campus consulting services to assist you. Alternatively, there are many private consulting firms that are also available for hire.

The challenge in this case is to find the sample size  $n$  so that this margin of error is less than or equal to 4, which we write as an inequality:

$$1.96 \times \frac{\sigma_{employee}}{\sqrt{n}} \leq 4$$

In the above equation we wish to solve for the appropriate value of  $n$ , but we need a value for  $\sigma_{employee}$  before we can proceed. However, we haven't yet collected any data, so we have no direct estimate! Instead, we use the best estimate available to us: the approximate standard deviation for the U.S. population, 25. To proceed and solve for  $n$ , we substitute 25 for  $\sigma_{employee}$ :

$$\begin{aligned} 1.96 \times \frac{\sigma_{employee}}{\sqrt{n}} &\approx 1.96 \times \frac{25}{\sqrt{n}} \leq 4 \\ 1.96 \times \frac{25}{4} &\leq \sqrt{n} \\ \left(1.96 \times \frac{25}{4}\right)^2 &\leq n \\ 150.06 &\leq n \end{aligned}$$

This suggests we should choose a sample size of at least 151 employees. We round up because the sample size must be *greater than or equal to 150.06*.

A potentially controversial part of Example 4.53 is the use of the U.S. standard deviation for the employee standard deviation. Usually the standard deviation is not known. In such cases, it is reasonable to review scientific literature or market research to make an educated guess about the standard deviation.

#### Identify a sample size for a particular margin of error

To estimate the necessary sample size for a maximum margin of error  $m$ , we set up an equation to represent this relationship:

$$m \geq ME = z^* \frac{\sigma}{\sqrt{n}}$$

where  $z^*$  is chosen to correspond to the desired confidence level, and  $\sigma$  is the standard deviation associated with the population. Solve for the sample size,  $n$ .

Sample size computations are helpful in planning data collection, and they require careful forethought. Next we consider another topic important in planning data collection and setting a sample size: the Type 2 Error rate.

### 4.6.2 Power and the Type 2 Error rate

Consider the following two hypotheses:

$H_0$ : The average blood pressure of employees is the same as the national average,  $\mu = 130$ .

$H_A$ : The average blood pressure of employees is different than the national average,  $\mu \neq 130$ .

Suppose the alternative hypothesis is actually true. Then we might like to know, what is the chance we make a Type 2 Error? That is, what is the chance we will fail to reject the null hypothesis even though we should reject it? The answer is not obvious! If the average blood pressure of the employees is 132 (just 2 mmHg from the null value), it might be very difficult to detect the difference unless we use a large sample size. On the other hand, it would be easier to detect a difference if the real average of employees was 140.

- **Example 4.54** Suppose the actual employee average is 132 and we take a sample of 100 individuals. Then the true sampling distribution of  $\bar{x}$  is approximately  $N(132, 2.5)$  (since  $SE = \frac{25}{\sqrt{100}} = 2.5$ ). What is the probability of successfully rejecting the null hypothesis?

This problem can be divided into two normal probability questions. First, we identify what values of  $\bar{x}$  would represent sufficiently strong evidence to reject  $H_0$ . Second, we use the hypothetical sampling distribution for that has center  $\mu = 132$  to find the probability of observing sample means in the areas we found in the first step.

**Step 1.** The null distribution could be represented by  $N(130, 2.5)$ , the same standard deviation as the true distribution but with the null value as its center. Then we can find the two tail areas by identifying the  $Z$  score corresponding to the 2.5% tails ( $\pm 1.96$ ), and solving for  $x$  in the  $Z$  score equation:

$$\begin{aligned} -1.96 = Z_1 = \frac{x_1 - 130}{2.5} & & +1.96 = Z_2 = \frac{x_2 - 130}{2.5} \\ x_1 = 125.1 & & x_2 = 134.9 \end{aligned}$$

(An equally valid approach is to recognize that  $x_1$  is  $1.96 \times SE$  below the mean and  $x_2$  is  $1.96 \times SE$  above the mean to compute the values.) Figure 4.23 shows the null distribution on the left with these two dotted cutoffs.

**Step 2.** Next, we compute the probability of rejecting  $H_0$  if  $\bar{x}$  actually came from  $N(132, 2.5)$ . This is the same as finding the two shaded tails for the second distribution in Figure 4.23. We use the  $Z$  score method:

$$\begin{aligned} Z_{left} = \frac{125.1 - 132}{2.5} = -2.76 & & Z_{right} = \frac{134.9 - 132}{2.5} = 1.16 \\ area_{left} = 0.003 & & area_{right} = 0.123 \end{aligned}$$

The probability of rejecting the null mean, if the true mean is 132, is the sum of these areas:  $0.003 + 0.123 = 0.126$ .

The probability of rejecting the null hypothesis is called the **power**. The power varies depending on what we suppose the truth might be. In Example 4.54, the difference between the null value and the supposed true mean was relatively small, so the power was also small: only 0.126. However, when the truth is far from the null value, where we use the standard error as a measure of what is far, the power tends to increase.

- ⊙ **Exercise 4.55** Suppose the true sampling distribution of  $\bar{x}$  is centered at 140. That is,  $\bar{x}$  comes from  $N(140, 2.5)$ . What would the power be under this scenario? It may be helpful to draw  $N(140, 2.5)$  and shade the area representing power on Figure 4.23; use the same cutoff values identified in Example 4.54.<sup>39</sup>

<sup>39</sup>Draw the distribution  $N(140, 2.5)$ , then find the area below 125.1 (about zero area) and above 134.9 (about 0.979). If the true mean is 140, the power is about 0.979.

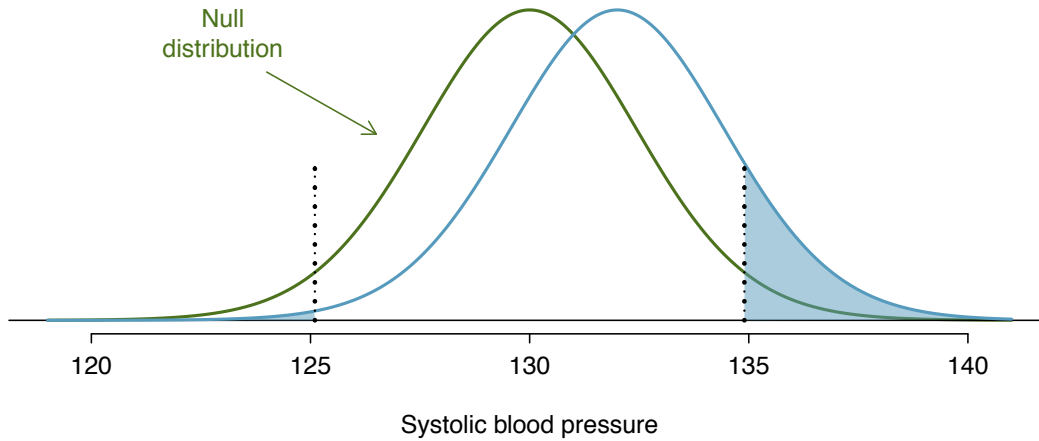


Figure 4.23: The sampling distribution of  $\bar{x}$  under two scenarios. Left:  $N(130, 2.5)$ . Right:  $N(132, 2.5)$ , and the shaded areas in this distribution represent the power of the test.

- ⊙ **Exercise 4.56** If the power of a test is 0.979 for a particular mean, what is the Type 2 Error rate for this mean?<sup>40</sup>
- ⊙ **Exercise 4.57** Provide an intuitive explanation for why we are more likely to reject  $H_0$  when the true mean is further from the null value.<sup>41</sup>

### 4.6.3 Statistical significance versus practical significance

When the sample size becomes larger, point estimates become more precise and any real differences in the mean and null value become easier to detect and recognize. Even a very small difference would likely be detected if we took a large enough sample. Sometimes researchers will take such large samples that even the slightest difference is detected. While we still say that difference is **statistically significant**, it might not be **practically significant**.

Statistically significant differences are sometimes so minor that they are not practically relevant. This is especially important to research: if we conduct a study, we want to focus on finding a meaningful result. We don't want to spend lots of money finding results that hold no practical value.

The role of a statistician in conducting a study often includes planning the size of the study. The statistician might first consult experts or scientific literature to learn what would be the smallest meaningful difference from the null value. She also would obtain some reasonable estimate for the standard deviation. With these important pieces of information, she would choose a sufficiently large sample size so that the power for the meaningful difference is perhaps 80% or 90%. While larger sample sizes may still be used, she might advise against using them in some cases, especially in sensitive areas of research.

<sup>40</sup>The Type 2 Error rate represents the probability of failing to reject the null hypothesis. Since the power is the probability we do reject, the Type 2 Error rate will be  $1 - 0.979 = 0.021$ .

<sup>41</sup>Answers may vary a little. When the truth is far from the null value, the point estimate also tends to be far from the null value, making it easier to detect the difference and reject  $H_0$ .

## 4.7 Exercises

### 4.7.1 Variability in estimates

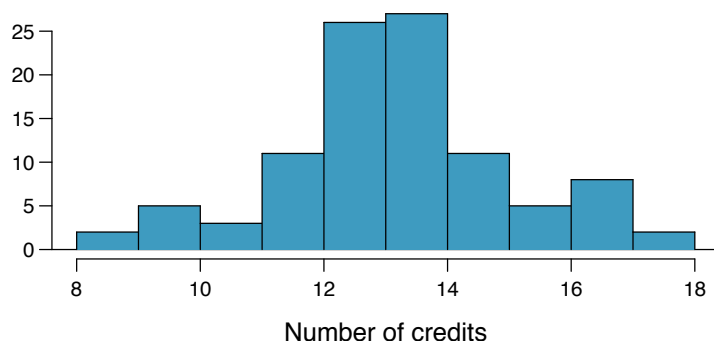
**4.1 Identify the parameter, Part I.** For each of the following situations, state whether the parameter of interest is a mean or a proportion. It may be helpful to examine whether individual responses are numerical or categorical.

- In a survey, one hundred college students are asked how many hours per week they spend on the Internet.
- In a survey, one hundred college students are asked: “What percentage of the time you spend on the Internet is part of your course work?”
- In a survey, one hundred college students are asked whether or not they cited information from Wikipedia in their papers.
- In a survey, one hundred college students are asked what percentage of their total weekly spending is on alcoholic beverages.
- In a sample of one hundred recent college graduates, it is found that 85 percent expect to get a job within one year of their graduation date.

**4.2 Identify the parameter, Part II.** For each of the following situations, state whether the parameter of interest is a mean or a proportion.

- A poll shows that 64% of Americans personally worry a great deal about federal spending and the budget deficit.
- A survey reports that local TV news has shown a 17% increase in revenue between 2009 and 2011 while newspaper revenues decreased by 6.4% during this time period.
- In a survey, high school and college students are asked whether or not they use geolocation services on their smart phones.
- In a survey, internet users are asked whether or not they purchased any Groupon coupons.
- In a survey, internet users are asked how many Groupon coupons they purchased over the last year.

**4.3 College credits.** A college counselor is interested in estimating how many credits a student typically enrolls in each semester. The counselor decides to randomly sample 100 students by using the registrar’s database of students. The histogram below shows the distribution of the number of credits taken by these students. Sample statistics for this distribution are also provided.

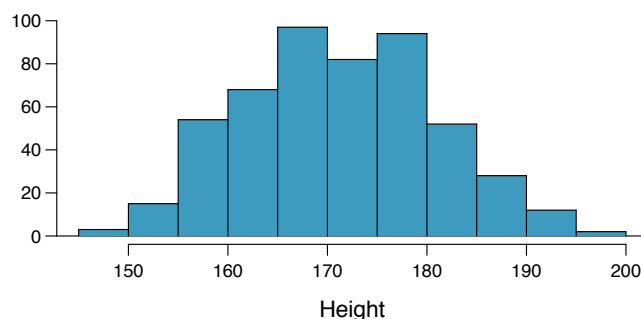


|        |       |
|--------|-------|
| Min    | 8     |
| Q1     | 13    |
| Median | 14    |
| Mean   | 13.65 |
| SD     | 1.91  |
| Q3     | 15    |
| Max    | 18    |

- What is the point estimate for the average number of credits taken per semester by students at this college? What about the median?

- (b) What is the point estimate for the standard deviation of the number of credits taken per semester by students at this college? What about the IQR?
- (c) Is a load of 16 credits unusually high for this college? What about 18 credits? Explain your reasoning. *Hint:* Observations farther than two standard deviations from the mean are usually considered to be unusual.
- (d) The college counselor takes another random sample of 100 students and this time finds a sample mean of 14.02 units. Should she be surprised that this sample statistic is slightly different than the one from the original sample? Explain your reasoning.
- (e) The sample means given above are point estimates for the mean number of credits taken by all students at that college. What measures do we use to quantify the variability of this estimate? Compute this quantity using the data from the original sample.

**4.4 Heights of adults.** Researchers studying anthropometry collected body girth measurements and skeletal diameter measurements, as well as age, weight, height and gender, for 507 physically active individuals. The histogram below shows the sample distribution of heights in centimeters.<sup>42</sup>



|        |       |
|--------|-------|
| Min    | 147.2 |
| Q1     | 163.8 |
| Median | 170.3 |
| Mean   | 171.1 |
| SD     | 9.4   |
| Q3     | 177.8 |
| Max    | 198.1 |

- (a) What is the point estimate for the average height of active individuals? What about the median?
- (b) What is the point estimate for the standard deviation of the heights of active individuals? What about the IQR?
- (c) Is a person who is 1m 80cm (180 cm) tall considered unusually tall? And is a person who is 1m 55cm (155cm) considered unusually short? Explain your reasoning.
- (d) The researchers take another random sample of physically active individuals. Would you expect the mean and the standard deviation of this new sample to be the ones given above. Explain your reasoning.
- (e) The samples means obtained are point estimates for the mean height of all active individuals, if the sample of individuals is equivalent to a simple random sample. What measure do we use to quantify the variability of such an estimate? Compute this quantity using the data from the original sample under the condition that the data are a simple random sample.

**4.5 Wireless routers.** John is shopping for wireless routers and is overwhelmed by the number of available options. In order to get a feel for the average price, he takes a random sample of 75 routers and finds that the average price for this sample is \$75 and the standard deviation is \$25.

- (a) Based on this information, how much variability should he expect to see in the mean prices of repeated samples, each containing 75 randomly selected wireless routers?
- (b) A consumer website claims that the average price of routers is \$80. Is a true average of \$80 consistent with John's sample?

<sup>42</sup>G. Heinz et al. "Exploring relationships in body dimensions". In: *Journal of Statistics Education* 11.2 (2003).



**4.6 Chocolate chip cookies.** Students are asked to count the number of chocolate chips in 22 cookies for a class activity. They found that the cookies on average had 14.77 chocolate chips with a standard deviation of 4.37 chocolate chips.

- (a) Based on this information, about how much variability should they expect to see in the mean number of chocolate chips in random samples of 22 chocolate chip cookies?
- (b) The packaging for these cookies claims that there are at least 20 chocolate chips per cookie. One student thinks this number is unreasonably high since the average they found is much lower. Another student claims the difference might be due to chance. What do you think?

### 4.7.2 Confidence intervals

**4.7 Relaxing after work.** The General Social Survey (GSS) is a sociological survey used to collect data on demographic characteristics and attitudes of residents of the United States. In 2010, the survey collected responses from 1,154 US residents. The survey is conducted face-to-face with an in-person interview of a randomly-selected sample of adults. One of the questions on the survey is “After an average work day, about how many hours do you have to relax or pursue activities that you enjoy?” A 95% confidence interval from the 2010 GSS survey is 3.53 to 3.83 hours.<sup>43</sup>

- (a) Interpret this interval in the context of the data.
- (b) What does a 95% confidence level mean in this context?
- (c) Suppose the researchers think a 90% confidence level would be more appropriate for this interval. Will this new interval be smaller or larger than the 95% confidence interval? Assume the standard deviation has remained constant since 2010.

**4.8 Mental health.** Another question on the General Social Survey introduced in Exercise 4.7 is “For how many days during the past 30 days was your mental health, which includes stress, depression, and problems with emotions, not good?” Based on responses from 1,151 US residents, the survey reported a 95% confidence interval of 3.40 to 4.24 days in 2010.

- (a) Interpret this interval in context of the data.
- (b) What does a 95% confidence level mean in this context?
- (c) Suppose the researchers think a 99% confidence level would be more appropriate for this interval. Will this new interval be smaller or larger than the 95% confidence interval?
- (d) If a new survey asking the same questions was to be done with 500 Americans, would the standard error of the estimate be larger, smaller, or about the same. Assume the standard deviation has remained constant since 2010.

**4.9 Width of a confidence interval.** Earlier in Chapter 4, we calculated the 99% confidence interval for the average age of runners in the 2012 Cherry Blossom Run as (32.7, 37.4) based on a sample of 100 runners. How could we decrease the width of this interval without losing confidence?

**4.10 Confidence levels.** If a higher confidence level means that we are more confident about the number we are reporting, why don’t we always report a confidence interval with the highest possible confidence level?

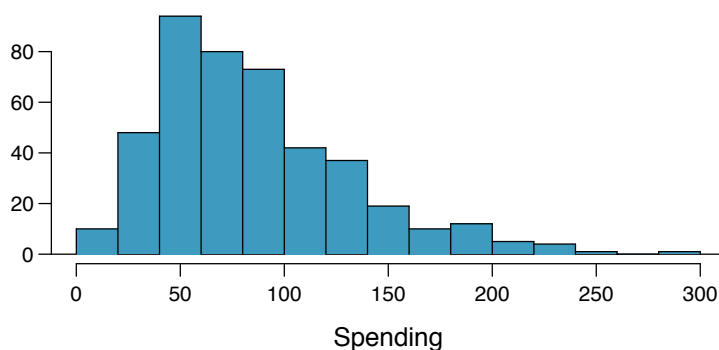
---

<sup>43</sup>National Opinion Research Center, General Social Survey, 2010.

**4.11 Waiting at an ER, Part I.** A hospital administrator hoping to improve wait times decides to estimate the average emergency room waiting time at her hospital. He collects a simple random sample of 64 patients and determines the time (in minutes) between when they checked in to the ER until they were first seen by a doctor. A 95% confidence interval based on this sample is (128 minutes, 147 minutes), which is based on the normal model for the mean. Determine whether the following statements are true or false, and explain your reasoning for those statements you identify as false.

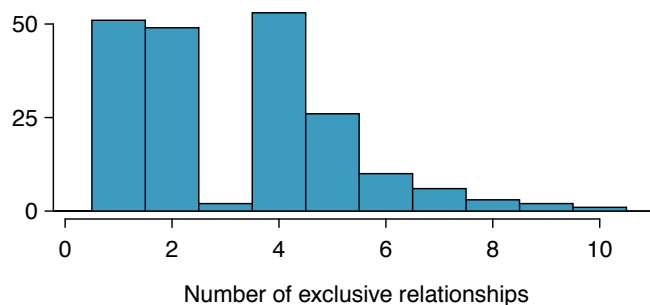
- This confidence interval is not valid since we do not know if the population distribution of the ER wait times is nearly normal.
- We are 95% confident that the average waiting time of these 64 emergency room patients is between 128 and 147 minutes.
- We are 95% confident that the average waiting time of all patients at this hospital's emergency room is between 128 and 147 minutes.
- 95% of such random samples would have a sample mean between 128 and 147 minutes.
- A 99% confidence interval would be narrower than the 95% confidence interval since we need to be more sure of our estimate.
- The margin of error is 9.5 and the sample mean is 137.5.
- In order to decrease the margin of error of a 95% confidence interval to half of what it is now, we would need to double the sample size.

**4.12 Thanksgiving spending, Part I.** The 2009 holiday retail season, which kicked off on November 27, 2009 (the day after Thanksgiving), had been marked by somewhat lower self-reported consumer spending than was seen during the comparable period in 2008. To get an estimate of consumer spending, 436 randomly sampled American adults were surveyed. Daily consumer spending for the six-day period after Thanksgiving, spanning the Black Friday weekend and Cyber Monday, averaged \$84.71. A 95% confidence interval based on this sample is (\$80.31, \$89.11). Determine whether the following statements are true or false, and explain your reasoning.



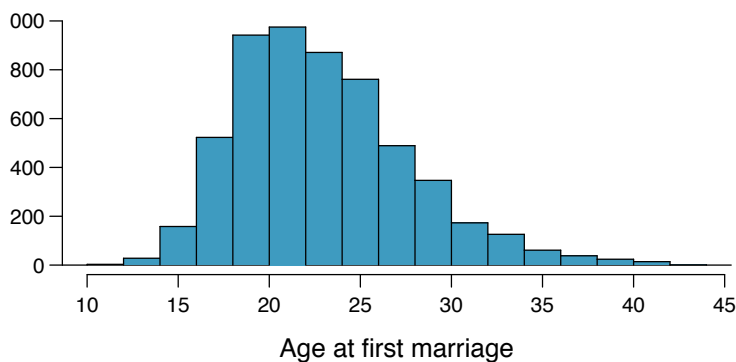
- We are 95% confident that the average spending of these 436 American adults is between \$80.31 and \$89.11.
- This confidence interval is not valid since the distribution of spending in the sample is right skewed.
- 95% of such random samples would have a sample mean between \$80.31 and \$89.11.
- We are 95% confident that the average spending of all American adults is between \$80.31 and \$89.11.
- A 90% confidence interval would be narrower than the 95% confidence interval since we don't need to be as sure about capturing the parameter.
- In order to decrease the margin of error of a 95% confidence interval to a third of what it is now, we would need to use a sample 3 times larger.
- The margin of error for the reported interval is 4.4.

**4.13 Exclusive relationships.** A survey was conducted on 203 undergraduates from Duke University who took an introductory statistics course in Spring 2012. Among many other questions, this survey asked them about the number of exclusive relationships they have been in. The histogram below shows the distribution of the data from this sample. The sample average is 3.2 with a standard deviation of 1.97.



Estimate the average number of exclusive relationships Duke students have been in using a 90% confidence interval and interpret this interval in context. Check any conditions required for inference, and note any assumptions you must make as you proceed with your calculations and conclusions.

**4.14 Age at first marriage, Part I.** The National Survey of Family Growth conducted by the Centers for Disease Control gathers information on family life, marriage and divorce, pregnancy, infertility, use of contraception, and men's and women's health. One of the variables collected on this survey is the age at first marriage. The histogram below shows the distribution of ages at first marriage of 5,534 randomly sampled women between 2006 and 2010. The average age at first marriage among these women is 23.44 with a standard deviation of 4.72.<sup>44</sup>



Estimate the average age at first marriage of women using a 95% confidence interval, and interpret this interval in context. Discuss any relevant assumptions.

<sup>44</sup>National Survey of Family Growth, 2006-2010 Cycle.

### 4.7.3 Hypothesis testing

**4.15 Identify hypotheses, Part I.** Write the null and alternative hypotheses in words and then symbols for each of the following situations.

- New York is known as “the city that never sleeps”. A random sample of 25 New Yorkers were asked how much sleep they get per night. Do these data provide convincing evidence that New Yorkers on average sleep less than 8 hours a night?
- Employers at a firm are worried about the effect of March Madness, a basketball championship held each spring in the US, on employee productivity. They estimate that on a regular business day employees spend on average 15 minutes of company time checking personal email, making personal phone calls, etc. They also collect data on how much company time employees spend on such non-business activities during March Madness. They want to determine if these data provide convincing evidence that employee productivity decreases during March Madness.

**4.16 Identify hypotheses, Part II.** Write the null and alternative hypotheses in words and using symbols for each of the following situations.

- Since 2008, chain restaurants in California have been required to display calorie counts of each menu item. Prior to menus displaying calorie counts, the average calorie intake of diners at a restaurant was 1100 calories. After calorie counts started to be displayed on menus, a nutritionist collected data on the number of calories consumed at this restaurant from a random sample of diners. Do these data provide convincing evidence of a difference in the average calorie intake of a diners at this restaurant?
- Based on the performance of those who took the GRE exam between July 1, 2004 and June 30, 2007, the average Verbal Reasoning score was calculated to be 462. In 2011 the average verbal score was slightly higher. Do these data provide convincing evidence that the average GRE Verbal Reasoning score has changed since 2004?<sup>45</sup>

**4.17 Online communication.** A study suggests that the average college student spends 2 hours per week communicating with others online. You believe that this is an underestimate and decide to collect your own sample for a hypothesis test. You randomly sample 60 students from your dorm and find that on average they spent 3.5 hours a week communicating with others online. A friend of yours, who offers to help you with the hypothesis test, comes up with the following set of hypotheses. Indicate any errors you see.

$$H_0 : \bar{x} < 2 \text{ hours}$$

$$H_A : \bar{x} > 3.5 \text{ hours}$$

**4.18 Age at first marriage, Part II.** Exercise 4.14 presents the results of a 2006 - 2010 survey showing that the average age of women at first marriage is 23.44. Suppose a researcher believes that this value has increased in 2012, but he would also be interested if he found a decrease. Below is how he set up his hypotheses. Indicate any errors you see.

$$H_0 : \bar{x} = 23.44 \text{ years old}$$

$$H_A : \bar{x} > 23.44 \text{ years old}$$

**4.19 Waiting at an ER, Part II.** Exercise 4.11 provides a 95% confidence interval for the mean waiting time at an emergency room (ER) of (128 minutes, 147 minutes).

- A local newspaper claims that the average waiting time at this ER exceeds 3 hours. What do you think of this claim?
- The Dean of Medicine at this hospital claims the average wait time is 2.2 hours. What do you think of this claim?
- Without actually calculating the interval, determine if the claim of the Dean from part (b) would be considered reasonable based on a 99% confidence interval?

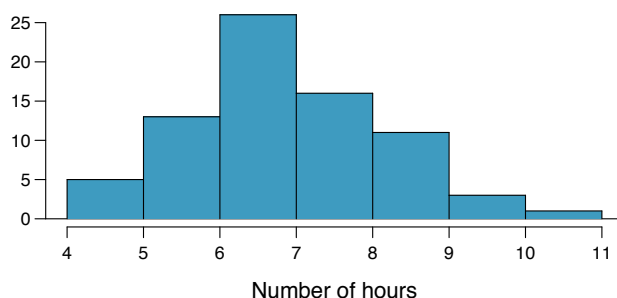
---

<sup>45</sup>ETS, Interpreting your GRE Scores.

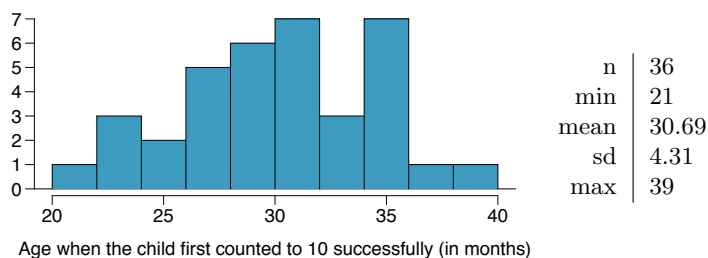
**4.20 Thanksgiving spending, Part II.** Exercise 4.12 provides a 95% confidence interval for the average spending by American adults during the six-day period after Thanksgiving 2009: (\$80.31, \$89.11).

- A local news anchor claims that the average spending during this period in 2009 was \$100. What do you think of this claim?
- Would the news anchor's claim be considered reasonable based on a 90% confidence interval? Why or why not?

**4.21 Ball bearings.** A manufacturer claims that bearings produced by their machine last 7 hours on average under harsh conditions. A factory worker randomly samples 75 ball bearings, and records their lifespans under harsh conditions. He calculates a sample mean of 6.85 hours, and the standard deviation of the data is 1.25 working hours. The following histogram shows the distribution of the lifespans of the ball bearings in this sample. Conduct a formal hypothesis test of this claim. Make sure to check that relevant conditions are satisfied.



**4.22 Gifted children, Part I.** Researchers investigating characteristics of gifted children collected data from schools in a large city on a random sample of thirty-six children who were identified as gifted children soon after they reached the age of four. The following histogram shows the distribution of the ages (in months) at which these children first counted to 10 successfully. Also provided are some sample statistics.<sup>46</sup>



- Are conditions for inference satisfied?
- Suppose you read on a parenting website that children first count to 10 successfully when they are 32 months old, on average. Perform a hypothesis test to evaluate if these data provide convincing evidence that the average age at which gifted children first count to 10 successfully is different than the general average of 32 months. Use a significance level of 0.10.
- Interpret the p-value in context of the hypothesis test and the data.
- Calculate a 90% confidence interval for the average age at which gifted children first count to 10 successfully.
- Do your results from the hypothesis test and the confidence interval agree? Explain.

<sup>46</sup>F.A. Graybill and H.K. Iyer. *Regression Analysis: Concepts and Applications*. Duxbury Press, 1994, pp. 511–516.

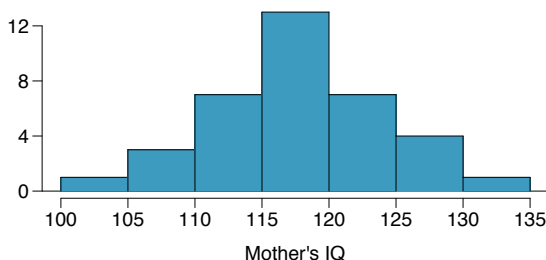
**4.23 Waiting at an ER, Part III.** The hospital administrator mentioned in Exercise 4.11 randomly selected 64 patients and measured the time (in minutes) between when they checked in to the ER and the time they were first seen by a doctor. The average time is 137.5 minutes and the standard deviation is 39 minutes. He is getting grief from his supervisor on the basis that the wait times in the ER increased greatly from last year's average of 127 minutes. However, the administrator claims that the increase is probably just due to chance.

- Are conditions for inference met? Note any assumptions you must make to proceed.
- Using a significance level of  $\alpha = 0.05$ , is the change in wait times statistically significant? Use a two-sided test since it seems the supervisor had to inspect the data before he suggested an increase occurred.
- Would the conclusion of the hypothesis test change if the significance level was changed to  $\alpha = 0.01$ ?

**4.24 Gifted children, Part II.** Exercise 4.22 describes a study on gifted children. In this study, along with variables on the children, the researchers also collected data on the mother's and father's IQ of the 36 randomly sampled gifted children. The histogram below shows the distribution of mother's IQ. Also provided are some sample statistics.

- Perform a hypothesis test to evaluate if these data provide convincing evidence that the average IQ of mothers of gifted children is different than the average IQ for the population at large, which is 100. Use a significance level of 0.10.
- Calculate a 90% confidence interval for the average IQ of mothers of gifted children.
- Do your results from the hypothesis test and the confidence interval agree? Explain.

|      |       |
|------|-------|
| n    | 36    |
| min  | 101   |
| mean | 118.2 |
| sd   | 6.5   |
| max  | 131   |



**4.25 Nutrition labels.** The nutrition label on a bag of potato chips says that a one ounce (28 gram) serving of potato chips has 130 calories and contains ten grams of fat, with three grams of saturated fat. A random sample of 35 bags yielded a sample mean of 134 calories with a standard deviation of 17 calories. Is there evidence that the nutrition label does not provide an accurate measure of calories in the bags of potato chips? We have verified the independence, sample size, and skew conditions are satisfied.

**4.26 Find the sample mean.** You are given the following hypotheses:  $H_0: \mu = 34$ ,  $H_A: \mu > 34$ . We know that the sample standard deviation is 10 and the sample size is 65. For what sample mean would the p-value be equal to 0.05? Assume that all conditions necessary for inference are satisfied.

**4.27 Testing for Fibromyalgia.** A patient named Diana was diagnosed with Fibromyalgia, a long-term syndrome of body pain, and was prescribed anti-depressants. Being the skeptic that she is, Diana didn't initially believe that anti-depressants would help her symptoms. However after a couple months of being on the medication she decides that the anti-depressants are working, because she feels like her symptoms are in fact getting better.

- Write the hypotheses in words for Diana's skeptical position when she started taking the anti-depressants.
- What is a Type 1 error in this context?
- What is a Type 2 error in this context?
- How would these errors affect the patient?

**4.28 Testing for food safety.** A food safety inspector is called upon to investigate a restaurant with a few customer reports of poor sanitation practices. The food safety inspector uses a hypothesis testing framework to evaluate whether regulations are not being met. If he decides the restaurant is in gross violation, its license to serve food will be revoked.

- (a) Write the hypotheses in words.
- (b) What is a Type 1 error in this context?
- (c) What is a Type 2 error in this context?
- (d) Which error is more problematic for the restaurant owner? Why?
- (e) Which error is more problematic for the diners? Why?
- (f) As a diner, would you prefer that the food safety inspector requires strong evidence or very strong evidence of health concerns before revoking a restaurant's license? Explain your reasoning.

**4.29 Errors in drug testing.** Suppose regulators monitored 403 drugs last year, each for a particular adverse response. For each drug they conducted a single hypothesis test with a significance level of 5% to determine if the adverse effect was higher in those taking the drug than those who did not take the drug; the regulators ultimately rejected the null hypothesis for 42 drugs.

- (a) Describe the error the regulators might have made for a drug where the null hypothesis was rejected.
- (b) Describe the error regulators might have made for a drug where the null hypothesis was not rejected.
- (c) Suppose the vast majority of the 403 drugs do not have adverse effects. Then, if you picked one of the 42 suspect drugs at random, about how sure would you be that the drug really has an adverse effect?
- (d) Can you also say how sure you are that a particular drug from the 361 where the null hypothesis was not rejected does not have the corresponding adverse response?

**4.30 Car insurance savings, Part I.** A car insurance company advertises that customers switching to their insurance save, on average, \$432 on their yearly premiums. A market researcher at a competing insurance discounter is interested in showing that this value is an overestimate so he can provide evidence to government regulators that the company is falsely advertising their prices. He randomly samples 82 customers who recently switched to this insurance and finds an average savings of \$395, with a standard deviation of \$102.

- (a) Are conditions for inference satisfied?
- (b) Perform a hypothesis test and state your conclusion.
- (c) Do you agree with the market researcher that the amount of savings advertised is an overestimate? Explain your reasoning.
- (d) Calculate a 90% confidence interval for the average amount of savings of all customers who switch their insurance.
- (e) Do your results from the hypothesis test and the confidence interval agree? Explain.

**4.31 Happy hour.** A restaurant owner is considering extending the happy hour at his restaurant since he would like to see if it increases revenue. If it does, he will permanently extend happy hour. He estimates that the current average revenue per customer is \$18 during happy hour. He runs the extended happy hour for a week and finds an average revenue of \$19.25 with a standard deviation \$3.02 based on a simple random sample of 70 customers.

- Are conditions for inference satisfied?
- Perform a hypothesis test. Suppose the customers and their buying habits this week were no different than in any other week for this particular bar. (This may not always be a reasonable assumption.)
- Calculate a 90% confidence interval for the average revenue per customer.
- Do your results from the hypothesis test and the confidence interval agree? Explain.
- If your hypothesis test and confidence interval suggest a significant increase in revenue per customer, why might you still not recommend that the restaurant owner extend the happy hour based on this criterion? What may be a better measure to consider?

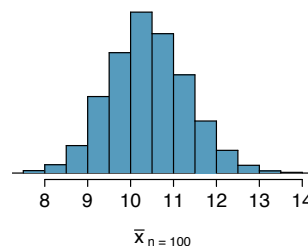
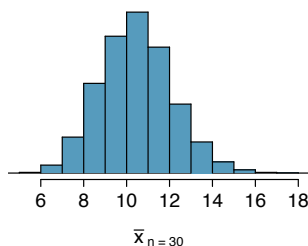
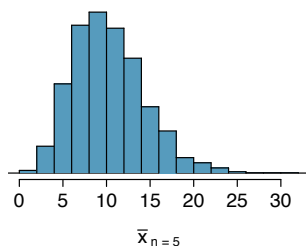
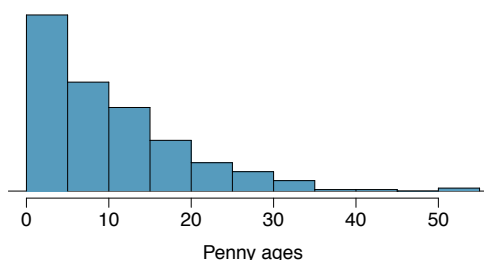
**4.32 Speed reading, Part I.** A company offering online speed reading courses claims that students who take their courses show a 5 times (500%) increase in the number of words they can read in a minute without losing comprehension. A random sample of 100 students yielded an average increase of 415% with a standard deviation of 220%. Is there evidence that the company's claim is false?

- Are conditions for inference satisfied?
- Perform a hypothesis test evaluating if the company's claim is reasonable or if the true average improvement is less than 500%. Make sure to interpret your response in context of the hypothesis test and the data. Use  $\alpha = 0.025$ .
- Calculate a 95% confidence interval for the average increase in the number of words students can read in a minute without losing comprehension.
- Do your results from the hypothesis test and the confidence interval agree? Explain.

### 4.7.4 Examining the Central Limit Theorem

**4.33 Ages of pennies, Part I.** The histogram below shows the distribution of ages of pennies at a bank.

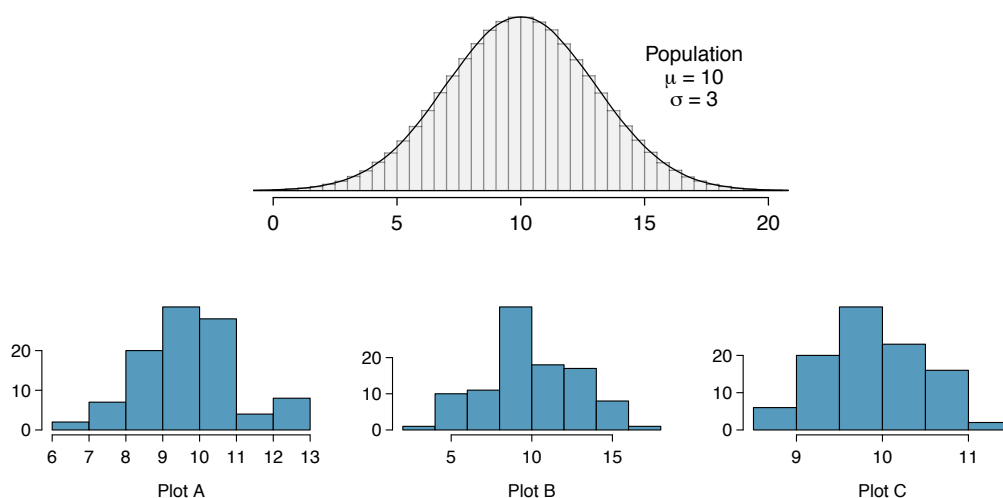
- Describe the distribution.
- Sampling distributions for means from simple random samples of 5, 30, and 100 pennies is shown in the histograms below. Describe the shapes of these distributions and comment on whether they look like what you would expect to see based on the Central Limit Theorem.



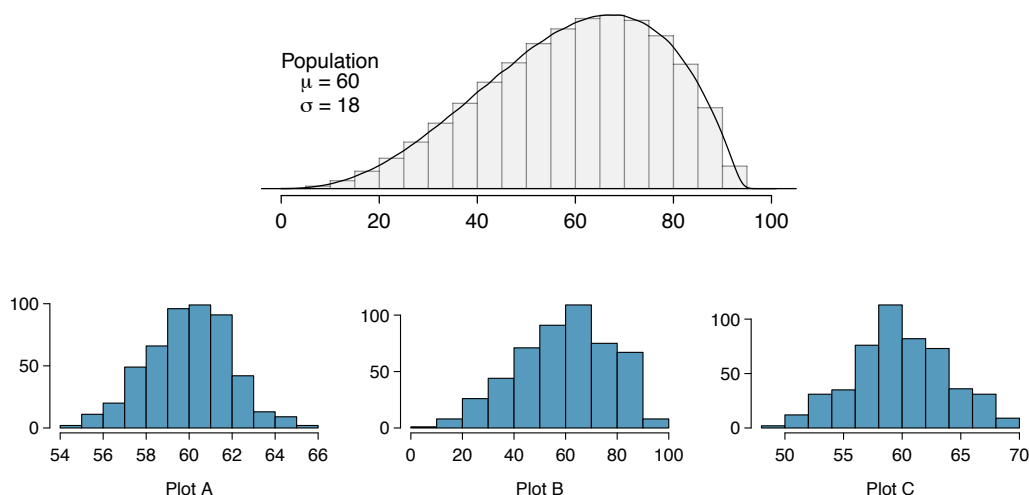


**4.34 Ages of pennies, Part II.** The mean age of the pennies from Exercise 4.33 is 10.44 years with a standard deviation of 9.2 years. Using the Central Limit Theorem, calculate the means and standard deviations of the distribution of the mean from random samples of size 5, 30, and 100. Comment on whether the sampling distributions shown in Exercise 4.33 agree with the values you compute.

**4.35 Identify distributions, Part I.** Four plots are presented below. The plot at the top is a distribution for a population. The mean is 10 and the standard deviation is 3. Also shown below is a distribution of (1) a single random sample of 100 values from this population, (2) a distribution of 100 sample means from random samples with size 5, and (3) a distribution of 100 sample means from random samples with size 25. Determine which plot (A, B, or C) is which and explain your reasoning.



**4.36 Identify distributions, Part II.** Four plots are presented below. The plot at the top is a distribution for a population. The mean is 60 and the standard deviation is 18. Also shown below is a distribution of (1) a single random sample of 500 values from this population, (2) a distribution of 500 sample means from random samples of each size 18, and (3) a distribution of 500 sample means from random samples of each size 81. Determine which plot (A, B, or C) is which and explain your reasoning.



**4.37 Housing prices, Part I.** A housing survey was conducted to determine the price of a typical home in Topanga, CA. The mean price of a house was roughly \$1.3 million with a standard deviation of \$300,000. There were no houses listed below \$600,000 but a few houses above \$3 million.

- (a) Is the distribution of housing prices in Topanga symmetric, right skewed, or left skewed? *Hint:* Sketch the distribution.
- (b) Would you expect most houses in Topanga to cost more or less than \$1.3 million?
- (c) Can we estimate the probability that a randomly chosen house in Topanga costs more than \$1.4 million using the normal distribution?
- (d) What is the probability that the mean of 60 randomly chosen houses in Topanga is more than \$1.4 million?
- (e) How would doubling the sample size affect the standard error of the mean?

**4.38 Stats final scores.** Each year about 1500 students take the introductory statistics course at a large university. This year scores on the final exam are distributed with a median of 74 points, a mean of 70 points, and a standard deviation of 10 points. There are no students who scored above 100 (the maximum score attainable on the final) but a few students scored below 20 points.

- (a) Is the distribution of scores on this final exam symmetric, right skewed, or left skewed?
- (b) Would you expect most students to have scored above or below 70 points?
- (c) Can we calculate the probability that a randomly chosen student scored above 75 using the normal distribution?
- (d) What is the probability that the average score for a random sample of 40 students is above 75?
- (e) How would cutting the sample size in half affect the standard error of the mean?

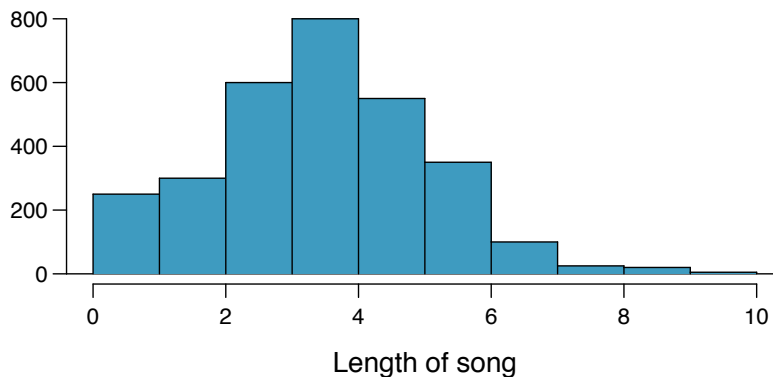
**4.39 Weights of pennies.** The distribution of weights of US pennies is approximately normal with a mean of 2.5 grams and a standard deviation of 0.03 grams.

- (a) What is the probability that a randomly chosen penny weighs less than 2.4 grams?
- (b) Describe the sampling distribution of the mean weight of 10 randomly chosen pennies.
- (c) What is the probability that the mean weight of 10 pennies is less than 2.4 grams?
- (d) Sketch the two distributions (population and sampling) on the same scale.
- (e) Could you estimate the probabilities from (a) and (c) if the weights of pennies had a skewed distribution?

**4.40 CFLs.** A manufacturer of compact fluorescent light bulbs advertises that the distribution of the lifespans of these light bulbs is nearly normal with a mean of 9,000 hours and a standard deviation of 1,000 hours.

- (a) What is the probability that a randomly chosen light bulb lasts more than 10,500 hours?
- (b) Describe the distribution of the mean lifespan of 15 light bulbs.
- (c) What is the probability that the mean lifespan of 15 randomly chosen light bulbs is more than 10,500 hours?
- (d) Sketch the two distributions (population and sampling) on the same scale.
- (e) Could you estimate the probabilities from parts (a) and (c) if the lifespans of light bulbs had a skewed distribution?

**4.41 Songs on an iPod.** Suppose an iPod has 3,000 songs. The histogram below shows the distribution of the lengths of these songs. We also know that, for this iPod, the mean length is 3.45 minutes and the standard deviation is 1.63 minutes.



- (a) Calculate the probability that a randomly selected song lasts more than 5 minutes.
- (b) You are about to go for an hour run and you make a random playlist of 15 songs. What is the probability that your playlist lasts for the entire duration of your run? *Hint:* If you want the playlist to last 60 minutes, what should be the minimum average length of a song?
- (c) You are about to take a trip to visit your parents and the drive is 6 hours. You make a random playlist of 100 songs. What is the probability that your playlist lasts the entire drive?

**4.42 Spray paint.** Suppose the area that can be painted using a single can of spray paint is slightly variable and follows a nearly normal distribution with a mean of 25 square feet and a standard deviation of 3 square feet.

- (a) What is the probability that the area covered by a can of spray paint is more than 27 square feet?
- (b) Suppose you want to spray paint an area of 540 square feet using 20 cans of spray paint. On average, how many square feet must each can be able to cover to spray paint all 540 square feet?
- (c) What is the probability that you can cover a 540 square feet area using 20 cans of spray paint?
- (d) If the area covered by a can of spray paint had a slightly skewed distribution, could you still calculate the probabilities in parts (a) and (c) using the normal distribution?

### 4.7.5 Inference for other estimators

**4.43 Spam mail, Part I.** The 2004 National Technology Readiness Survey sponsored by the Smith School of Business at the University of Maryland surveyed 418 randomly sampled Americans, asking them how many spam emails they receive per day. The survey was repeated on a new random sample of 499 Americans in 2009.<sup>47</sup>

- (a) What are the hypotheses for evaluating if the average spam emails per day has changed from 2004 to 2009.
- (b) In 2004 the mean was 18.5 spam emails per day, and in 2009 this value was 14.9 emails per day. What is the point estimate for the difference between the two population means?
- (c) A report on the survey states that the observed difference between the sample means is not statistically significant. Explain what this means in context of the hypothesis test and the data.
- (d) Would you expect a confidence interval for the difference between the two population means to contain 0? Explain your reasoning.

**4.44 Nearsightedness.** It is believed that nearsightedness affects about 8% of all children. In a random sample of 194 children, 21 are nearsighted.

- (a) Construct hypotheses appropriate for the following question: do these data provide evidence that the 8% value is inaccurate?
- (b) What proportion of children in this sample are nearsighted?
- (c) Given that the standard error of the sample proportion is 0.0195 and the point estimate follows a nearly normal distribution, calculate the test statistic (the Z statistic).
- (d) What is the p-value for this hypothesis test?
- (e) What is the conclusion of the hypothesis test?

**4.45 Spam mail, Part II.** The National Technology Readiness Survey from Exercise 4.43 also asked Americans how often they delete spam emails. 23% of the respondents in 2004 said they delete their spam mail once a month or less, and in 2009 this value was 16%.

- (a) What are the hypotheses for evaluating if the proportion of those who delete their email once a month or less (or never) has changed from 2004 to 2009?
- (b) What is the point estimate for the difference between the two population proportions?
- (c) A report on the survey states that the observed decrease from 2004 to 2009 is statistically significant. Explain what this means in context of the hypothesis test and the data.
- (d) Would you expect a confidence interval for the difference between the two population proportions to contain 0? Explain your reasoning.

**4.46 Unemployment and relationship problems.** A USA Today/Gallup poll conducted between 2010 and 2011 asked a group of unemployed and underemployed Americans if they have had major problems in their relationships with their spouse or another close family member as a result of not having a job (if unemployed) or not having a full-time job (if underemployed). 27% of the 1,145 unemployed respondents and 25% of the 675 underemployed respondents said they had major problems in relationships as a result of their employment status.

- (a) What are the hypotheses for evaluating if the proportions of unemployed and underemployed people who had relationship problems were different?
- (b) The p-value for this hypothesis test is approximately 0.35. Explain what this means in context of the hypothesis test and the data.

---

<sup>47</sup>Rockbridge, 2009 National Technology Readiness Survey SPAM Report.

### 4.7.6 Sample size and power

**4.47 Which is higher?** In each part below, there is a value of interest and two scenarios (I and II). For each part, report if the value of interest is larger under scenario I, scenario II, or whether the value is equal under the scenarios.

- (a) The standard error of  $\bar{x}$  when  $s = 120$  and (I)  $n = 25$  or (II)  $n = 125$ .
- (b) The margin of error of a confidence interval when the confidence level is (I) 90% or (II) 80%.
- (c) The p-value for a Z statistic of 2.5 when (I)  $n = 500$  or (II)  $n = 1000$ .
- (d) The probability of making a Type 2 error when the alternative hypothesis is true and the significance level is (I) 0.05 or (II) 0.10.

**4.48 True or false.** Determine if the following statements are true or false, and explain your reasoning. If false, state how it could be corrected.

- (a) If a given value (for example, the null hypothesized value of a parameter) is within a 95% confidence interval, it will also be within a 99% confidence interval.
- (b) Decreasing the significance level ( $\alpha$ ) will increase the probability of making a Type 1 error.
- (c) Suppose the null hypothesis is  $\mu = 5$  and we fail to reject  $H_0$ . Under this scenario, the true population mean is 5.
- (d) If the alternative hypothesis is true, then the probability of making a Type 2 error and the power of a test add up to 1.
- (e) With large sample sizes, even small differences between the null value and the true value of the parameter, a difference often called the effect size, will be identified as statistically significant.
- (f) A cutoff of  $\alpha = 0.05$  is the ideal value for all hypothesis tests.

**4.49 Car insurance savings, Part II.** The market researcher from Exercise 4.30 collected data about the savings of 82 customers at a competing car insurance company. The mean and standard deviation of this sample are \$395 and \$102, respectively. He would like to conduct another survey but have a margin of error of no more than \$10 at a 99% confidence level. How large of a sample should he collect?

**4.50 Speed reading, Part II.** A random sample of 100 students who took online speed reading courses from the company described in Exercise 4.32 yielded an average increase in reading speed of 415% and a standard deviation of 220%. We would like to calculate a 95% confidence interval for the average increase in reading speed with a margin of error of no more than 15%. How many students should we sample?

**4.51 Waiting at the ER, Part IV.** Exercise 4.23 introduced us to a hospital where ER wait times were being analyzed. The previous year's average was 128 minutes. Suppose that this year's average wait time is 135 minutes.

- (a) Provide the hypotheses for this situation in plain language.
- (b) If we plan to collect a sample size of  $n = 64$ , what values could  $\bar{x}$  take so that we reject  $H_0$ ? Suppose the sample standard deviation from the earlier exercise (39 minutes) is the population standard deviation. You may assume that the conditions for the nearly normal model for  $\bar{x}$  are satisfied.
- (c) Calculate the probability of a Type 2 error.