+

# Chapter 9: Testing a Claim

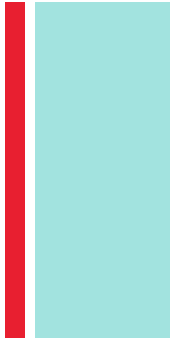**Section 9.3**
**Tests About a Population Mean**

**The Practice of Statistics, 4th edition – For AP***
**STARNES, YATES, MOORE**

**+**

# Chapter 9
# Testing a Claim

# Section 9.3
# Tests About a Population Mean

After this section, you should be able to…

- ✓ CHECK conditions for carrying out a test about a population mean.

- ✓ CONDUCT a one-sample $t$ test about a population mean.

- ✓ CONSTRUCT a confidence interval to draw a conclusion for a two-sided test about a population mean.

- ✓ PERFORM significance tests for paired data.

# ■ **Introduction**

Confidence intervals and significance tests for a population proportion $p$ are based on $z$-values from the standard Normal distribution.

Inference about a population mean $\mu$ uses a $t$ distribution with $n - 1$ degrees of freedom, except in the rare case when the population standard deviation $\sigma$ is known.

We learned how to construct confidence intervals for a population mean in Section 8.3. Now we'll examine the details of testing a claim about an unknown parameter $\mu$.

# ■ Carrying Out a Significance Test for *μ*

In an earlier example, a company claimed to have developed a new AAA battery that lasts longer than its regular AAA batteries. Based on years of experience, the company knows that its regular AAA batteries last for 30 hours of continuous use, on average. An SRS of 15 new batteries lasted an average of 33.9 hours with a standard deviation of 9.8 hours. Do these data give *convincing evidence* that the new batteries last longer on average?

To find out, we must perform a significance test of

$$H_0: \mu = 30 \text{ hours}$$
$$H_a: \mu > 30 \text{ hours}$$

where $\mu$ = the true mean lifetime of the new deluxe AAA batteries.

## Check Conditions:

Three conditions should be met before we perform inference for an unknown population mean: Random, Normal, and Independent. The Normal condition for means is

Population distribution is Normal or sample size is large ($n \geq 30$)

We often don't know whether the population distribution is Normal. But if the sample size is large ($n \geq 30$), we can safely carry out a significance test (due to the central limit theorem). If the sample size is small, we should examine the sample data for any obvious departures from Normality, such as skewness and outliers.
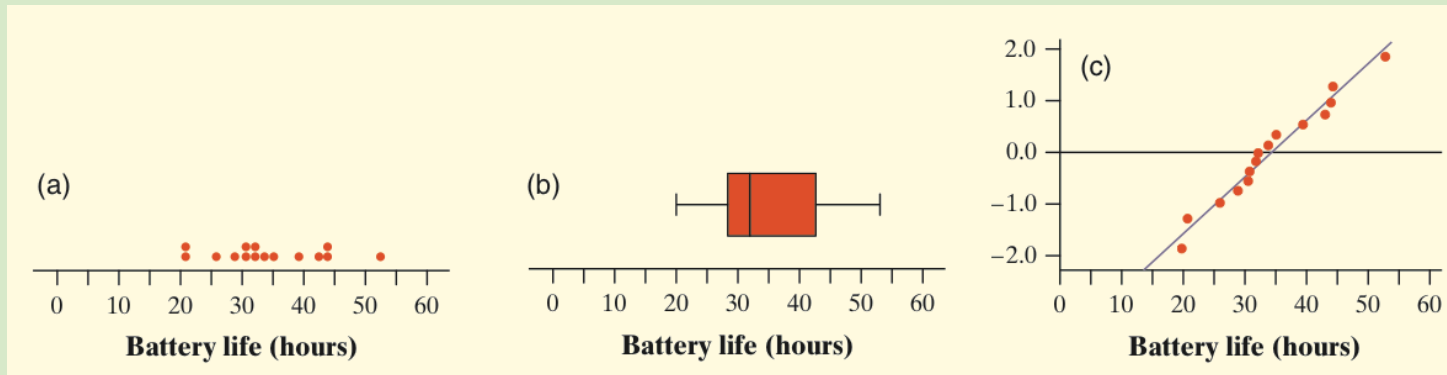
# ■ Carrying Out a Significance Test for *μ*

## Check Conditions:

Three conditions should be met before we perform inference for an unknown population mean: Random, Normal, and Independent.

✓ **Random** The company tests an SRS of 15 new AAA batteries.

✓ **Normal** We don't know if the population distribution of battery lifetimes for the company's new AAA batteries is Normal. With such a small sample size (*n* = 15), we need to inspect the data for any departures from Normality.
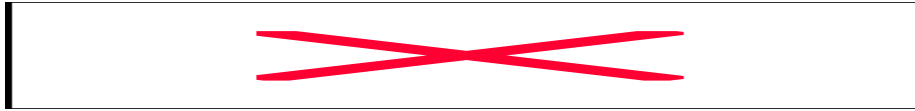


The dotplot and boxplot show slight right-skewness but no outliers. The Normal probability plot is close to linear. We should be safe performing a test about the population mean lifetime *μ*.

✓ **Independent** Since the batteries are being sampled without replacement, we need to check the *10% condition*: there must be at least 10(15) = 150 new AAA batteries. This seems reasonable to believe.
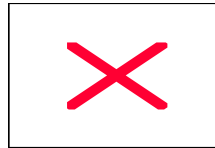
# ■ **Carrying Out a Significance Test**

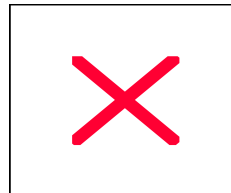**Calculations: Test statistic and P-value**

When performing a significance test, we do calculations assuming that the null hypothesis $H_0$ is true. The test statistic measures how far the sample result diverges from the parameter value specified by $H_0$, in standardized units. As before,

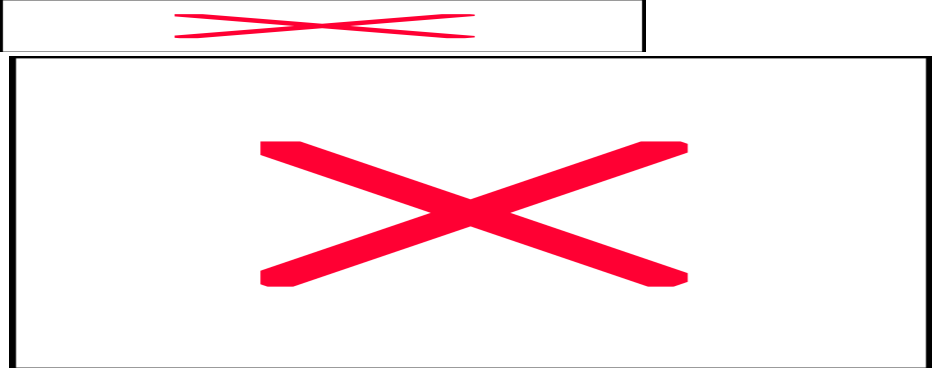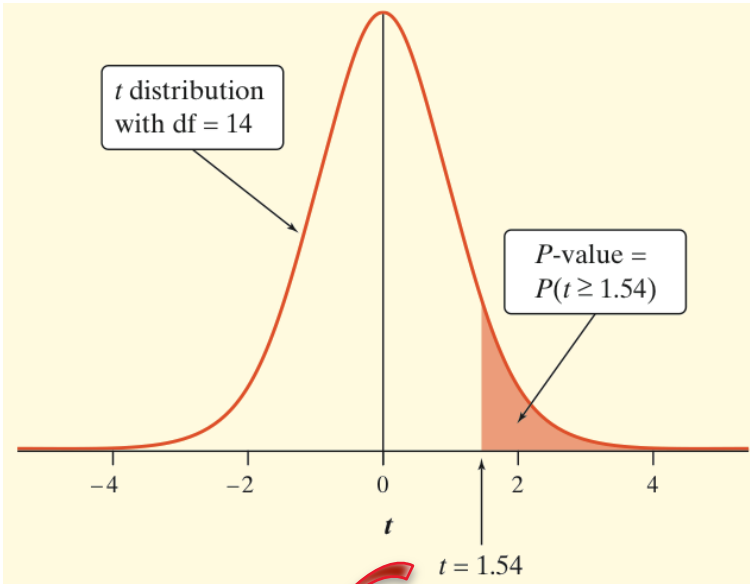For a test of $H_0: \mu = \mu_0$, our statistic is the sample mean. Its standard deviation is

Because the population standard deviation σ is usually unknown, we use the sample standard deviation $s_x$ in its place. The resulting test statistic has the standard error of the sample mean in the denominator

When the Normal condition is met, this statistic has a $t$ distribution with $n - 1$ degrees of freedom.

# ■ **Carrying Out a Hypothesis Test**

The battery company wants to test $H_0$: $\mu = 30$ versus $H_a$: $\mu > 30$ based on an SRS of 15 new AAA batteries with mean lifetime and standard deviation



The *P*-value is the probability of getting a result this large or larger in the direction indicated by $H_a$, that is, $P(t \geq 1.54)$.

**Upper-tail probability *p***

| df | .10 | .05 | .025 |
|----|------|------|------|
| 13 | 1.350 | 1.771 | 2.160 |
| 14 | 1.345 | 1.761 | 2.145 |
| 15 | 1.341 | 1.753 | 3.131 |
| | 80% | 90% | 95% |

**Confidence level *C***

✓ Go to the *df* = 14 row.

✓ Since the *t* statistic falls between the values 1.345 and 1.761, the "Upper-tail probability *p*" is between 0.10 and 0.05.

✓ The *P*-value for this test is between 0.05 and 0.10.

**Because the *P*-value exceeds our default α = 0.05 significance level, we can't conclude that the company's new AAA batteries last longer than 30 hours, on average.**
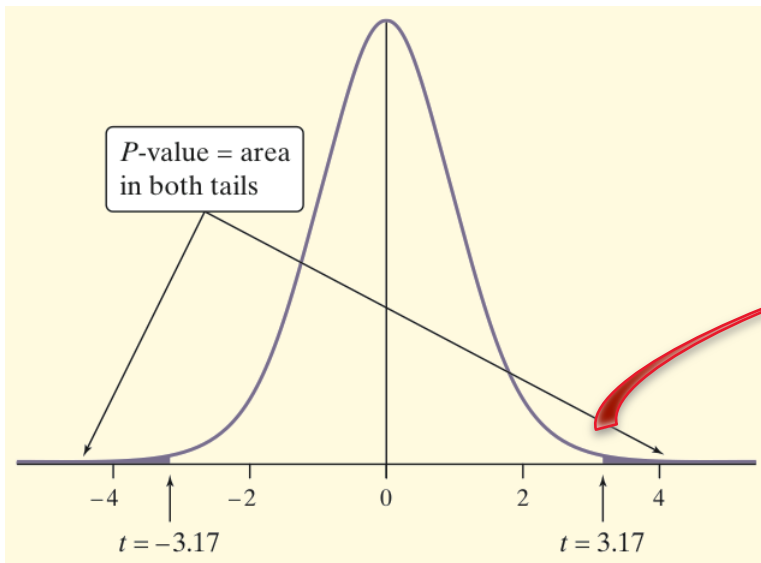
Tests About a Population Mean

# ■ Using Table B Wisely

• Table B gives a range of possible *P*-values for a significance. We can still draw a conclusion from the test in much the same way as if we had a single probability by comparing the range of possible *P*-values to our desired significance level.

• Table B has other limitations for finding P-values. It includes probabilities only for *t* distributions with degrees of freedom from 1 to 30 and then skips to *df* = 40, 50, 60, 80, 100, and 1000. (The bottom row gives probabilities for *df* = ∞, which corresponds to the standard Normal curve.) *Note: If the df you need isn't provided in Table B, use the next lower df that is available.*

• Table B shows probabilities only for positive values of *t*. To find a *P*-value for a negative value of *t*, we use the symmetry of the *t* distributions.

# ■ **Using Table B Wisely**

Suppose you were performing a test of $H_0$: $\mu = 5$ versus $H_a$: $\mu \neq 5$ based on a sample size of $n = 37$ and obtained $t = -3.17$. Since this is a two-sided test, you are interested in the probability of getting a value of $t$ less than -3.17 or greater than 3.17.

Due to the symmetric shape of the density curve, $P(t \leq -3.17) = P(t \geq 3.17)$. Since Table B shows only positive $t$-values, we must focus on $t = 3.17$.



**Upper-tail probability $p$**

| df | .005 | .0025 | .001 |
|----|------|-------|------|
| 29 | 2.756 | 3.038 | 3.396 |
| 30 | 2.750 | 3.030 | 3.385 |
| 40 | 2.704 | 2.971 | 3.307 |
|    | 99% | 99.5% | 99.8% |

**Confidence level $C$**

Since $df = 37 - 1 = 36$ is not available on the table, move across the $df = 30$ row and notice that $t = 3.17$ falls between 3.030 and 3.385.
The corresponding "Upper-tail probability $p$" is between 0.0025 and 0.001. For this two-sided test, the corresponding $P$-value would be between 2(0.001) = 0.002 and 2(0.0025) = 0.005.

# The One-Sample *t* Test

When the conditions are met, we can test a claim about a population mean $\mu$ using a **one-sample t test**.

| One-Sample *t* Test |
| --- |

Choose an SRS of size *n* from a large population that contains an unknown mean $\mu$. To test the hypothesis $H_0 : \mu = \mu_0$, compute the one-sample *t* statistic

Find the ... large or la ... distribu ...

**Use this test only when
(1) the population distribution is
Normal or the sample is large
($n \geq 30$), and (2) the population is at
least 10 times as large as the
sample.**

$H_a : \mu > \mu_0$    $H_a : \mu \neq \mu_0$

$-|t|$    $|t|$

# Example: Healthy Streams

The level of dissolved oxygen (DO) in a stream or river is an important indicator of the water's ability to support aquatic life. A researcher measures the DO level at 15 randomly chosen locations along a stream. Here are the results in milligrams per liter:

| 4.53 | 5.04 | 3.29 | 5.23 | 4.13 | 5.50 | 4.83 | 4.40 |
|------|------|------|------|------|------|------|------|
| 5.42 | 6.38 | 4.01 | 4.66 | 2.87 | 5.73 | 5.55 | |

A dissolved oxygen level below 5 mg/l puts aquatic life at risk.

**State:** We want to perform a test at the $\alpha$ = 0.05 significance level of

$$H_0: \mu = 5$$
$$H_a: \mu < 5$$

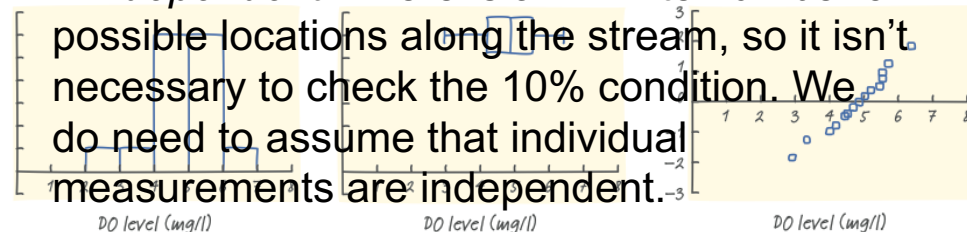where $\mu$ is the actual mean dissolved oxygen level in this stream.

**Plan:** If conditions are met, we should do a one-sample $t$ test for $\mu$.

✓*Random* The researcher measured the DO level at 15 randomly chosen locations.

✓*Normal* We don't know whether the population distribution of DO levels at all points along the stream is Normal. With such a small sample size ($n$ = 15), we need to look at the data to see if it's safe to use $t$ procedures.

The histogram looks roughly symmetric; the boxplot shows no outliers; and the Normal probability plot is fairly linear. With no outliers or strong skewness, the $t$ procedures should be pretty accurate even if the population distribution isn't Normal.

✓*Independent* There is an infinite number of possible locations along the stream, so it isn't necessary to check the 10% condition. We do need to assume that individual measurements are independent.

# Example: Healthy Streams

**Do:** The sample mean and standard deviation are

**P-value** The *P*-value is the area to the left of $t = -0.94$ under the *t* distribution curve with df = 15 − 1 = 14.

**Conclude:** The *P*-value, is between 0.15 and 0.20. Since this is greater than our α = 0.05 significance level, we fail to reject $H_0$. We don't have enough evidence to conclude that the mean DO level in the stream is less than 5 mg/l.



**Upper-tail probability *p***

| df | .25 | .20 | .15 |
|----|------|------|-------|
| 13 | .694 | .870 | 1.079 |
| 14 | .692 | .868 | 1.076 |
| 15 | .691 | .866 | 1.074 |
|    | 50%  | 60%  | 70%   |

**Confidence level *C***

*Since we decided not to reject $H_0$, we could have made a Type II error (failing to reject $H_0$ when $H_0$ is false). If we did, then the mean dissolved oxygen level μ in the stream is actually less than 5 mg/l, but we didn't detect that with our significance test.*

# ■ Two-Sided Tests

At the Hawaii Pineapple Company, managers are interested in the sizes of the pineapples grown in the company's fields. Last year, the mean weight of the pineapples harvested from one large field was 31 ounces. A new irrigation system was installed in this field after the growing season. Managers wonder whether this change will affect the mean weight of future pineapples grown in the field. To find out, they select and weigh a random sample of 50 pineapples from this year's crop. The Minitab output below summarizes the data. Determine whether there are any outliers.

## Descriptive Statistics: Weight (oz)

| Variable | N | Mean | SE Mean | StDev | Minimum | Q1 | Median | Q3 | Maximum |
|----------|---|------|---------|-------|---------|----|--------|----|---------|
| Weight (oz) | 50 | 31.935 | 0.339 | 2.394 | 26.491 | 29.990 | 31.739 | 34.115 | 35.547 |

✓ $IQR = Q_3 - Q_1 = 34.115 - 29.990 = 4.125$

✓ Any data value greater than $Q_3 + 1.5(IQR)$ or less than $Q_1 - 1.5(IQR)$ is considered an outlier.

$$Q_3 + 1.5(IQR) = 34.115 + 1.5(4.125) = 40.3025$$
$$Q_1 - 1.5(IQR) = 29.990 - 1.5(4.125) = 23.0825$$

✓ Since the maximum value 35.547 is less than 40.3025 and the minimum value 26.491 is greater than 23.0825, there are no outliers.

# ■ **Two-Sided Tests**

**State:** We want to test the hypotheses

$$H_0: \mu = 31$$
$$H_a: \mu \neq 31$$

where $\mu$ = the mean weight (in ounces) of all pineapples grown in the field this year.  Since no significance level is given, we'll use $\alpha = 0.05$.

**Plan:** If conditions are met, we should do a one-sample $t$ test for $\mu$.

✓*Random*  The data came from a random sample of 50 pineapples from this year's crop.

✓*Normal*  We don't know whether the population distribution of pineapple weights this year is Normally distributed. But $n = 50 \geq 30$, so the large sample size (and the fact that there are no outliers) makes it OK to use $t$ procedures.

✓*Independent*  There need to be at least 10(50) = 500 pineapples in the field because managers are sampling without replacement (*10% condition*). We would expect many more than 500 pineapples in a "large field."

# Two-Sided Tests

**Do:** The sample mean and standard deviation are



t distribution, 49 degrees of freedom

P-value = 0.0081

Values of t

t = −2.762          t = 2.762

**Upper-tail probability p**

| df | .005 | .0025 | .001 |
|---|---|---|---|
| **30** | 2.750 | 3.030 | 3.385 |
| **40** | 2.704 | 2.971 | 3.307 |
| **50** | 2.678 | 2.937 | 3.261 |
| | 99% | 99.5% | 99.8% |

**Confidence level C**

**P-value** The P-value for this two-sided test is the area under the t distribution curve with 50 - 1 = 49 degrees of freedom. Since Table B does not have an entry for df = 49, we use the more conservative df = 40. The upper tail probability is between 0.005 and 0.0025 so the desired P-value is between 0.01 and 0.005.

**Conclude:** Since the P-value is between 0.005 and 0.01, it is less than our $\alpha$ = 0.05 significance level, so we have enough evidence to reject $H_0$ and conclude that the mean weight of the pineapples in this year's crop is not 31 ounces.

# ■ Confidence Intervals Give More Information

Minitab output for a significance test and confidence interval based on the pineapple data is shown below. The test statistic and *P*-value match what we got earlier (up to rounding).

### One-Sample T: Weight (oz)

```
Test of mu = 31 vs not = 31

Variable       N    Mean   StDev  SE Mean        95% CI          T     P
Weight (oz)   50  31.935   2.394   0.339  (31.255, 32.616)  2.76  0.008
```
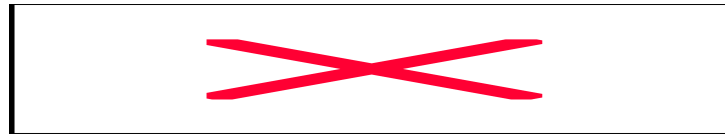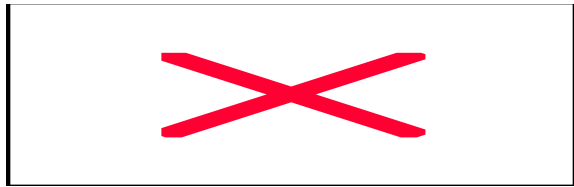
*The 95% confidence interval for the mean weight of all the pineapples grown in the field this year is 31.255 to 32.616 ounces. We are 95% confident that this interval captures the true mean weight μ of this year's pineapple crop.*

**As with proportions, there is a link between a two-sided test at significance level α and a 100(1 – α)% confidence interval for a population mean μ.**

For the pineapples, the two-sided test at $\alpha = 0.05$ rejects $H_0$: $\mu = 31$ in favor of $H_a$: $\mu \neq 31$. The corresponding 95% confidence interval does not include 31 as a plausible value of the parameter $\mu$. In other words, the test and interval lead to the same conclusion about $H_0$. But the confidence interval provides much more information: *a set of plausible values for the population mean*.

# ■ Confidence Intervals and Two-Sided Tests

The connection between two-sided tests and confidence intervals is even stronger for means than it was for proportions. That's because both inference methods for means use the standard error of the sample mean in the calculations.

✓ A two-sided test at significance level $\alpha$ (say, $\alpha$ = 0.05) and a 100(1 − $\alpha$)% confidence interval (a 95% confidence interval if $\alpha$ = 0.05) give similar information about the population parameter.

✓ When the two-sided significance test at level α rejects $H_0$: $\mu = \mu_0$, the 100(1 − $\alpha$)% confidence interval for $\mu$ will not contain the hypothesized value $\mu_0$ .

✓ When the two-sided significance test at level $\alpha$ fails to reject the null hypothesis, the confidence interval for $\mu$ will contain $\mu_0$ .

# ■ Inference for Means: Paired Data

Comparative studies are more convincing than single-sample investigations. For that reason, one-sample inference is less common than comparative inference. Study designs that involve making two observations on the same individual, or one observation on each of two similar individuals, result in **paired data**.

When paired data result from measuring the same quantitative variable twice, as in the job satisfaction study, we can make comparisons by analyzing the differences in each pair. If the conditions for inference are met, we can use one-sample *t* procedures to perform inference about the mean difference $\mu_d$.

These methods are sometimes called **paired *t* procedures**.

# ■ **Paired *t* Test**

Researchers designed an experiment to study the effects of caffeine withdrawal. They recruited 11 volunteers who were diagnosed as being caffeine dependent to serve as subjects. Each subject was barred from coffee, colas, and other substances with caffeine for the duration of the experiment. During one two-day period, subjects took capsules containing their normal caffeine intake. During another two-day period, they took placebo capsules. The order in which subjects took caffeine and the placebo was randomized. At the end of each two-day period, a test for depression was given to all 11 subjects. Researchers wanted to know whether being deprived of caffeine would lead to an increase in depression.

| Results of a caffeine deprivation study | | | |
|---|---|---|---|
| Subject | Depression (caffeine) | Depression (placebo) | Difference (placebo – caffeine) |
| 1 | 5 | 16 | 11 |
| 2 | 5 | 23 | 18 |
| 3 | 4 | 5 | 1 |
| 4 | 3 | 7 | 4 |
| 5 | 8 | 14 | 6 |
| 6 | 5 | 24 | 19 |
| 7 | 0 | 6 | 6 |
| 8 | 0 | 3 | 3 |
| 9 | 2 | 15 | 13 |
| 10 | 11 | 12 | 1 |
| 11 | 1 | 0 | - 1 |

**State:** If caffeine deprivation has no effect on depression, then we would expect the actual mean difference in depression scores to be 0. We want to test the hypotheses
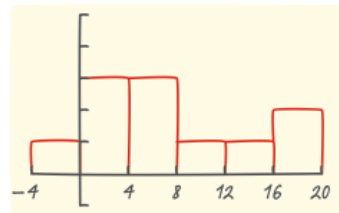
$$H_0: \mu_d = 0$$
$$H_a: \mu_d > 0$$

where $\mu_d$ = the true mean difference (placebo – caffeine) in depression score. Since no significance level is given, we'll use $\alpha = 0.05$.
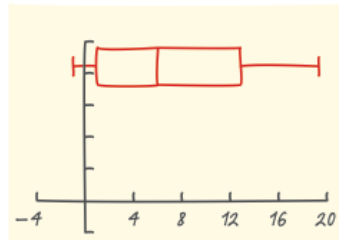
# ■ **Paired *t* Test**

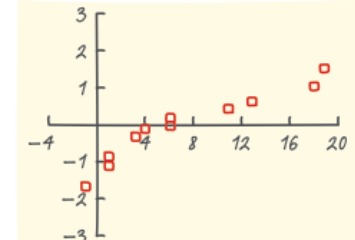**Plan:** If conditions are met, we should do a paired *t* test for $\mu_d$.

✓*Random* researchers randomly assigned the treatment order—placebo then caffeine, caffeine then placebo—to the subjects.

✓*Normal* We don't know whether the actual distribution of difference in depression scores (placebo - caffeine) is Normal. With such a small sample size (*n* = 11), we need to examine the data to see if it's safe to use *t* procedures.



Change in depression
(placebo – caffeine)
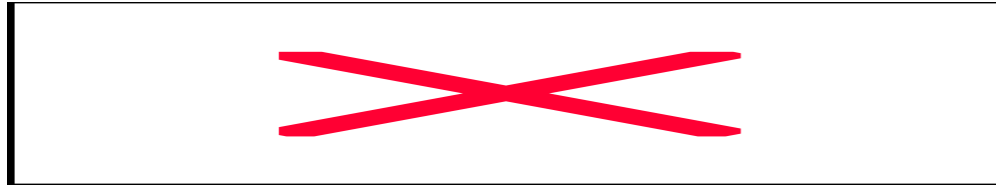
Change in depression
(placebo – caffeine)

Change in depression
(placebo – caffeine)

The histogram has an irregular shape with so few values; the boxplot shows some right-skewness but not outliers; and the Normal probability plot looks fairly linear. With no outliers or strong skewness, the *t* procedures should be pretty accurate.

✓*Independent* We aren't sampling, so it isn't necessary to check the *10% condition*. We will assume that the changes in depression scores for individual subjects are independent. This is reasonable if the experiment is conducted properly.

# ■ **Paired _t_ Test**

**Do:** The sample mean and standard deviation are

**P-value** According to technology, the area to the right of _t_ = 3.53 on the _t_ distribution curve with df = 11 – 1 = 10 is 0.0027.

**Conclude:** With a _P_-value of 0.0027, which is much less than our chosen $\alpha$ = 0.05, we have convincing evidence to reject $H_0$: $\mu_d$ = 0. We can therefore conclude that depriving these caffeine-dependent subjects of caffeine caused an average increase in depression scores.

# ■ Using Tests Wisely

Significance tests are widely used in reporting the results of research in many fields. New drugs require significant evidence of effectiveness and safety. Courts ask about statistical significance in hearing discrimination cases. Marketers want to know whether a new ad campaign significantly outperforms the old one, and medical researchers want to know whether a new therapy performs significantly better. In all these uses, statistical significance is valued because it points to an effect that is unlikely to occur simply by chance.

Carrying out a significance test is often quite simple, especially if you use a calculator or computer. Using tests wisely is not so simple. Here are some points to keep in mind when using or interpreting significance tests.

**Statistical Significance and Practical Importance**
When a null hypothesis ("no effect" or "no difference") can be rejected at the usual levels ($\alpha = 0.05$ or $\alpha = 0.01$), there is good evidence of a difference. But that difference may be very small. When large samples are available, even tiny deviations from the null hypothesis will be significant.

# ■ Using Tests Wisely

**Don't Ignore Lack of Significance**
There is a tendency to infer that there is no difference whenever a *P*-value fails to attain the usual 5% standard. In some areas of research, small differences that are detectable only with large sample sizes can be of great practical significance. When planning a study, verify that the test you plan to use has a high probability (power) of detecting a difference of the size you hope to find.

**Statistical Inference Is Not Valid for All Sets of Data**
Badly designed surveys or experiments often produce invalid results. Formal statistical inference cannot correct basic flaws in the design. Each test is valid only in certain circumstances, with properly produced data being particularly important.

**Beware of Multiple Analyses**
Statistical significance ought to mean that you have found a difference that you were looking for. The reasoning behind statistical significance works well if you decide what difference you are seeking, design a study to search for it, and use a significance test to weigh the evidence you get. In other settings, significance may have little meaning.

# Section 9.3
# Tests About a Population Mean

**Summary**

In this section, we learned that…

- ✓ Significance tests for the mean $\mu$ of a Normal population are based on the sampling distribution of the sample mean. Due to the central limit theorem, the resulting procedures are approximately correct for other population distributions when the sample is large.

- ✓ If we somehow know σ, we can use a $z$ test statistic and the standard Normal distribution to perform calculations. In practice, we typically do not know σ. Then, we use the **one-sample $t$ statistic**



with $P$-values calculated from the $t$ distribution with $n$ - 1 degrees of freedom.

# Section 9.3
# Tests About a Population Mean

## Summary

✓ The **one-sample *t* test** is approximately correct when

  **Random** The data were produced by random sampling or a randomized experiment.

  **Normal** The population distribution is Normal OR the sample size is large ($n \geq 30$).

  **Independent** Individual observations are independent. When sampling without replacement, check that the population is at least 10 times as large as the sample.

✓ Confidence intervals provide additional information that significance tests do not—namely, a range of plausible values for the parameter $\mu$. A two-sided test of $H_0$: $\mu = \mu_0$ at significance level α gives the same conclusion as a $100(1 - \alpha)\%$ confidence interval for $\mu$.

✓ Analyze **paired data** by first taking the difference within each pair to produce a single sample. Then use one-sample *t* procedures.

**+**

# Section 9.3
# Tests About a Population Mean

## Summary

✓ Very small differences can be highly significant (small *P*-value) when a test is based on a large sample. A statistically significant difference need not be practically important.

✓ Lack of significance does not imply that $H_0$ is true. Even a large difference can fail to be significant when a test is based on a small sample.

✓ Significance tests are not always valid. Faulty data collection, outliers in the data, and other practical problems can invalidate a test. Many tests run at once will probably produce some significant results by chance alone, even if all the null hypotheses are true.

# **Looking Ahead…**

We'll learn how to compare two populations or groups.

We'll learn about
- ✓ **Comparing Two Proportions**
- ✓ **Comparing Two Means**