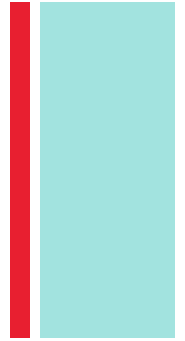


## Unit 5: Hypothesis Testing

The Practice of Statistics, 4<sup>th</sup> edition – For AP\*  
STARNES, YATES, MOORE

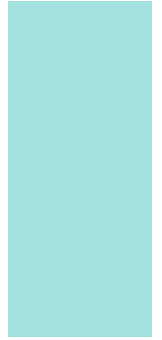
# + Unit 5: Hypothesis Testing



- **9.1**     **Significance Tests: The Basics**
- **9.2**     Tests about a Population Proportion
- **9.3**     Tests about a Population Mean
- **9.1 and 9.2** Errors and the Power of a Test



## Section 9.1 Significance Tests: The Basics



### Learning Objectives

After this section, you should be able to...

- ✓ STATE correct hypotheses for a significance test about a population proportion or mean.
- ✓ INTERPRET  $P$ -values in context.

## ■ Introduction

Two of the most common types of inference:

### **Confidence intervals**

Used when your goal is to estimate a population parameter.

### **Significance tests**

Used to assess the evidence provided by data about some claim concerning a population.

A **significance test** is a formal procedure for comparing observed data with a claim (also called a hypothesis) whose truth we want to assess. The claim is a statement about a parameter, like the population proportion  $p$  or the population mean  $\mu$ . We express the results of a significance test in terms of a probability that measures how well the data and the claim agree.

## ■ The Reasoning of Significance Tests

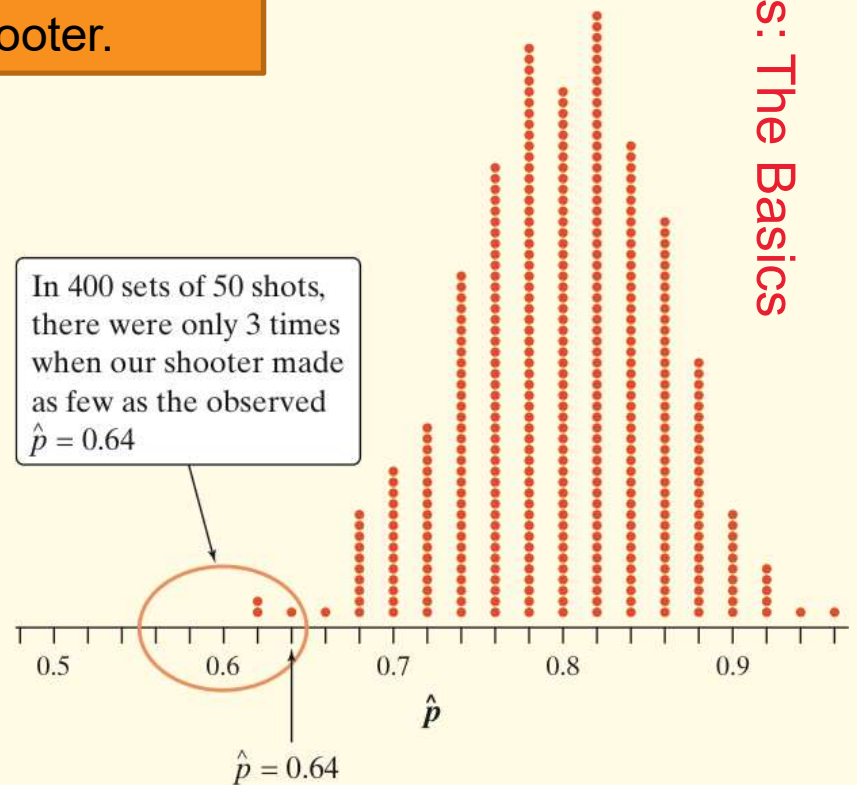
Suppose a basketball player claimed to be an 80% free-throw shooter. To test this claim, we have him attempt 50 free-throws. He makes 32 of them. His sample proportion of made shots is  $32/50 = 0.64$ .

What can we conclude about the claim based on this sample data?

We can use software to simulate 400 sets of 50 shots assuming that the player is really an 80% shooter.

You can say how strong the evidence against the player's claim is by giving the probability that he would make as few as 32 out of 50 free throws if he really makes 80% in the long run.

The observed statistic is so unlikely if the actual parameter value is  $p = 0.80$  that it gives convincing evidence that the player's claim is not true.



## ■ The Reasoning of Significance Tests

Based on the evidence, we might conclude the player's claim is incorrect.

In reality, there are two possible explanations for the fact that he made only 64% of his free throws.

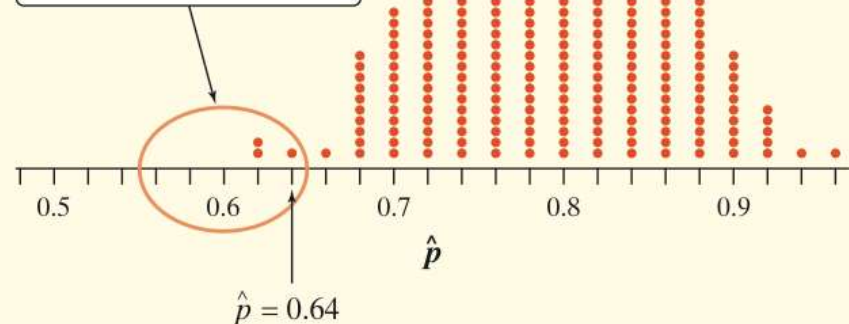
1) The player's claim is correct ( $p = 0.8$ ), and by bad luck, a very unlikely outcome occurred.

2) The population proportion is actually less than 0.8, so the sample result is not an unlikely outcome.

### Basic Idea

*An outcome that would rarely happen if a claim were true is good evidence that the claim is not true.*

In 400 sets of 50 shots, there were only 3 times when our shooter made as few as the observed  $\hat{p} = 0.64$



## ■ Stating Hypotheses

A significance test starts with a careful statement of claims.

The **null hypothesis**.

This claim is a statement of “no difference.”

The **alternative hypothesis**.

The claim we hope or suspect to be true instead of the null.

### Definition:

The claim tested by a statistical test is called the **null hypothesis ( $H_0$ )**. The test is designed to assess the strength of the evidence against the null hypothesis. Often the null hypothesis is a statement of “no difference.”

The claim about the population that we are trying to find evidence for is the **alternative hypothesis ( $H_a$ )**.

In the free-throw shooter example, our hypotheses are

$$H_0 : p = 0.80$$

$$H_a : p < 0.80$$

where  $p$  is the long-run proportion of made free throws.

## ■ Stating Hypotheses

In any significance test, the null hypothesis has the form

$$H_0 : \text{parameter} = \text{value}$$

The alternative hypothesis has one of the forms

$$H_a : \text{parameter} < \text{value}$$

$$H_a : \text{parameter} > \text{value}$$

$$H_a : \text{parameter} \neq \text{value}$$

To determine the correct form of  $H_a$ , read the problem carefully.

### Definition:

The alternative hypothesis is **one-sided** if it states that a parameter is *larger than* the null hypothesis value or if it states that the parameter is *smaller than* the null value.

It is **two-sided** if it states that the parameter is *different* from the null hypothesis value (it could be either larger or smaller).

- ✓ Hypotheses always refer to a *population*, not to a sample. Be sure to state  $H_0$  and  $H_a$  in terms of *population parameters*.
- ✓ It is *never* correct to write a hypothesis about a sample statistic, such as  $\hat{p} = 0.64$  or  $\bar{x} = 85$ .



## ■ Example: Studying Job Satisfaction

Does the job satisfaction of assembly-line workers differ when their work is machine-paced rather than self-paced? One study chose 18 subjects at random from a company with over 200 workers who assembled electronic devices. Half of the workers were assigned at random to each of two groups. Both groups did similar assembly work, but one group was allowed to pace themselves while the other group used an assembly line that moved at a fixed pace. After two weeks, all the workers took a test of job satisfaction. Then they switched work setups and took the test again after two more weeks. The response variable is the difference in satisfaction scores, self-paced minus machine-paced.

### a) Describe the parameter of interest in this setting.

The parameter of interest is the mean  $\mu$  of the differences (*self-paced minus machine-paced*) in job satisfaction scores in the population of all assembly-line workers at this company.

### b) State appropriate hypotheses for performing a significance test.

Because the initial question asked whether job satisfaction differs, the alternative hypothesis is two-sided; that is, either  $\mu < 0$  or  $\mu > 0$ . For simplicity, we write this as  $\mu \neq 0$ . That is,

$$H_0: \mu = 0$$

$$H_a: \mu \neq 0$$

## ■ Interpreting $P$ -Values

The null hypothesis  $H_0$  states the claim that we are seeking evidence against. The probability that measures the strength of the evidence against a null hypothesis is called a  **$P$ -value**.

### Definition:

The probability, computed assuming  $H_0$  is true, that the statistic would take a value as extreme as or more extreme than the one actually observed is called the  **$P$ -value** of the test. The smaller the  $P$ -value, the stronger the evidence against  $H_0$  provided by the data.

- ✓ Small  $P$ -values are evidence against  $H_0$  because they say that the observed result is unlikely to occur when  $H_0$  is true.
- ✓ Large  $P$ -values fail to give convincing evidence against  $H_0$  because they say that the observed result is likely to occur by chance when  $H_0$  is true.

## ■ Example: Studying Job Satisfaction

For the job satisfaction study, the hypotheses are

$$H_0: \mu = 0$$

$$H_a: \mu \neq 0$$

Data from the 18 workers gave  $\bar{x} = 17$  and  $s_x = 60$ . That is, these workers rated the self-paced environment, on average, 17 points higher. Researchers performed a significance test using the sample data and found a  $P$ -value of 0.2302.

### a) Explain what it means for the null hypothesis to be true in this setting.

In this setting,  $H_0: \mu = 0$  says that the mean difference in satisfaction scores (*self-paced - machine-paced*) for the entire population of assembly-line workers at the company is 0. If  $H_0$  is true, then the workers don't favor one work environment over the other, on average.

### b) Interpret the $P$ -value in context.

✓ The  $P$ -value is the probability of observing a sample result as extreme or more extreme in the direction specified by  $H_a$  just by chance when  $H_0$  is actually true.

Because the alternative hypothesis is two-sided, the  $P$ -value is the probability of getting a value of  $\bar{x}$  as far from 0 in either direction as the observed  $\bar{x} = 17$  when  $H_0$  is true. That is, an average difference of 17 or more points between the two work environments would happen 23% of the time just by chance in random samples of 18 assembly-line workers when the true population mean is  $\mu = 0$ .

## ■ Example: Studying Job Satisfaction

For the job satisfaction study, the hypotheses are

$$H_0: \mu = 0$$

$$H_a: \mu \neq 0$$

Data from the 18 workers gave  $\bar{x} = 17$  and  $s_x = 60$ . That is, these workers rated the self-paced environment, on average, 17 points higher. Researchers performed a significance test using the sample data and found a  $P$ -value of 0.2302.

### a) Explain what it means for the null hypothesis to be true in this setting.

In this setting,  $H_0: \mu = 0$  says that the mean difference in satisfaction scores (*self-paced - machine-paced*) for the entire population of assembly-line workers at the company is 0. If  $H_0$  is true, then the workers don't favor one work environment over the other, on average.

### b) Interpret the $P$ -value in context.

**Results as extreme or more as the ones we got could have happened by chance alone 23% of the time; when it is true that there is no difference between job satisfaction for self-paced versus machine-paced work.**

## ■ Statistical Significance

The final step is to draw a conclusion about the competing claims.

Make one of two decisions:

### **Reject $H_0$**

If your sample result is too unlikely to have happened by chance alone when assuming  $H_0$  is true.

### **Fail to reject $H_0$ .**

If your sample result is likely to have happened by chance alone when assuming  $H_0$  is true.

**Note:** A fail-to-reject  $H_0$  decision in a significance test doesn't mean that  $H_0$  is true. For that reason, you should never “accept  $H_0$ ” or use language implying that you believe  $H_0$  is true.

In a nutshell, our conclusion in a significance test comes down to

$P$ -value small  $\rightarrow$  reject  $H_0 \rightarrow$  conclude  $H_a$  (in context)

$P$ -value large  $\rightarrow$  fail to reject  $H_0 \rightarrow$  cannot conclude  $H_a$  (in context)

## ■ Statistical Significance

There is no rule for how small a  $P$ -value we should require in order to reject  $H_0$  — it's a matter of judgment and depends on the specific circumstances. But we can compare the  $P$ -value with a fixed value that we regard as decisive, called the **significance level**. We write it as  $\alpha$ , the Greek letter alpha. When our  $P$ -value is less than the chosen  $\alpha$ , we say that the result is **statistically significant**.

### Definition:

If the  $P$ -value is smaller than alpha, we say that the data are **statistically significant at level  $\alpha$** . In that case, we reject the null hypothesis  $H_0$  and conclude that there is convincing evidence in favor of the alternative hypothesis  $H_a$ .

When we use a fixed level of significance to draw a conclusion in a significance test,

$P\text{-value} < \alpha \rightarrow \text{reject } H_0 \rightarrow \text{conclude } H_a \text{ (in context)}$

$P\text{-value} \geq \alpha \rightarrow \text{fail to reject } H_0 \rightarrow \text{cannot conclude } H_a \text{ (in context)}$

## ■ Example: Better Batteries

A company has developed a new deluxe AAA battery that is supposed to last longer than its regular AAA battery. However, these new batteries are more expensive to produce, so the company would like to be convinced that they really do last longer. Based on years of experience, the company knows that its regular AAA batteries last for 30 hours of continuous use, on average. The company selects an SRS of 15 new batteries and uses them continuously until they are completely drained. A significance test is performed using the hypotheses

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

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

where  $\mu$  is the true mean lifetime of the new deluxe AAA batteries. The resulting  $P$ -value is 0.0276.

### a) What conclusion can you make for the significance level $\alpha = 0.05$ ?

Since the  $P$ -value, 0.0276, is less than  $\alpha = 0.05$ , the sample result is statistically significant at the 5% level. We have sufficient evidence to reject  $H_0$  and conclude that the company's deluxe AAA batteries last longer than 30 hours, on average.

### b) What conclusion can you make for the significance level $\alpha = 0.01$ ?

Since the  $P$ -value, 0.0276, is greater than  $\alpha = 0.01$ , the sample result is not statistically significant at the 1% level. We do not have enough evidence to reject  $H_0$  in this case. therefore, we cannot conclude that the deluxe AAA batteries last longer than 30 hours, on average.