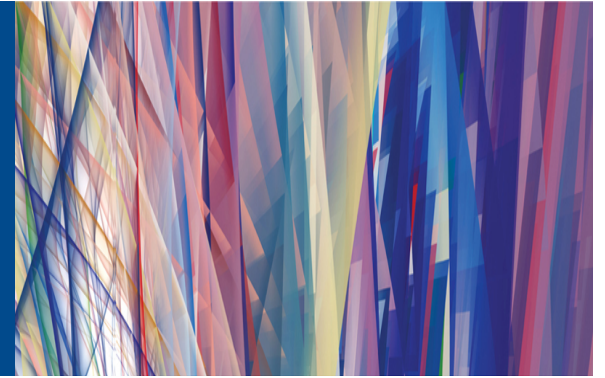# Chapter 6

# Introduction to Inference

**Introduction to the Practice of**
**STATISTICS** SEVENTH EDITION

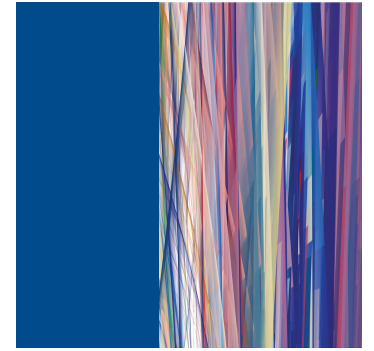**Moore / McCabe / Craig**

**Lecture Presentation Slides**

# Chapter 6
# Introduction to Inference

**6.1 Estimating with Confidence** *Confidence Intervals*

**6.2 Tests of Significance** *(Hypothesis Tests)*

**6.3 Use and Abuse of Tests**

**6.4 Power and Inference as a Decision**
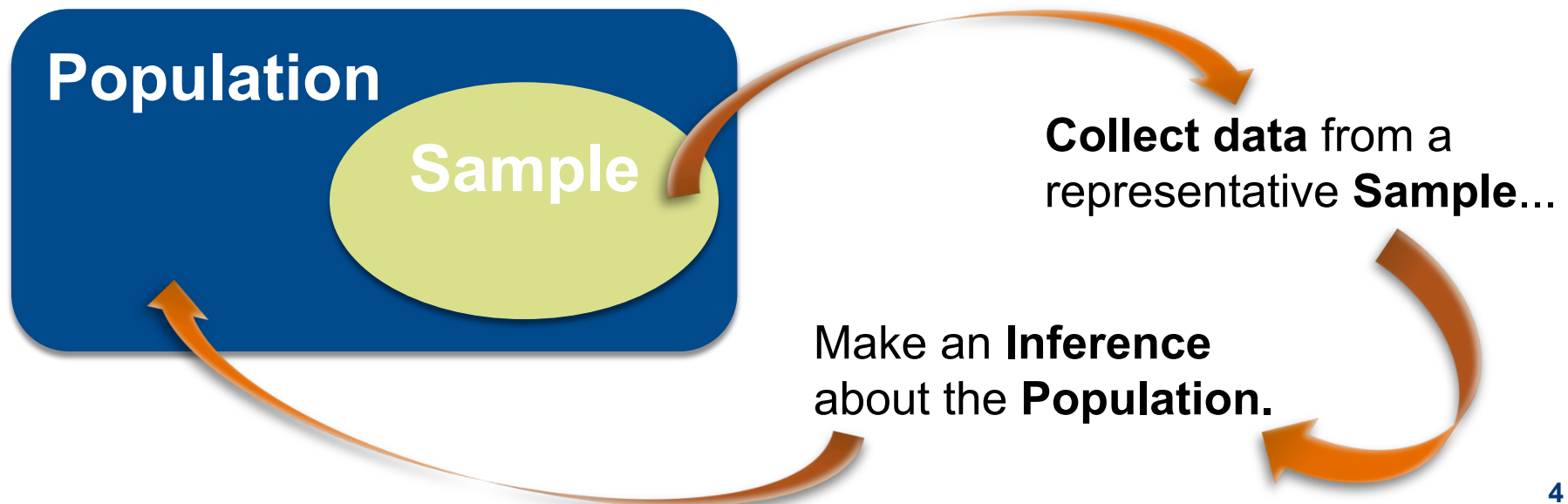
# 6.1 Estimating with Confidence

- Inference

- Statistical Confidence

- Confidence Intervals

- Confidence Interval for a Population Mean

- Choosing the Sample Size

# Statistical Inference

After we have selected a sample, we know the responses of the individuals in the sample. However, the reason for taking the sample is to infer from that data some conclusion about the wider population represented by the sample.

**Statistical inference** provides methods for drawing conclusions about a population from sample data.

**Population**

**Sample**

**Collect data** from a representative **Sample**...

Make an **Inference** about the **Population.**

# Simple Conditions for Inference About a Mean

This chapter presents the basic reasoning of statistical inference. We start with a setting that is too simple to be realistic.

**Simple Conditions for Inference About a Mean**

1. We have an SRS from the population of interest. There is no nonresponse or other practical difficulty.
2. The variable we measure has an exactly Normal distribution $N(\mu,\sigma)$ in the population.
3. We don't know the population mean $\mu$, but we do know the population standard deviation $\sigma$.

**Note:** The conditions that we have a perfect SRS, that the population is exactly Normal, and that we know the population standard deviation are all unrealistic.

# Statistical Estimation

Suppose your instructor has selected a "Mystery Mean" value $\mu$ and stored it as "M" in their calculator. The following command was executed on their calculator: `mean(randNorm(M,20,16))`

```
mean(randNorm(M,
20,16))
              240.79
```

The result was `240.79.` This tells us the calculator chose an SRS of 16 observations from a Normal population with mean M and standard deviation 20. The resulting sample mean of those 16 values was 240.79.

We want to determine an interval of *reasonable* values for the population mean $\mu$. We can use the result above and what we learned about sampling distributions in the previous chapters.
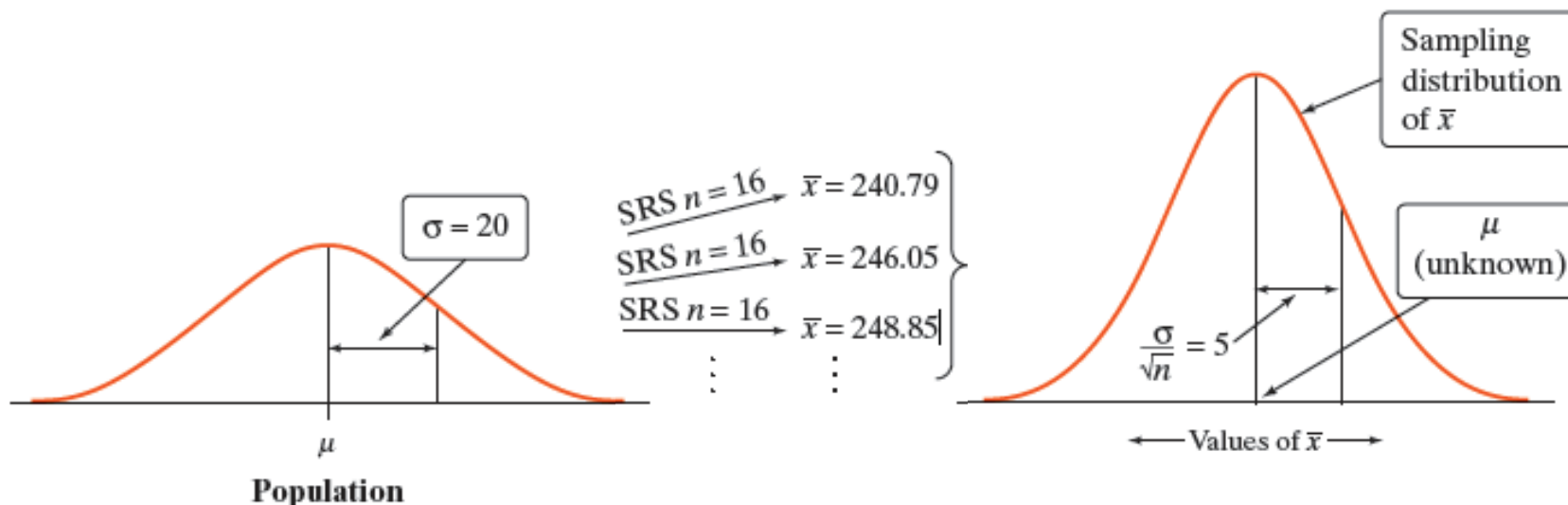
# Statistical Estimation

Since the sample mean is 240.79, we could guess that $\mu$ is "somewhere" around 240.79.  **How close to 240.79 is $\mu$ likely to be?**
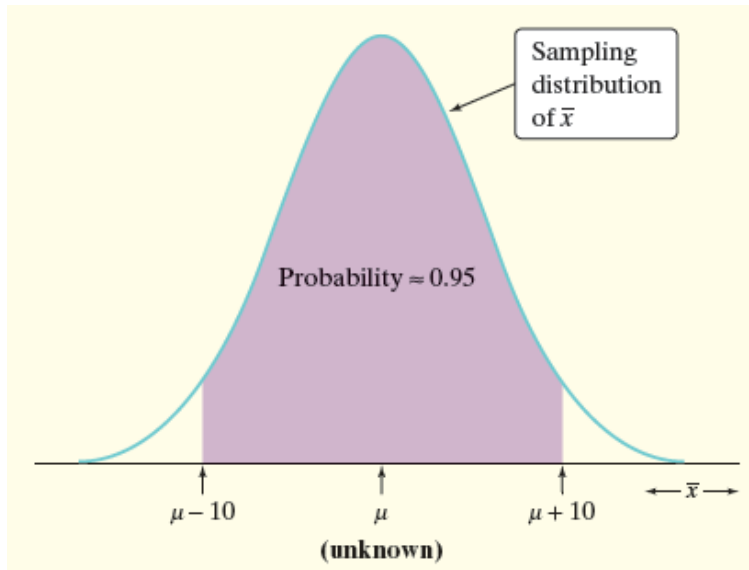
```
mean(randNorm(M,
20,16))
            240.79
```

To answer this question, we must ask another:

**How would the sample mean $\bar{x}$ vary if we took many SRSs of size 16 from the population?**

$\sigma = 20$

SRS $n = 16$ → $\bar{x} = 240.79$

SRS $n = 16$ → $\bar{x} = 246.05$

SRS $n = 16$ → $\bar{x} = 248.85$

Sampling distribution of $\bar{x}$

$\mu$ (unknown)

$\dfrac{\sigma}{\sqrt{n}} = 5$

← Values of $\bar{x}$ →

$\mu$

**Population**

# Statistical Estimation



✓ In repeated samples, the values of the sample mean will follow a Normal distribution with mean $\mu$ and standard deviation 5.

✓ The 68-95-99.7 Rule tells us that in 95% of all samples of size 16, the sample mean will be within 10 (two standard deviations) of $\mu$.

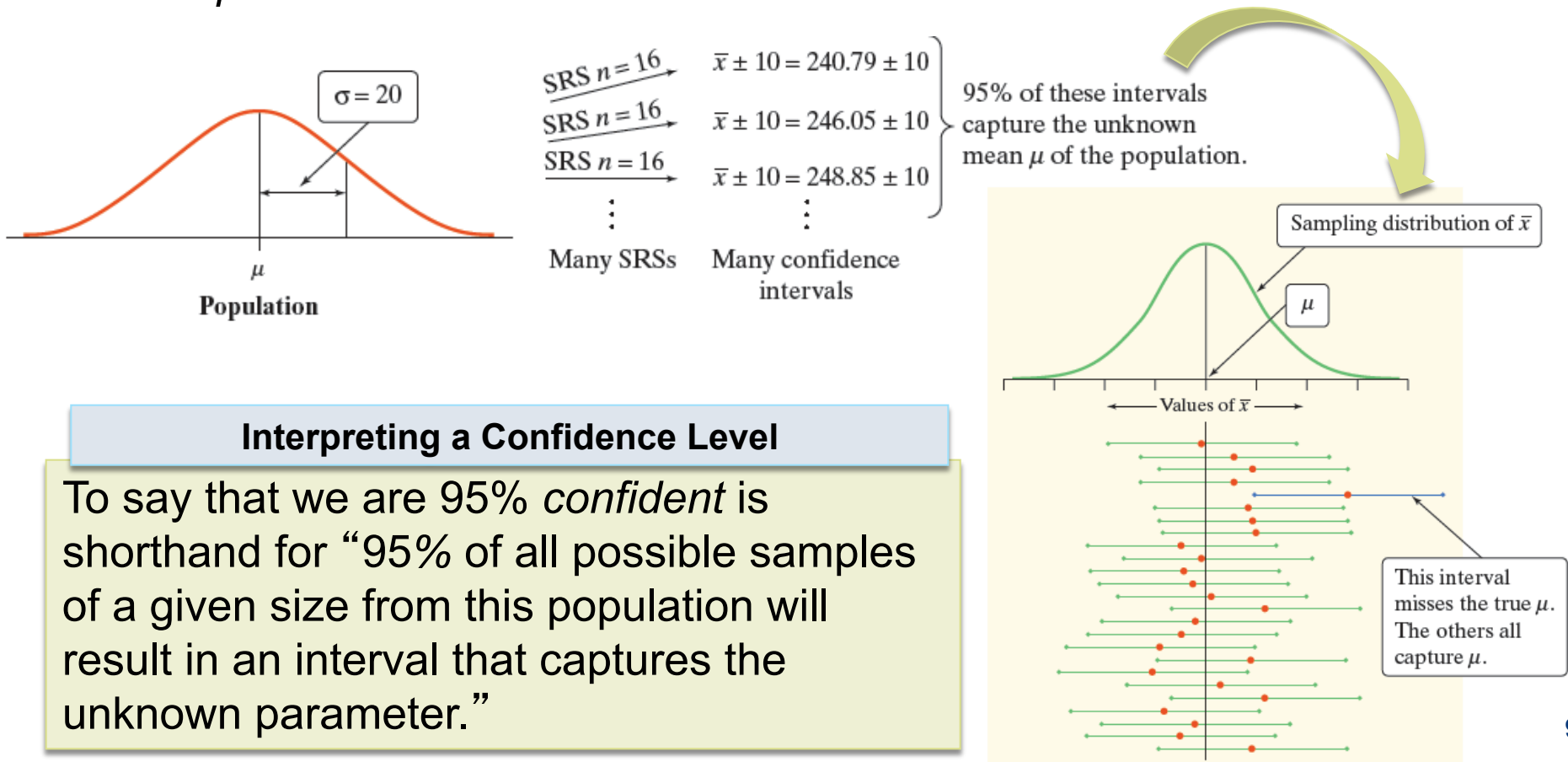✓ If the sample mean is within 10 points of $\mu$, then $\mu$ is within 10 points of the sample mean.

✓ Therefore, the interval from 10 points below to 10 points above the sample mean will "capture" $\mu$ in about 95% of all samples of size 16.

If we estimate that $\mu$ lies somewhere in the interval **230.79** to **250.79,** we'd be calculating an interval using a method that captures the true $\mu$ in about 95% of all possible samples of this size.

# Confidence Level

The confidence level is the overall capture rate if the method is used many times. The sample mean will vary from sample to sample, but when we use the method *estimate ± margin of error* to get an interval based on each sample, C% of these intervals capture the unknown population mean *μ*.



$\sigma = 20$

$\mu$

Population

$$\text{SRS } n = 16 \quad \bar{x} \pm 10 = 240.79 \pm 10$$
$$\text{SRS } n = 16 \quad \bar{x} \pm 10 = 246.05 \pm 10$$
$$\text{SRS } n = 16 \quad \bar{x} \pm 10 = 248.85 \pm 10$$

Many SRSs    Many confidence intervals

95% of these intervals capture the unknown mean $\mu$ of the population.

Sampling distribution of $\bar{x}$

$\mu$

Values of $\bar{x}$

This interval misses the true $\mu$. The others all capture $\mu$.

### Interpreting a Confidence Level

To say that we are 95% *confident* is shorthand for "95% of all possible samples of a given size from this population will result in an interval that captures the unknown parameter."

9

Pop'n : $\mu$ (unknown)
$\sigma$ (known)
$\quad\quad\longrightarrow$ SRS of $n \longrightarrow \bar{X} = 50$

(4

# Confidence Interval

**The Big Idea:** The sampling distribution of $\bar{x}$ tells us how close to $\mu$ the sample mean $\bar{x}$ is likely to be. All confidence intervals we construct will have a form similar to this:

**estimate ± margin of error**

A **level C confidence interval** for a parameter has two parts:
- An **interval** calculated from the data, which has the form:
estimate ± margin of error

- A **confidence level *C*,** which gives the probability that the interval will capture the true parameter value in repeated samples. That is, the confidence level is the success rate for the method.
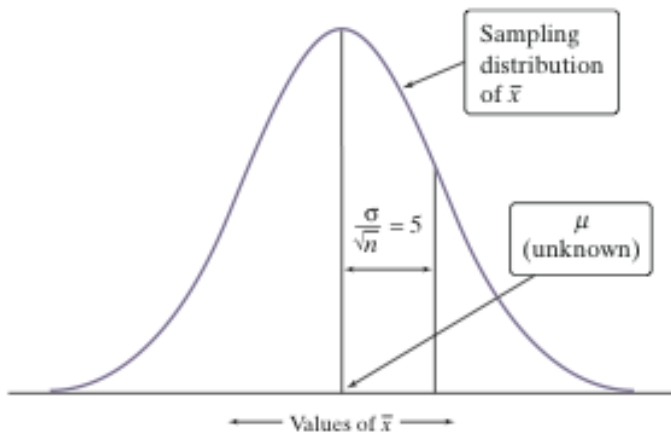
We usually choose a confidence level of 90% or higher because we want to be quite sure of our conclusions. The most common confidence level is 95%.

# Confidence Interval for a Population Mean

Previously, we estimated the "mystery mean" $\mu$ by constructing a confidence interval using the sample mean = 240.79.

To calculate a 95% confidence interval for $\mu$, we use the formula:
estimate ± (critical value) • (standard deviation of statistic)



$$\bar{x} \pm z^* \cdot \frac{\sigma}{\sqrt{n}} = 240.79 \pm 1.96 \cdot \frac{20}{\sqrt{16}}$$

$$= 240.79 \pm 9.8$$

$$= (230.99, 250.59)$$

**Confidence Interval for the Mean of a Normal Population**

Choose an SRS of size $n$ from a population having unknown mean $\mu$ and known standard deviation $\sigma$. A level **C confidence interval for $\mu$** is:

$$\bar{x} \pm z^* \frac{\sigma}{\sqrt{n}}$$

The critical value $z^*$ is found from the standard Normal distribution.

# Finding Specific *z\** Values

We can use a table of *z/t* values (Table D). For a particular confidence level, *C*, the appropriate *z\** value is just above it.

| *z\** | 0.674 | 0.841 | 1.036 | 1.282 | 1.645 | 1.960 | 2.054 | 2.326 | 2.576 | 2.807 | 3.091 | 3.291 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 50% | 60% | 70% | 80% | 90% | 95% | 96% | 98% | 99% | 99.5% | 99.8% | 99.9% |

Confidence level *C*

**Example:** For a 98% confidence level, *z\** = 2.326.

We can also use software. In **Excel:**

=NORMINV(probability,mean,standard_dev)

*gives z for a given cumulative probability.*

Since we want the middle *C* probability, the probability we require is (1 - C)/2

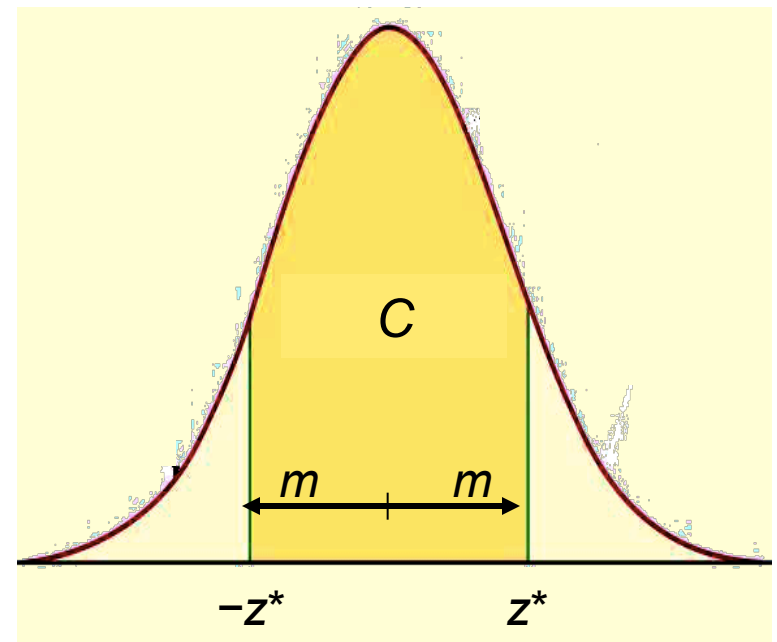**Example:** For a 98% confidence level, =NORMINV(.01,0,1) = −2.32635

# The Margin of Error

The confidence level *C* determines the value of *z\** (in table C).

The margin of error also depends on *z\**.

$$m = z^* \sigma / \sqrt{n}$$

Higher confidence **C** implies a larger margin of error **m** (thus less precision in our estimates).

A lower confidence level **C** produces a smaller margin of error **m** (thus better precision in our estimates).

# How Confidence Intervals Behave

The *z* confidence interval for the mean of a Normal population illustrates several important properties that are shared by all confidence intervals in common use.

- The user chooses the confidence level and the margin of error follows.
- We would like high confidence and a small margin of error.
  - High confidence suggests our method almost always gives correct answers.
  - A small margin of error suggests we have pinned down the parameter precisely.

The margin of error for the *z* confidence interval is:

$$z * \cdot \frac{\sigma}{\sqrt{n}}$$

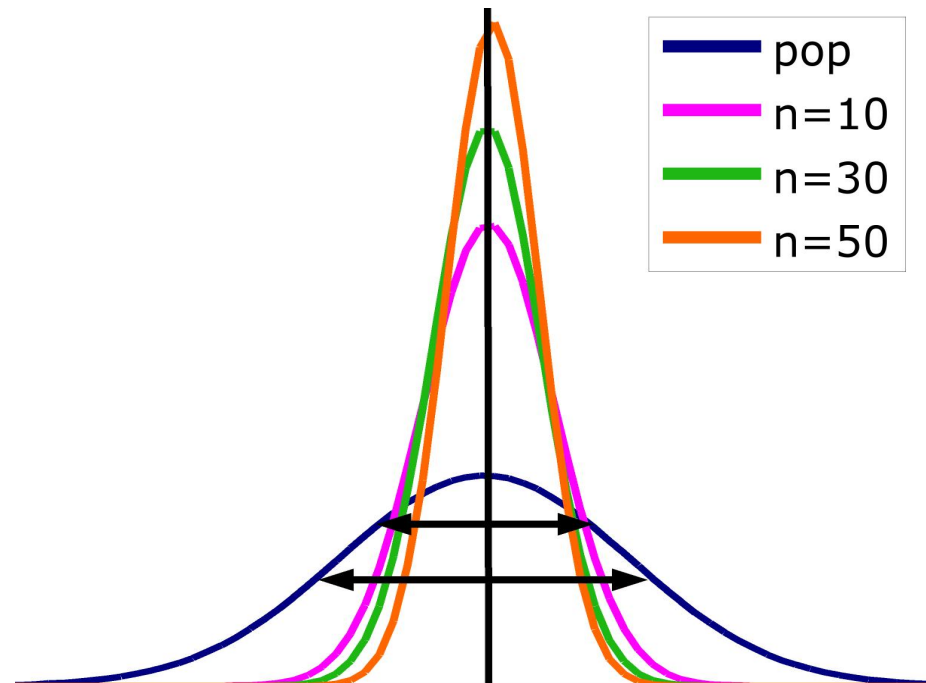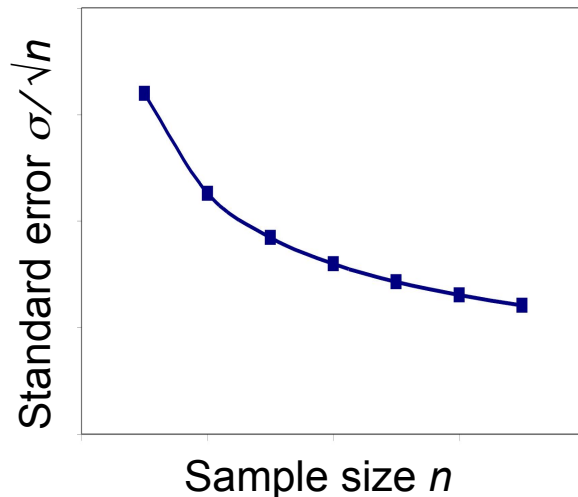The margin of error gets smaller when:

- *z*\* gets smaller (the same as a lower confidence level *C*).
- *σ* is smaller. It is easier to pin down *μ* when *σ* is smaller.
- *n* gets larger. Since *n* is under the square root sign, we must take four times as many observations to cut the margin of error in half.

# Impact of sample size

The spread in the sampling distribution of the mean is a function of the number of individuals per sample.

- The larger the sample size, the smaller the standard deviation (spread) of the sample mean distribution.

- The spread decreases at a rate equal to $\sqrt{n}$.

# Choosing the Sample Size

You may need a certain margin of error (e.g., drug trial, manufacturing specs). In many cases, the population variability ($\sigma$) is fixed, but we can choose the number of measurements ($n$).

The confidence interval for a population mean will have a specified margin of error $m$ when the sample size is:

$$m = z^* \frac{\sigma}{\sqrt{n}} \qquad \Leftrightarrow \qquad n = \left( \frac{z^* \sigma}{m} \right)^2$$

*Remember, though, that sample size is not always stretchable at will. There are typically costs and constraints associated with large samples. The best approach is to use the smallest sample size that can give you useful results.*

# Example

**Density of bacteria in solution:**

Measurement equipment has standard deviation $\sigma$ = 1 * 10⁶ bacteria/ml fluid.

How many measurements should you make to obtain a margin of error of at most 0.5 * 10⁶ bacteria/ml with a confidence level of 95%?

For a 95% confidence interval, $z*$ = 1.96.

$$n = \left(\frac{z*\sigma}{m}\right)^2 \quad \Rightarrow \quad n = \left(\frac{1.96*1}{0.5}\right)^2 = 3.92^2 = 15.3664$$

Using only 15 measurements will not be enough to ensure that $m$ is no more than 0.5 * 10⁶. Therefore, we need at least 16 measurements.
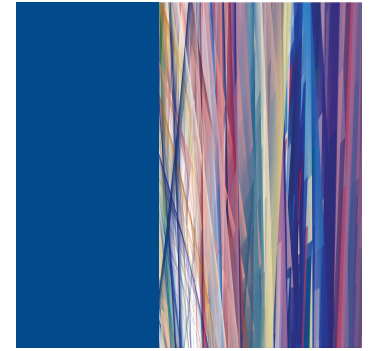
| $z*$ | 0.674 | 0.841 | 1.036 | 1.282 | 1.645 | 1.960 | 2.054 | 2.326 | 2.576 | 2.807 | 3.091 | 3.291 |
|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
|      | 50%   | 60%   | 70%   | 80%   | 90%   | 95%   | 96%   | 98%   | 99%   | 99.5% | 99.8% | 99.9% |

Confidence level C

# Some Cautions

- The data should be a **SRS** from the population.

- The formula is **not** correct for other sampling designs.

- Inference **cannot** rescue badly produced data.

- Confidence intervals are **not resistant** to outliers.

- If *n* is small **(<15)** and the population is not normal, the true confidence level will be **different** from *C.*

- The standard deviation σ of the population must be known.

- The margin of error in a confidence interval covers only random sampling errors!

# 6.2 Tests of Significance

- The Reasoning of Tests of Significance

- Stating Hypotheses

- Test Statistics

- *P*-values

- Statistical Significance

- Test for a Population Mean

- Two-Sided Significance Tests and Confidence Intervals

# Statistical Inference

Confidence intervals are one of the two most common types of statistical inference. Use a confidence interval when your goal is to estimate a population parameter. The second common type of inference, called *tests of significance*, has a different goal: to assess the evidence provided by data about some claim concerning a population.

A **test of significance** 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.
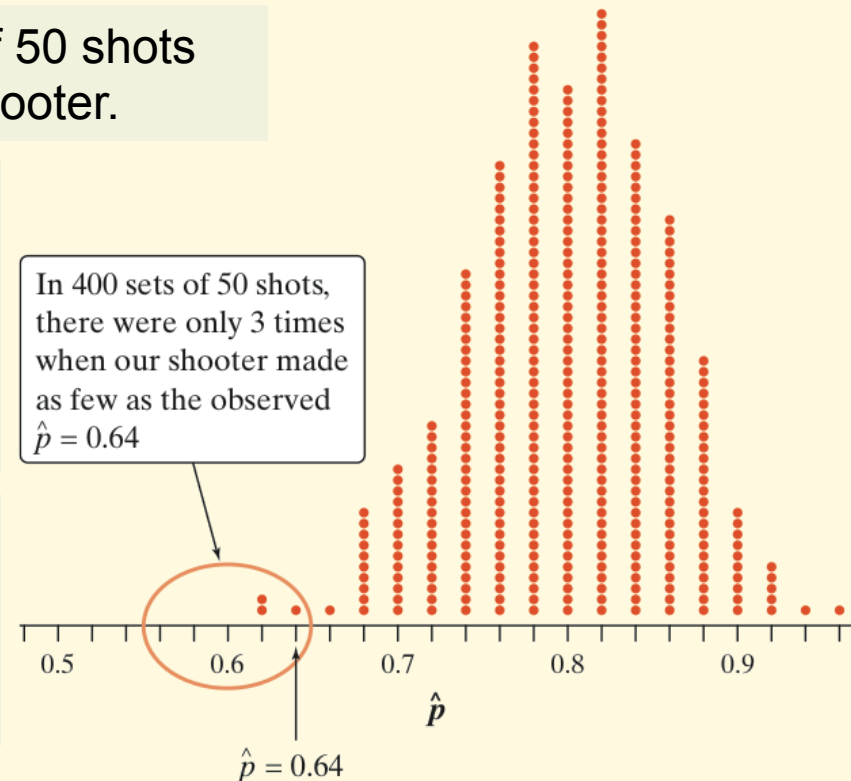
# The Reasoning of Tests of Significance

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.

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

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.

0.5    0.6    0.7    0.8    0.9

$\hat{p}$

$\hat{p} = 0.64$

# Stating Hypotheses

A significance test starts with a careful statement of the claims we want to compare.

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 effect" or "no difference in the true means."

The claim about the population that we are trying to find evidence for is the **alternative hypothesis ($H_a$).** The alternative is **one-sided** if it states that a parameter is *larger* or *smaller than* the null hypothesis value. It is **two-sided** if it states that the parameter is *different from* the null value (it could be either smaller or larger).

In the free-throw shooter example, our hypotheses are:

$$H_0: p = 0.80$$
$$H_a: p < 0.80$$

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

# Example

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 set-ups and took the test again after two more weeks. The response variable is the difference in satisfaction scores, self-paced minus machine-paced.

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.

**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:

$$H_0: \mu = 0$$
$$H_a: \mu \neq 0$$

# Test Statistic

A test of significance is based on a statistic that estimates the parameter that appears in the hypotheses. When $H_0$ is true, we expect the estimate to take a value near the parameter value specified in $H_0$.

Values of the estimate far from the parameter value specified by $H_0$ give evidence against $H_0$.

A **test statistic** calculated from the sample data measures how far the data diverge from what we would expect if the null hypothesis $H_0$ were true.

$$z = \frac{\text{estimate - hypothesized value}}{\text{standard deviation of the estimate}}$$

Large values of the statistic show that the data are not consistent with $H_0$.

# *P*-Value

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.**

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.

# Statistical Significance

The final step in performing a significance test is to draw a conclusion about the competing claims you were testing. We will make one of two decisions based on the strength of the evidence against the null hypothesis (and in favor of the alternative hypothesis)—**reject $H_0$ or fail to reject $H_0$.**

- If our sample result is too unlikely to have happened by chance assuming $H_0$ is true, then we'll reject $H_0$.
- Otherwise, we will fail to reject $H_0$.

**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 α, the Greek letter alpha. When our *P*-value is less than the chosen αwe say that the result is **statistically significant.**

If the *P*-value is smaller than alpha, we say that the data are **statistically significant at level** $\alpha$**.** The quantity $\alpha$ is called the **significance level** or the **level of significance.**

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

*P*-value < $\alpha \to$ reject $H_0 \to$ conclude $H_a$ (in context)

*P*-value ≥ $\alpha \to$ fail to reject $H_0 \to$ cannot conclude $H_a$ (in context)

# Four Steps of Tests of Significance

**Tests of Significance: Four Steps**

1. State the null and alternative **hypotheses.**

2. Calculate the value of the **test statistic.**

3. Find the ***P*-value** for the observed data.

4. State a **conclusion.**

We will learn the details of many tests of significance in the following chapters. The proper test statistic is determined by the hypotheses and the data collection design.

① Is there a increase in avg SAT score for students taking a prep class? For students w/o the class, the average score is 980 w/ std dev 100.

Let $\mu$ = mean SAT score of students taking prep class

$$H_0 : \mu = 980$$

$$H_A : \mu > 980$$

If $\bar{x} = 1050$, what is our conclusion?

for $n = 100$

# Tests for a Population Mean

## z TEST FOR A POPULATION MEAN

Draw an SRS of size $n$ from a Normal population that has unknown mean $\mu$ and known standard deviation $\sigma$. To test the null hypothesis that $\mu$ has a specified value,
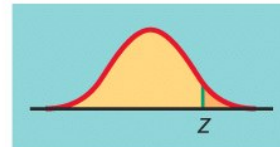
$$H_0: \mu = \mu_0$$
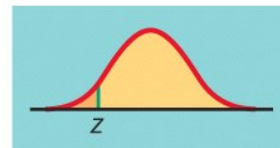
calculate the **one-sample $z$ statistic**

$$z = \frac{\overline{x} - \mu_0}{\sigma/\sqrt{n}}$$

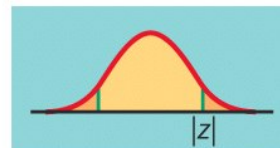In terms of a variable $Z$ having the standard Normal distribution, the $P$-value for a test of $H_0$ against

$H_a: \mu > \mu_0$ is $P(Z \geq z)$

$H_a: \mu < \mu_0$ is $P(Z \leq z)$

$H_a: \mu \neq \mu_0$ is $2P(Z \geq |z|)$

$H_0: \mu = 980$

$H_A: \mu > 980$

$\sigma = 100$

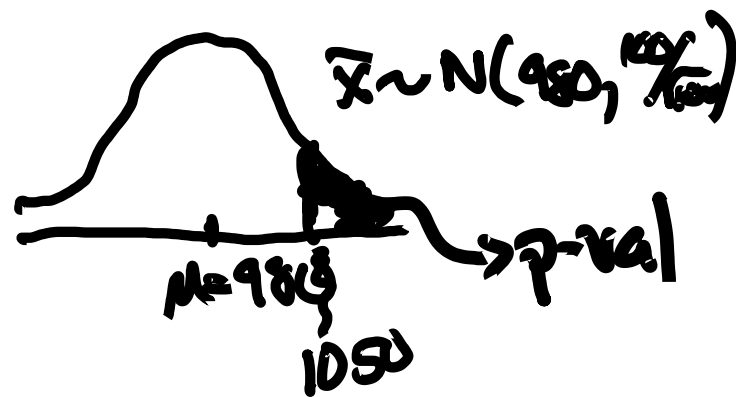$\overline{X} = 1050 \, , \, n = 100$



$\overline{X} \sim N(980, \frac{100^2}{100})$

$\mu = 980$

$1050$

$\to$ p-val

$z = \dfrac{\overline{X} - 980}{100/\sqrt{100}}$

$= \dfrac{1050 - 980}{10}$

$= 70/10$

$= 7$

P-val $\approx 0$

$\Rightarrow$ Reject $H_0$ & conclude
that the class leads to an
increase in mean SAT score

Is there a difference between the air speeds of African & European sparrows?

Let: $\mu_A$ = mean speed for African sparrows

$\mu_E$ = "  "  "  "  Euro  "

$$H_0: \mu_A - \mu_E = 0$$

$$H_A: \mu_A - \mu_E \neq 0$$

# Example

Does the job satisfaction of assembly-line workers differ when their work is machine-paced rather than self-paced? A matched pairs study was performed on a sample of workers, and each worker's satisfaction was assessed after working in each setting. The response variable is the difference in satisfaction scores, self-paced minus machine-paced.

The null hypothesis is no average difference in scores in the population of assembly-line workers, while the alternative hypothesis (that which we want to show is likely to be true) is that there is an average difference in scores in the population of assembly workers.

$$H_0: \mu = 0 \qquad H_a: \mu \neq 0$$

This is considered a two-sided test because we are interested in determining if a difference exists (the direction of the difference is not of interest in this study).

# Example

Suppose job satisfaction scores follow a Normal distribution with standard deviation $\sigma = 60$. Data from 18 workers gave a sample mean score of 17. If the null hypothesis of no average difference in job satisfaction is true, the test statistic would be:

$$z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} = \frac{17 - 0}{60 / \sqrt{18}} \approx 1.20$$
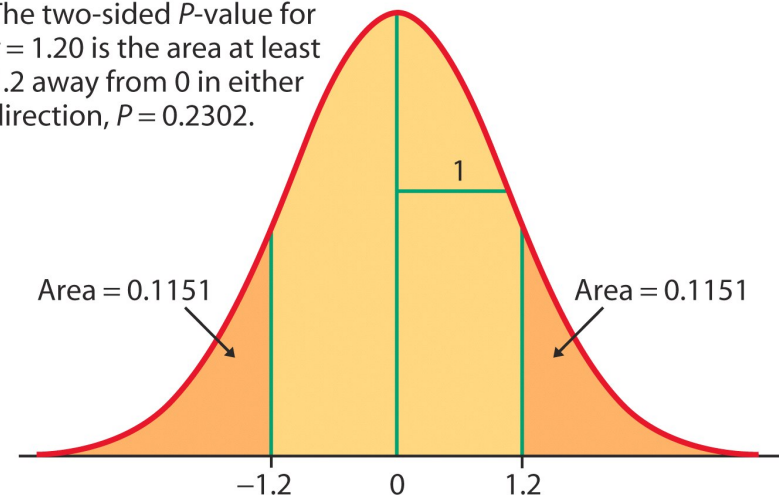
# Example

For the test statistic $z = 1.20$ and alternative hypothesis $H_a$: $\mu \neq 0$, the $P$-value would be:

$P$-value = $P(Z < -1.20$ or $Z > 1.20)$
      = $2\ P(Z < -1.20) = 2\ P(Z > 1.20)$
      = $(2)(0.1151) = 0.2302$

The two-sided $P$-value for $z = 1.20$ is the area at least 1.2 away from 0 in either direction, $P = 0.2302$.

Area = 0.1151

Area = 0.1151

−1.2    0    1.2

If $H_0$ is true, there is a 0.2302 (23.02%) chance that we would see results at least as extreme as those in the sample; thus, because we saw results that are likely if $H_0$ is true, we therefore do not have good evidence against $H_0$ and in favor of $H_a$.
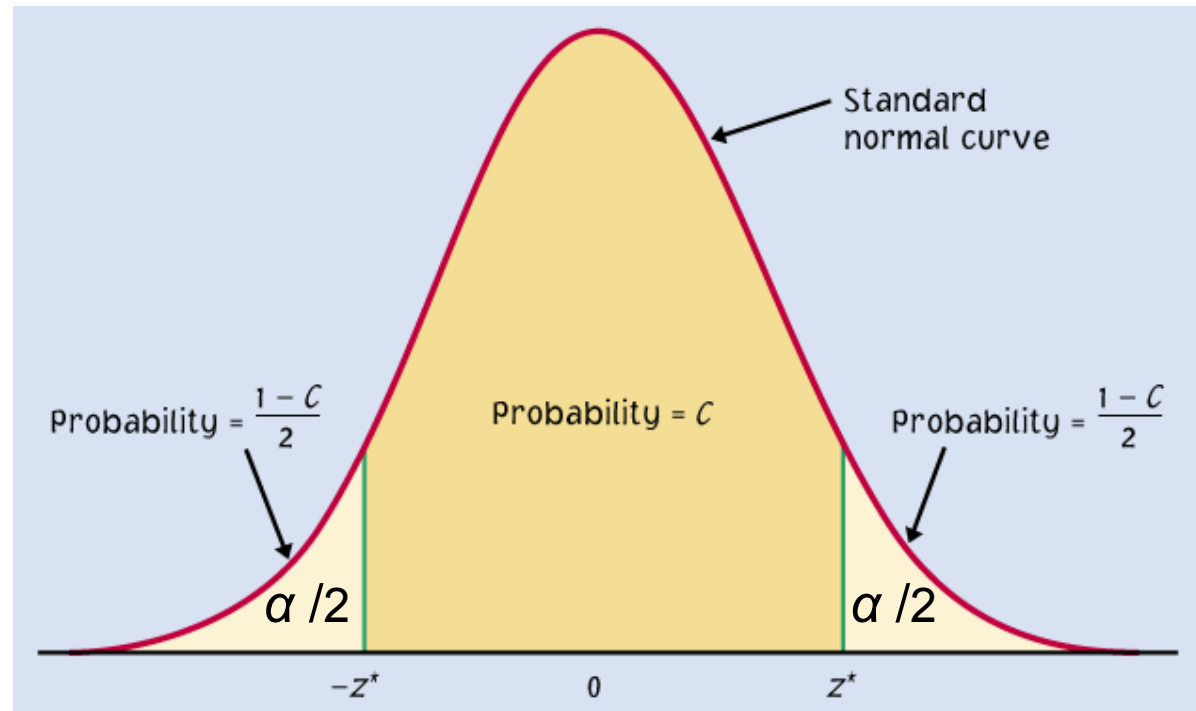
# Two-Sided Significance Tests and Confidence Intervals

Because a two-sided test is symmetrical, you can also use a $1 - \alpha$ confidence interval to test a two-sided hypothesis at level $\alpha$.

In a two-sided test,

$C = 1 - \alpha$

$C$ confidence level

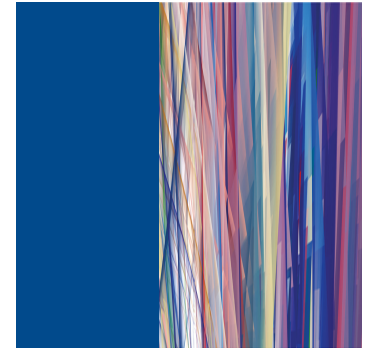$\alpha$ significance level

# *P*-Values versus Fixed $\alpha$

Statistics in practice uses technology to get *P*-values quickly and accurately. In the absence of suitable technology, you can get approximate *P*-values by comparing your test statistic with **critical values** from a table.

To find the approximate *P*-value for any *z* statistic, compare *z* (ignoring its sign) with the **critical values *z*\*** at the bottom of Table C. If *z* falls between two values of *z*\*, the *P*-value falls between the two corresponding values of *P* in the "One-sided *P*" or the "Two-sided *P*" row of Table C.

A confidence interval gives a black and white answer: Reject or don't reject $H_0$. But it also estimates a range of likely values for the true population mean *μ*.

A *P*-value quantifies how strong the evidence is against the $H_0$. But if you reject $H_0$, it doesn't provide any information about the true population mean *μ*.

# 6.3 Use and Abuse of Tests

- Choosing a Significance Level

- What Statistical Significance Does Not Mean

- Don't Ignore Lack of Significance

- Beware of Searching for Significance

# Cautions About Significance Tests

**Choosing the significance level $\alpha$**

Factors often considered:

- What are the consequences of rejecting the null hypothesis (e.g., global warming, convicting a person for life with DNA evidence)?
- Are you conducting a preliminary study? If so, you may want a larger $\alpha$ so that you will be less likely to miss an interesting result.

Some conventions:

- We typically use the standards of our field of work.

- There are no "sharp" cutoffs: for example, 4.9% versus 5.1%.

- It is the order of magnitude of the *P*-value that matters: "somewhat significant," "significant," or "very significant."

# Cautions About Significance Tests

**What statistical significance does not mean**

Statistical significance only says whether the effect observed is likely to be due to chance alone because of random sampling.

Statistical significance may not be practically important. That's because statistical significance doesn't tell you about the **magnitude** of the effect, only that there is one.

An effect could be too small to be relevant. And with a large enough sample size, significance can be reached even for the tiniest effect.

- A drug to lower temperature is found to reproducibly lower patient temperature by 0.4°Celsius (*P*-value < 0.01). But clinical benefits of temperature reduction only appear for a 1° decrease or larger.

# Cautions About Significance Tests

## Don't ignore lack of significance

- Consider this provocative title from the *British Medical Journal*: "Absence of evidence is not evidence of absence."
- Having no proof of who committed a murder does not imply that the murder was not committed.

Indeed, failing to find statistical significance in results is not rejecting the null hypothesis. This is very different from actually accepting it. The sample size, for instance, could be too small to overcome large variability in the population.

When comparing two populations, lack of significance does not imply that the two samples come from the same population. They could represent two very distinct populations with similar mathematical properties.
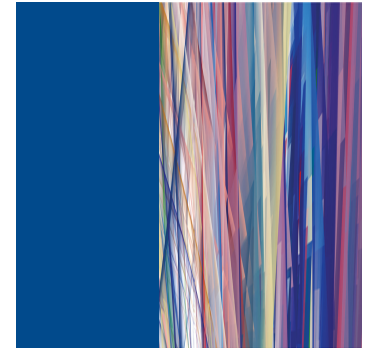
# Cautions About Significance Tests

**Beware of searching for significance**

There is no consensus on how big an effect has to be in order to be considered meaningful. In some cases, effects that may appear to be trivial can be very important.

- **Example:** Improving the format of a computerized test reduces the average response time by about 2 seconds. Although this effect is small, it is important because this is done millions of times a year. The *cumulative* time savings of using the better format is gigantic.

Always think about the context. Try to plot your results, and compare them with a baseline or results from similar studies.

# 6.4 Power and Inference as a Decision

- Power

- Increasing the Power

- Inference as a Decision

- Error Probabilities

# Power

When we draw a conclusion from a significance test, we hope our conclusion will be correct. But sometimes it will be wrong. There are two types of mistakes we can make.

If we reject $H_0$ when $H_0$ is true, we have committed a **Type I error**.

If we fail to reject $H_0$ when $H_0$ is false, we have committed a **Type II error**.

| | | Truth about the population | |
|---|---|---|---|
| | | $H_0$ true | $H_0$ false ($H_a$ true) |
| **Conclusion based on sample** | Reject $H_0$ | *False Positive* **Type I error** | *Correct conclusion* |
| | Fail to reject $H_0$ | *Correct conclusion* | *False Negative* **Type II error** |

# Power

The probability of a Type I error is the probability of rejecting $H_0$ when it is really true. This is exactly the significance level of the test.

> The significance level $\alpha$ of any fixed-level test is the probability of a Type I error. That is, $\alpha$ is the probability that the test will reject the null hypothesis $H_0$ when $H_0$ is in fact true. Consider the consequences of a Type I error before choosing a significance level.

A significance test makes a Type II error when it fails to reject a null hypothesis that really is false. There are many values of the parameter that satisfy the alternative hypothesis, so we concentrate on one value. We can calculate the probability that a test *does* reject $H_0$ when an alternative is true. This probability is called the **power** of the test against that specific alternative.

> The **power** of a test against a specific alternative is the probability that the test will reject $H_0$ at a chosen significance level $\alpha$ when the specified alternative value of the parameter is true.

46

# Power

A potato-chip producer wonders whether the significance test of $H_0$: $p = 0.08$ versus $H_a$: $p > 0.08$ based on a random sample of 500 potatoes has enough power to detect a shipment with, say, 11% blemished potatoes.

**What if $p = 0.11$?**

If $H_0$ is false, a decision to fail to reject $H_0$ based on the data is a Type II error.

Sampling distribution of $\hat{p}$ if $H_0$ is false and $p = 0.11$ is true

$P$(Type II error)
$= 1 - 0.7517$
$= 0.2483$

The power of the test to detect that $p = 0.11$

0.7517

|  | 0.0680 | 0.0820 | 0.0960 | 0.1100 | 0.1240 | 0.1380 | 0.1520 |

**Values of $\hat{p}$**

$\hat{p} = 0.0999$

**Power and Type II Error**

The power of a test against any alternative is 1 minus the probability of a Type II error for that alternative; that is, power = $1 - \beta$.

ld reject
= 0.05 if
ple
a

on to the
he
e.

Since we reject $H_0$ at $\alpha = 0.05$ if our sample yields a proportion > 0.0999, we'd correctly reject the shipment about 75% of the time.

# Power

How large a sample should we take when we plan to carry out a significance test? The answer depends on what alternative values of the parameter are important to detect.

Summary of influences on the question "How many observations do I need?"

•If you insist on a smaller significance level (such as 1% rather than 5%), you have to take a larger sample. A smaller significance level requires stronger evidence to reject the null hypothesis.

•If you insist on higher power (such as 99% rather than 90%), you will need a larger sample. Higher power gives a better chance of detecting a difference when it is really there.

•     At any significance level and desired power, detecting a small     difference requires a larger sample than detecting a large     difference.
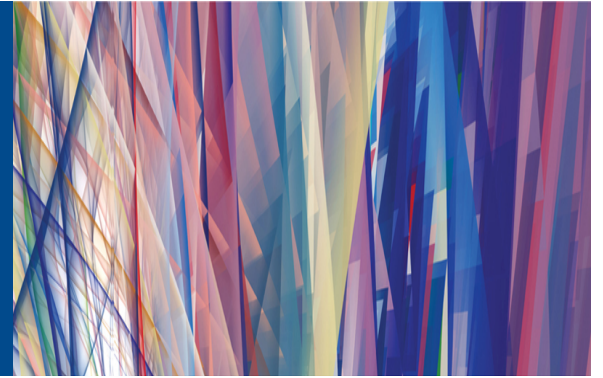
# The Common Practice of Testing Hypotheses

1. State $H_0$ and $H_a$ as in a test of significance.

2. Think of the problem as a decision problem, so the probabilities of Type I and Type II errors are relevant.

3. Consider only tests in which the probability of a Type I error is no greater than $\alpha$.

4. Among these tests, select a test that makes the probability of a Type II error as small as possible.

# Chapter 6
# Introduction to Inference

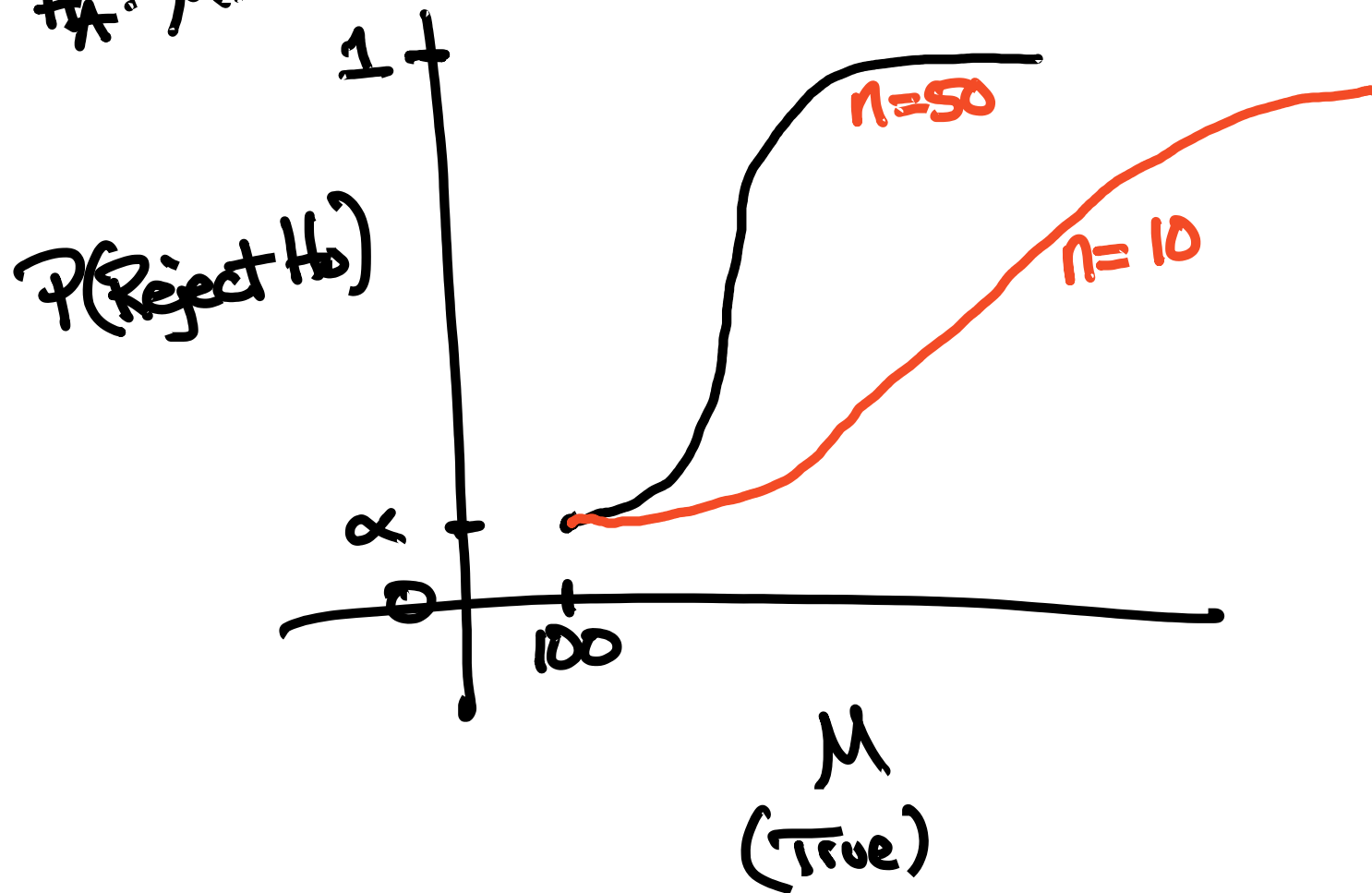**6.1 Estimating with Confidence**

**6.2 Tests of Significance**

**6.3 Use and Abuse of Tests**

**6.4 Power and Inference as a Decision**
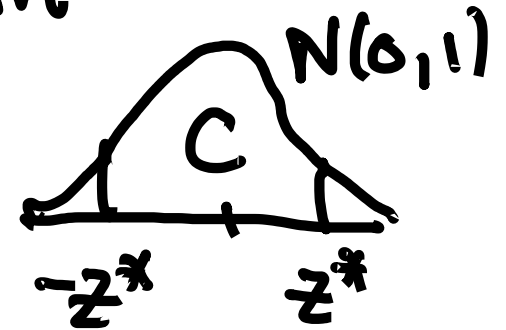
$H_0: \mu = 100$
$H_A: \mu > 100$

Power Curve



1

$P(\text{Reject } H_0)$

n=50

n=10

$\alpha$

0

100

$\mu$
(True)

$$\mu_{\hat{p}} = p$$

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

$$\hat{p} \sim N\left(\mu_{\hat{p}}, \sigma_{\hat{p}}\right)$$

CI for $\mu$:  $\bar{X} \pm z^* \, \sigma/\sqrt{n}$


$N(0,1)$
$C$
$-z^*$  $z^*$

Testing $H_0: \mu = \mu_0$

$$z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$


$\rightarrow$ p-val
$z$