

# **CHAPTER 8**

## **ESTIMATION**

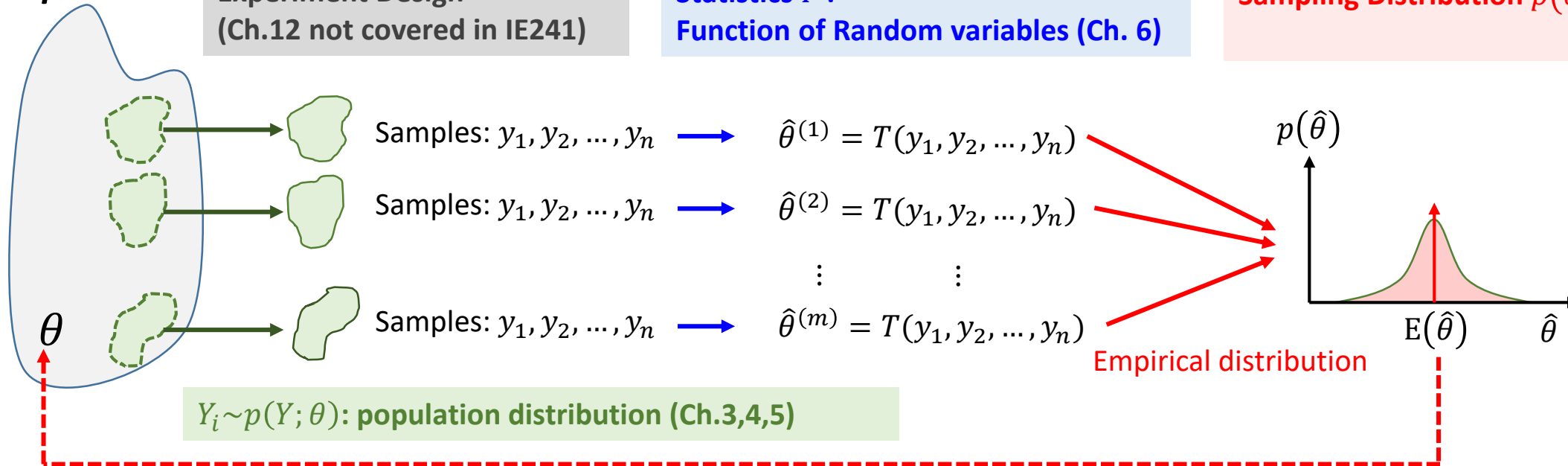
# Road Map on IE241 I

## Population

Experiment Design  
(Ch.12 not covered in IE241)

Statistics  $T$  :  
Function of Random variables (Ch. 6)

Sampling Distribution  $p(\hat{\theta})$  : (Ch. 7)



**Parameter Inference:**  
(with goodness measures)

- ✓ Estimation:  $E(\hat{\theta}) = \theta$ ? (Ch.8 & Ch.9)
- ✓ Hypothesis Testing:  $\hat{\theta} = \theta$ ? or  $\hat{\theta} > \theta$ ? (Ch.10)

- **Probability Theory (Ch.2 ~ Ch.6)** plays an important role in inference by computing the probability of the occurrence of the sample and connects the computed probability to the most probable target parameter.
- **Estimator**  $\hat{\theta} = T(Y_1, Y_2, \dots, Y_n)$  for a target parameter  $\theta$  is a function of the random variables observed in a sample and therefore itself is a random variable.
- **Sampling distribution**  $p(\hat{\theta})$  can be used to evaluate the goodness of the **estimator** (confidence interval) and the errors (i.e.,  $\alpha$  and  $\beta$  errors) of **hypothesis testing**.

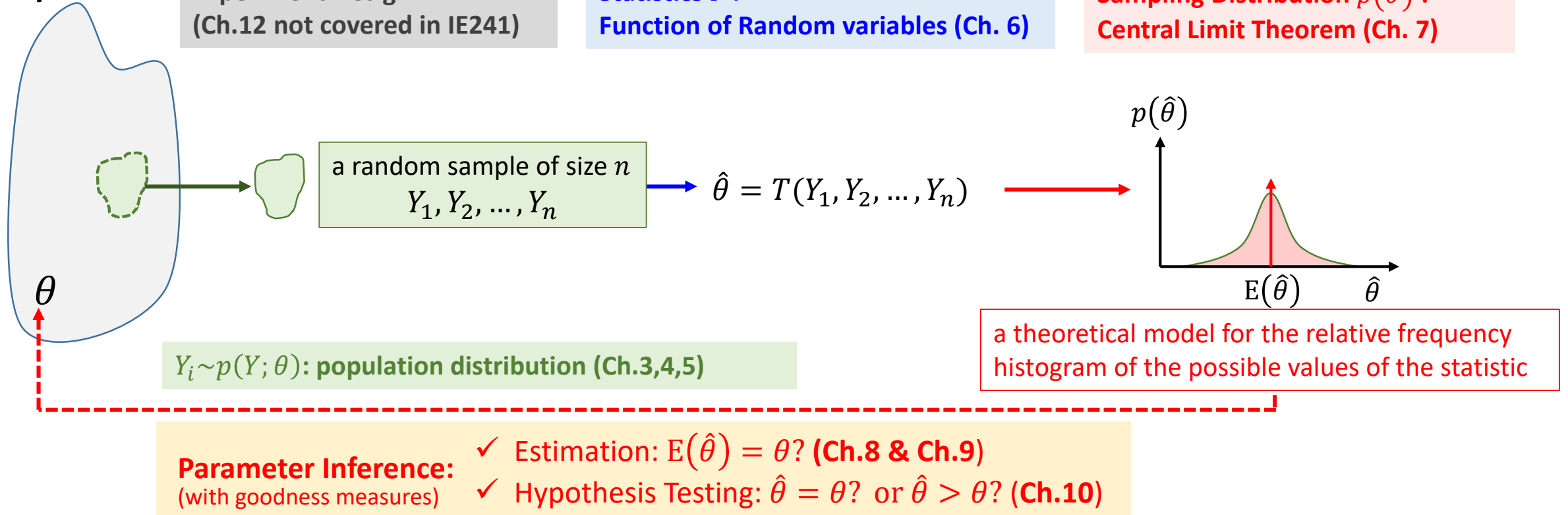
# Road Map on IE241 I

## Population

Experiment Design  
(Ch.12 not covered in IE241)

Statistics  $T$  :  
Function of Random variables (Ch. 6)

Sampling Distribution  $p(\hat{\theta})$  :  
Central Limit Theorem (Ch. 7)



- **Probability Theory (Ch.2 ~ Ch.6)** plays an important role in inference by computing the probability of the occurrence of the sample and connects the computed probability to the most probable target parameter.
- **Estimator**  $\hat{\theta} = T(Y_1, Y_2, \dots, Y_n)$  for a target parameter  $\theta$  is a function of the random variables observed in a sample and therefore itself is a random variable.
- **Sampling distribution**  $p(\hat{\theta})$  can be used to evaluate the goodness of the **estimator** (confidence interval) and the errors (i.e.,  $\alpha$  and  $\beta$  errors) of **hypothesis testing**.

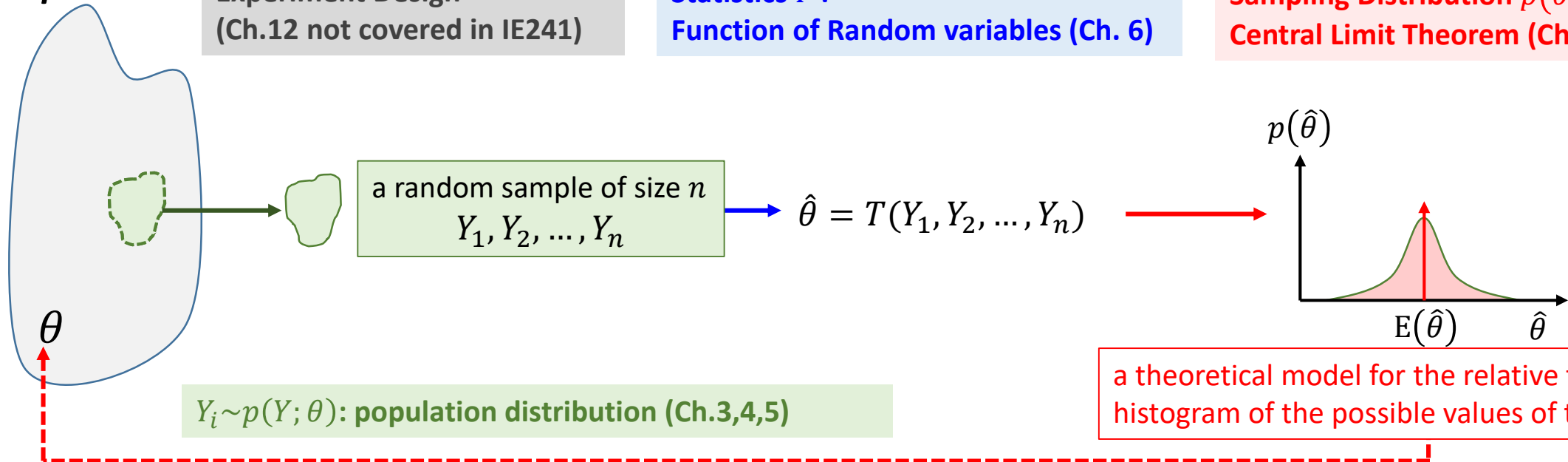
## Overview on Estimation

## Population

Experiment Design  
(Ch.12 not covered in IE241)

Statistics  $T$  :  
Function of Random variables (Ch. 6)

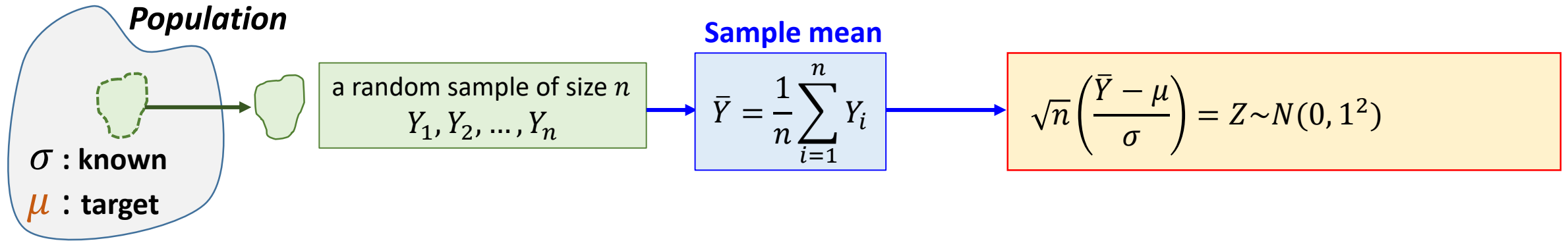
Sampling Distribution  $p(\hat{\theta})$  :  
Central Limit Theorem (Ch. 7)



- An estimator  $\hat{\theta} = E(Y_1, Y_2, \dots, Y_n)$  is a Statistics  $T(Y_1, Y_2, \dots, Y_n)$
- Type of estimator  $\hat{\theta} = E(Y_1, Y_2, \dots, Y_n)$ 
  - Point estimator vs. Interval estimator
- Goodness of estimator
  - Bias vs. variance
- Impact of the size of sample
  - Large sample size vs. Small sample size

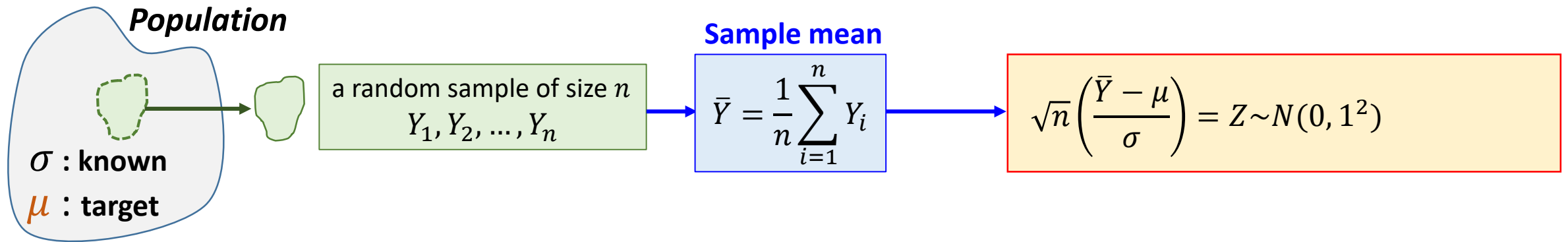
# Motivation

- An estimator  $\hat{\theta} = E(Y_1, Y_2, \dots, Y_n)$  is a Statistics  $T(Y_1, Y_2, \dots, Y_n)$
- Type of estimator  $\hat{\theta} = E(Y_1, Y_2, \dots, Y_n)$ 
  - Point estimator vs. Interval estimator
- Goodness of estimator
  - Bias vs. variance
- Impact of the size of sample
  - Large sample size vs. Small sample size

Estimation of the mean  $\mu$  for the population distribution

- When  $Y_i \sim N(\mu, \sigma^2)$
- When  $n$  is **small or large**
- When  $\sigma$  is assumed to be **known**

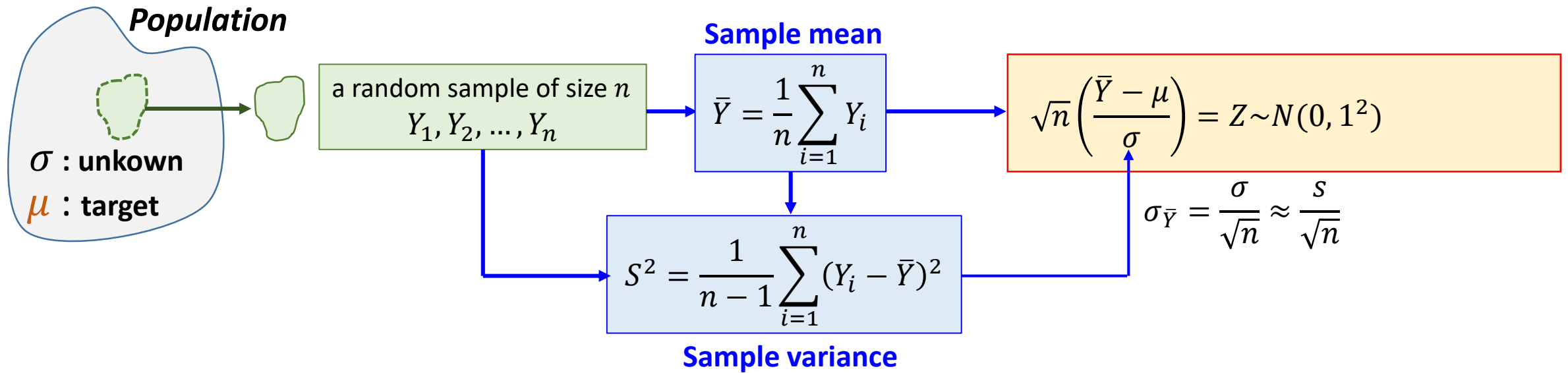
$$\sqrt{n} \left( \frac{\bar{Y} - \mu}{\sigma} \right) = Z \sim N(0, 1^2)$$

Estimation of the mean  $\mu$  for the population distribution

- When  $Y_i \sim$  Distribution with  $E(Y_i) = \mu$  and  $V(Y_i) = \sigma^2$
- When  $n$  is **large**
- When  $\sigma$  is assumed to be **known**

$$\sqrt{n} \left( \frac{\bar{Y} - \mu}{\sigma} \right) \sim Z \sim N(0, 1^2) \quad \text{Due to Central limit theorem}$$

# Estimation of the mean $\mu$ for the population distribution



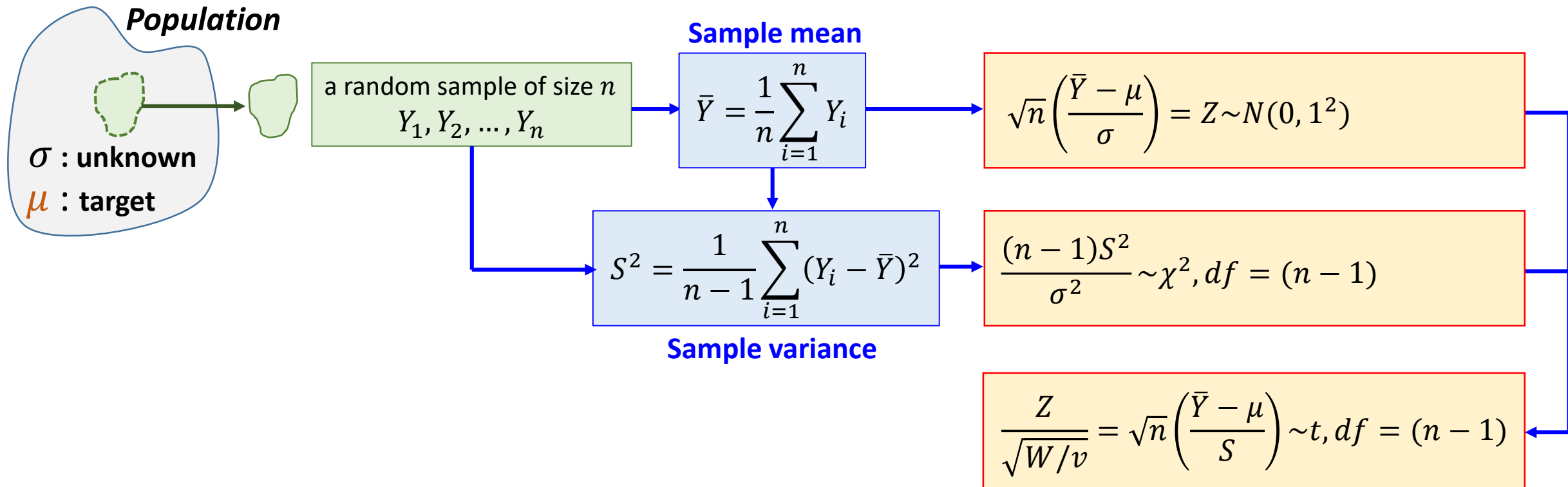
- When  $Y_i \sim$  Distribution with  $E(Y_i) = \mu$  and  $V(Y_i) = \sigma^2$
- When  $n$  is **large**
- When  $\sigma$  is assumed to be **unknown** ( $s$  is used instead of  $\sigma$ )

$$\left( \frac{\bar{Y} - \mu}{\sigma_{\bar{Y}}} \right) = \left( \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \right) \approx \left( \frac{\bar{Y} - \mu}{s/\sqrt{n}} \right) \sim Z \sim N(0, 1^2)$$

Due to Central limit theorem



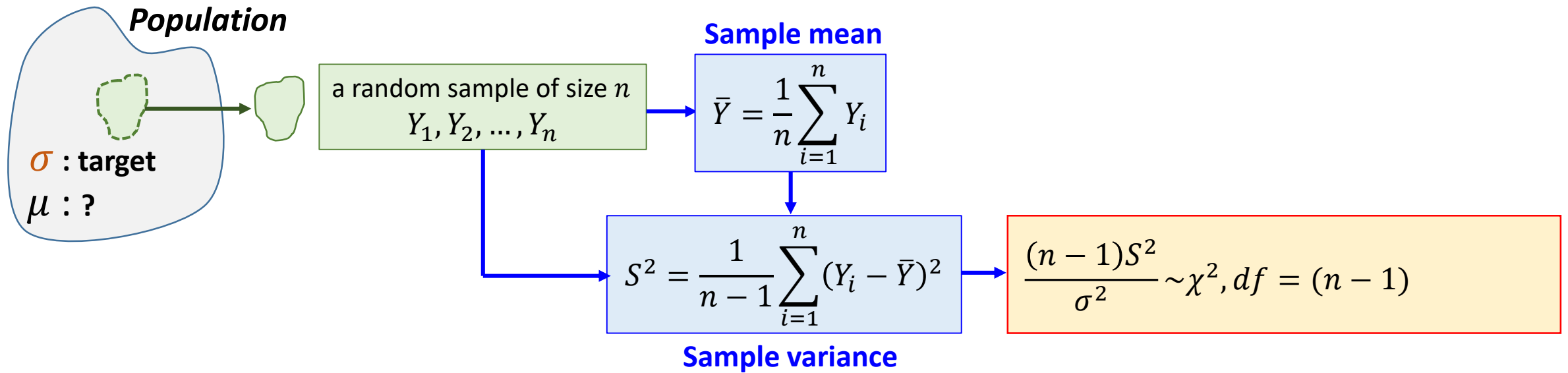
# Estimation of the mean $\mu$ for the population distribution



- When  $Y_i \sim N(\mu, \sigma^2)$
- When  $n$  is **small**
- When  $\sigma$  is assumed to be **unknown**

$$\sqrt{n} \left( \frac{\bar{Y} - \mu}{S} \right) \sim t, df = (n-1)$$

# Estimation of the variance $\sigma^2$ for the population distribution



- When  $Y_i \sim N(\mu, \sigma^2)$
- When  $n$  is **small or large**

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2, df = (n-1)$$

# Practical Uses of Estimation

- Estimation has many practical applications. For example,
  - ✓ estimating the proportion  $p$  of washers that can be expected to fail prior to the expiration of a 1-year guarantee time.
  - ✓ Other important population parameters are the population mean, variance, and standard deviation. For example,
    - the mean waiting time  $\mu$  at a supermarket checkout station or
    - the standard deviation of the error of measurement  $\sigma$  of an electronic instrument.

### Types of an Estimation

- Suppose that we wish to estimate the average amount of mercury  $\mu$  that a newly developed process can remove from 1 ounce of ore obtained at a geographic location
  - ✓ A point estimate : 13 ounce-that we think is close to the unknown population mean  $\mu$
  - ✓ An interval estimate : (.07, .19) that is intended to enclose the parameter of interests
- The information in the sample can be used to calculate the value of a point estimate, an interval estimate, or both.
- In any case, the actual estimation is accomplished by using an estimator for the target parameter.

# Definition of Estimator

### DEFINITION 8.1

An *estimator* is a rule, often expressed as a formula, that tells how to calculate the value of an estimate based on the measurements contained in a sample.

- Many different estimators (rules for estimating) may be obtained for the same population parameter.
  - ✓ Ten engineers, each assigned to estimate the cost of a large construction job, could use different methods of estimation and thereby arrive at different estimates of the total cost.
    - Each estimator represents a unique human subjective rule for obtaining a single estimate.
    - The management of a construction firm must define good and bad as they relate to the estimation of the cost of a job.
- How can we establish criteria of **goodness** to compare statistical estimators

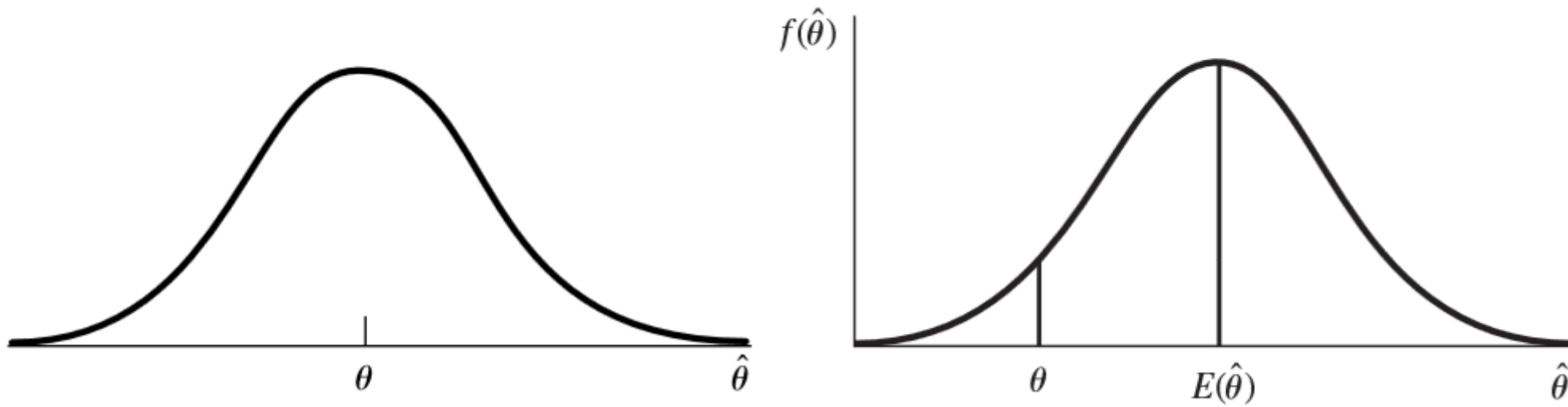
### Motivation



- Suppose that a man fires **a single shot** at a target and that shot pierces the bull's-eye.
  - we would not decide that the man is an expert Marksperson based on such a small amount of evidence.
- On the other hand, if **100 shots** in succession hit the bull's-eye,
  - we might acquire sufficient confidence in the marksperson and consider holding the target for the next shot if the compensation was adequate.
- The point is that we cannot evaluate the goodness of a point estimation procedure on the basis of the value of a single estimate;
  - rather, we must observe the results when the estimation procedure is used many, many times.

### Goodness of an Estimator

- Because the estimates are numbers, we evaluate **the goodness of the point estimator** by constructing a frequency distribution of the values of the estimates obtained **in repeated sampling** and note how closely this distribution clusters about the target parameter.
- Suppose that we wish to specify a point estimate for a population parameter that we will call  $\theta$ .
  - ✓ The estimator of  $\theta$  will be indicated by the symbol  $\hat{\theta}$ , read as “ $\theta$  hat.”



**Which one is a good estimator?**

### Unbiased Estimator

#### DEFINITION 8.2

Let  $\hat{\theta}$  be a point estimator for a parameter  $\theta$ . Then  $\hat{\theta}$  is an *unbiased estimator* if

$$E(\hat{\theta}) = \theta.$$

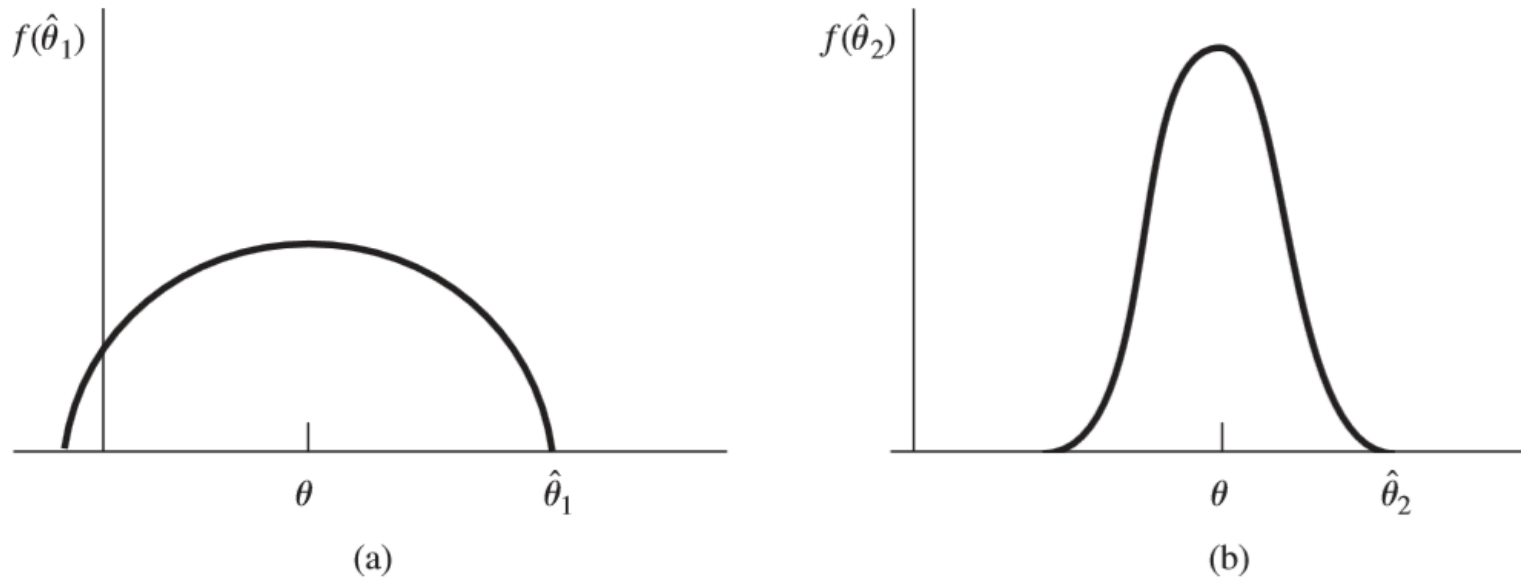
If  $E(\hat{\theta}) \neq \theta$ ,  $\hat{\theta}$  is said to be *biased*.

#### DEFINITION 8.3

The *bias* of a point estimator  $\hat{\theta}$  is given by  $B(\hat{\theta}) = E(\hat{\theta}) - \theta$ .



## Variance of Estimates



**Which one is a good estimator?**

- Figure shows two possible sampling distributions for **unbiased point estimators** for a target parameter  $\theta$ .
- We would prefer that our estimator have the type of distribution indicated in (b) because
  - ✓ the smaller variance guarantees that in repeated sampling a higher fraction of values of  $\hat{\theta}_2$  will be “close” to  $\theta$ .
- Thus, **in addition to preferring unbiasedness**, we want the variance of the distribution of the estimator  $V(\hat{\theta})$  to be as small as possible.

## Motivation

### DEFINITION 8.4

The *mean square error* of a point estimator  $\hat{\theta}$  is

$$\text{MSE}(\hat{\theta}) = E \left[ (\hat{\theta} - \theta)^2 \right].$$

- It can be shown that

$$\text{MSE}(\hat{\theta}) = V(\hat{\theta}) + [B(\hat{\theta})]^2$$

**Proof:**

$$\hat{\theta} - \theta = \hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta = \hat{\theta} - E(\hat{\theta}) + B(\hat{\theta})$$

## Motivation

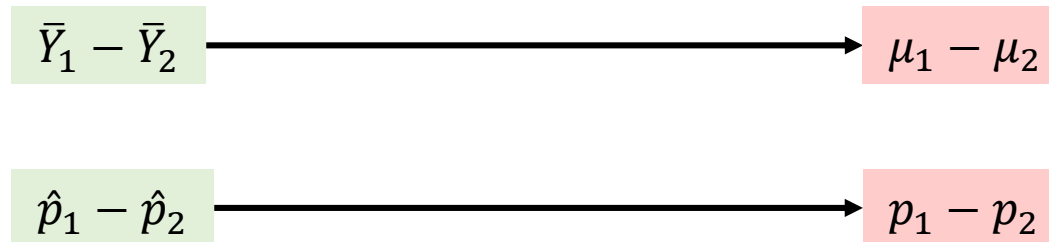
- It seems natural
  - to use the sample mean  $\bar{Y}$  to estimate the population mean  $\mu$  and
  - to use the sample proportion  $\hat{p} = Y/n$  to estimate a binomial parameter  $p$ .

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \longrightarrow \mu$$

$$\bar{p} = \frac{Y}{n} = \frac{1}{n} \sum_{i=1}^n X_i \longrightarrow p$$

## Motivation

- If an inference is to be based on independent random samples of  $n_1$  and  $n_2$  observations selected from two different populations,
  - to use the difference between sample means  $\bar{Y}_1 - \bar{Y}_2$  to estimate the difference between the population mean  $\mu_1 - \mu_2$
  - to use the difference between sample means  $\bar{p}_1 - \bar{p}_2$  to estimate the difference between the two binomial parameters  $p_1 - p_2$



## Contents

when random sampling has been employed (each sample is independent to each other)

$$E(\bar{Y}_1 - \bar{Y}_2) = E(\bar{Y}_1) - E(\bar{Y}_2) = \mu_1 - \mu_2$$

$$V(\bar{Y}_1 - \bar{Y}_2) = V(\bar{Y}_1) + V(\bar{Y}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

$$E(\hat{p}_1 - \hat{p}_2) = E(\hat{p}_1) - E(\hat{p}_2) = p_1 - p_2$$

$$V(\hat{p}_1 - \hat{p}_2) = V(\hat{p}_1) + V(\hat{p}_2) = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}$$

**Notations:**

$\sigma_{\hat{\theta}}^2$  : the variance of the sampling distribution of the estimator  $\hat{\theta}$ .

$\sigma_{\hat{\theta}} = \sqrt{\sigma_{\hat{\theta}}^2}$  : the standard deviation of the sampling distribution of the estimator  $\hat{\theta}$ .

usually called **the standard error** of the estimator  $\hat{\theta}$ .

## Common Point Estimates

Table 8.1 Expected values and standard errors of some common point estimators

Target Parameter $\theta$	Sample Size(s)	Point Estimator $\hat{\theta}$	$E(\hat{\theta})$	Standard Error $\sigma_{\hat{\theta}}$
$\mu$	$n$	$\bar{Y}$	$\mu$	$\frac{\sigma}{\sqrt{n}}$
$p$	$n$	$\hat{p} = \frac{Y}{n}$	$p$	$\sqrt{\frac{pq}{n}}$
$\mu_1 - \mu_2$	$n_1$ and $n_2$	$\bar{Y}_1 - \bar{Y}_2$	$\mu_1 - \mu_2$	$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}^{*\dagger}$
$p_1 - p_2$	$n_1$ and $n_2$	$\hat{p}_1 - \hat{p}_2$	$p_1 - p_2$	$\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}^{\dagger}$

\* $\sigma_1^2$  and  $\sigma_2^2$  are the variances of populations 1 and 2, respectively.

$\dagger$ The two samples are assumed to be independent.

- First, the expected values and standard errors for  $\bar{Y}$  and  $\bar{Y}_1 - \bar{Y}_2$  given in the table are valid regardless of the distribution of the population(s) from which the sample(s) is (are) taken.
- Second, all four estimators possess probability distributions that are approximately normal for large samples.

## Example

## EXAMPLE 8.1

Let  $Y_1, Y_2, \dots, Y_n$  be a random sample with  $E(Y_i) = \mu$  and  $V(Y_i) = \sigma^2$ .

Show that

$$S'^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

is a **biased estimator** for  $\sigma^2$  and that

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

is an **unbiased estimator** for  $\sigma^2$ .

## Example

## SOLUTION 8.1

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - \frac{1}{n} (\sum_{i=1}^n Y_i)^2 = \sum_{i=1}^n Y_i^2 - n\bar{Y}^2$$

Thus,

$$E[\sum_{i=1}^n (Y_i - \bar{Y})^2] = E[\sum_{i=1}^n Y_i^2] - nE[\bar{Y}^2] = \sum_{i=1}^n E[Y_i^2] - nE[\bar{Y}^2] = (n-1)\sigma^2$$

$$\left( \because E(Y_i^2) = V(Y_i) + E(Y_i)^2 = \sigma^2 + \mu^2, \quad E[\bar{Y}^2] = V(\bar{Y}) + E(\bar{Y})^2 = \frac{\sigma^2}{n} + \mu^2 \right)$$

It follows that

$$E(S'^2) = \frac{1}{n} E[\sum_{i=1}^n (Y_i - \bar{Y})^2] = \frac{1}{n} (n-1)\sigma^2 = \frac{n-1}{n} \sigma^2 \text{ (biased)}$$

$$E(S^2) = \frac{1}{n-1} E[\sum_{i=1}^n (Y_i - \bar{Y})^2] = \frac{1}{n-1} (n-1)\sigma^2 = \sigma^2 \text{ (unbiased)}$$



### Motivation

If we use an estimator once and acquire a single estimate, how good will this estimate be?

How much faith can we place in the validity of our inference

- One way to measure the goodness of any point estimation procedure is in terms of the distances between the estimates that it generates and the target parameter.
- This quantity, which varies randomly in repeated sampling, is called **the error of estimation**. Naturally we would like the error of estimation to be as small as possible.

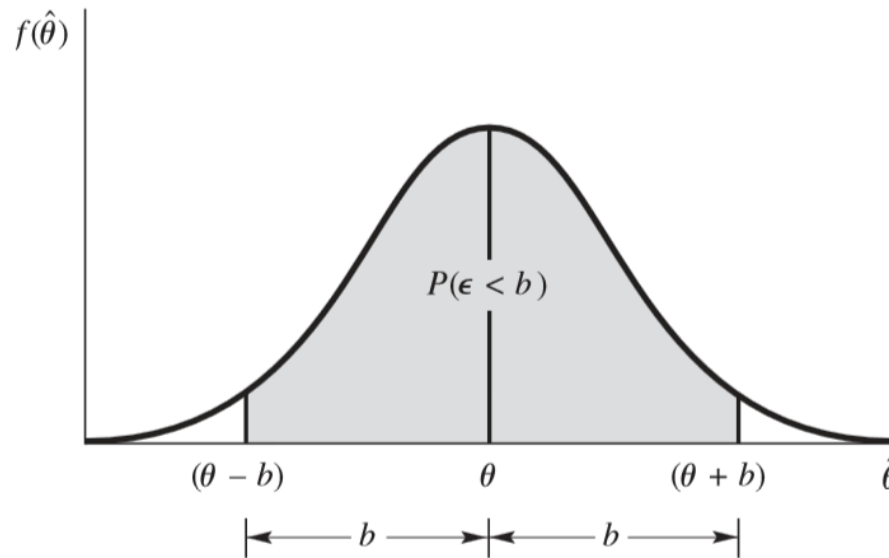
## Definition

## DEFINITION 8.5

The *error of estimation*  $\varepsilon$  is the distance between an estimator and its target parameter. That is,  $\varepsilon = |\hat{\theta} - \theta|$ .

- Because  $\hat{\theta}$  is a random variable, the error of estimation is also a random quantity, and
- We cannot say how large or small it will be for a particular estimate. However, we can make **probability statements** about it
  - ✓  $P(|\hat{\theta} - \theta| < b) = P[-b < \hat{\theta} - \theta < b] = P(\hat{\theta} - b < \theta < \hat{\theta} + b)$ 
    - We can think of  $b$  as a probabilistic bound on the error of estimation.

### Bound on Point Estimate



- Although we are not certain that a given error is less than  $b$  (because  $\epsilon$  is R.V.), the above figure indicates that  $P(\epsilon < b)$  is high.
  - If  $b$  can be regarded from a practical point of view as small, then  $P(\epsilon < b)$  provides a measure of the goodness of a single estimate.
  - This probability identifies the fraction of times, in repeated sampling, that the estimator  $\bar{\theta}$  falls within  $b$  units of  $\theta$ , the target parameter.

### Bound on Point Estimate

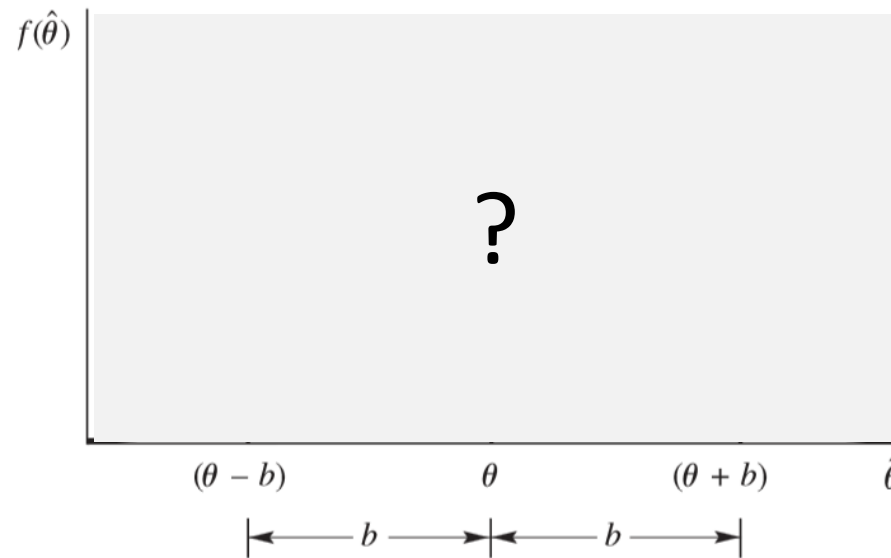
- Suppose that we want to find the value of  $b$  so that  $P(\varepsilon < b) = .90$ .
- If we know the probability density function of  $f(\hat{\theta})$  on  $\hat{\theta}$ , then we can find  $b$  such that

$$P(\varepsilon < b) = \int_{\theta-b}^{\theta+b} f(\hat{\theta}) d\hat{\theta} = 0.90$$

- Suppose that we want to find the value of  $b$  so that  $P(\varepsilon < b) = .90$ .
- If we know the probability density function of  $f(\hat{\theta})$  on  $\hat{\theta}$ , then we can find  $b$  such that

$$P(\varepsilon < b) = \int_{\theta-b}^{\theta+b} f(\hat{\theta}) d\hat{\theta} = 0.90$$

## Bound on Point Estimate



- Suppose that we want to find the value of  $b$  so that  $P(\epsilon < b) = .90$ .
- But whether we know the probability distribution of  $\hat{\theta}$  or not, if  $\hat{\theta}$  is unbiased we can find an approximate bound on  $\epsilon$  by expressing  $b$  as a multiple of the standard error of  $\hat{\theta}$

$$P(\epsilon < b) = P(\epsilon < k\sigma_{\hat{\theta}}) = P(|\hat{\theta} - \theta| < k\sigma_{\hat{\theta}}) \geq 1 - \frac{1}{k^2}$$

due to Tchebysheff's Theorem

## Recall

### THEOREM 4.13 (Tchebysheff's Theorem)

Let  $Y$  be a random variable with finite mean  $\mu$  and variance  $\sigma^2$ . Then, for any  $k > 0$ ,

$$P(|Y - \mu| < k\sigma) \geq 1 - \frac{1}{k^2} \quad \text{or} \quad P(|Y - \mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

## Bound on Point Estimate

Table 8.2 Probability that  $(\mu - 2\sigma) < Y < (\mu + 2\sigma)$ 

Distribution	Probability
Normal	.9544
Uniform	1.0000
Exponential	.9502

- The point is that  $b = 2\sigma_{\hat{\theta}}$  is a good approximate bound on the error of estimation in most practical situations.
- According to Tchebysheff's theorem, the probability that the error of estimation will be less than this bound  $(\theta \pm 2\sigma_{\hat{\theta}})$  is at least .75.
  - ✓ the bounds for probabilities provided by Tchebysheff's theorem are usually very conservative;
  - ✓ the actual probabilities usually **exceed** the Tchebysheff bounds by a considerable amount.

### Example

#### EXAMPLE 8.2

A sample of  $n = 1000$  voters, randomly selected from a city, showed  $y = 560$  in favor of candidate Jones. Estimate  $p$ , the fraction of voters in the population favoring Jones, and place a 2-standard-error bound on the error of estimation.



## Example

## SOLUTION 8.2

$$\hat{p} = \frac{y}{n} = \frac{560}{1000} = 0.56$$

The probability distribution of  $\hat{p}$  is very accurately approximated by a normal probability distribution for large samples. Since  $n = 1000$ , when  $b = 2\sigma_{\hat{p}}$ , the probability that  $\varepsilon$  will be less than  $b$  is approximately 0.95. From Table 8.1, the standard error of the estimator for  $p$  is given by  $\sigma_{\hat{p}} = \sqrt{\frac{pq}{n}}$ . Therefore,  $b = 2\sigma_{\hat{p}} = 2\sqrt{pq/n}$ .

Unfortunately, to calculate  $b$ , we need to know  $p$ , and estimating  $p$  was the objective of our sampling. This apparent stalemate is not a handicap, however, because  $\sigma_{\hat{p}}$  varies little for small changes in  $p$ . Hence, substitution of the estimate  $\hat{p}$  for  $p$  produces little error in calculating the exact value of  $b = 2\sigma_{\hat{p}}$ . Then, for our example, we have

$$b = 2\sigma_{\hat{p}} = 2\sqrt{pq/n} \approx 2\sqrt{\frac{(0.56)(0.44)}{1000}} = 0.03$$

What is the significance of our calculations? The probability that the error of estimation is less than 0.03 is approximately 0.95. Consequently, we can be reasonably confident that our estimate, 0.56, is within 0.03 of the true value of  $p$ , the proportion of voters in the population who favor Jones.

## Example

### EXAMPLE 8.3

A comparison of the durability of two types of automobile tires was obtained by road testing samples of  $n_1 = n_2 = 100$  tires of each type. The number of miles until wear-out was recorded, where wear-out was defined as the number of miles until the amount of remaining tread reached a prespecified small value. The measurements for the two types of tires were obtained independently, and the following means and variances were computed:

$$\begin{aligned}\bar{y}_1 &= 26,400 \text{ miles}, & \bar{y}_2 &= 25,100 \text{ miles}, \\ s_1^2 &= 1,440,000, & s_2^2 &= 1,960,000.\end{aligned}$$

Estimate the difference in mean miles to wear-out and place a 2-standard-error bound on the error of estimation.

## Example

### SOLUTION 8.3

The point estimate of  $(\mu_1 - \mu_2)$  is  $(\bar{y}_1 - \bar{y}_2) = 26,400 - 25,100 = 1300 \text{ miles}$ ,

and the standard error of the estimator (see Table 8.1) is  $\sigma_{\bar{Y}_1 - \bar{Y}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ .

We must know  $\sigma_1^2$  and  $\sigma_2^2$ , or have good approximate values for them, to calculate  $\sigma_{\bar{Y}_1 - \bar{Y}_2}$ . Fairly accurate values of  $\sigma_1^2$  and  $\sigma_2^2$  often can be calculated from similar experimental data collected at some prior time, or they can be obtained from the current sample data by using the unbiased estimators

$$\hat{\sigma}_i^2 = S_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2, \quad i = 1, 2$$

These estimates will be adequate if the sample sizes are reasonably large.

$$\sigma_{\bar{Y}_1 - \bar{Y}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \approx \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = \sqrt{\frac{1,440,000}{100} + \frac{1,960,000}{100}} = 184.4 \text{ miles}.$$

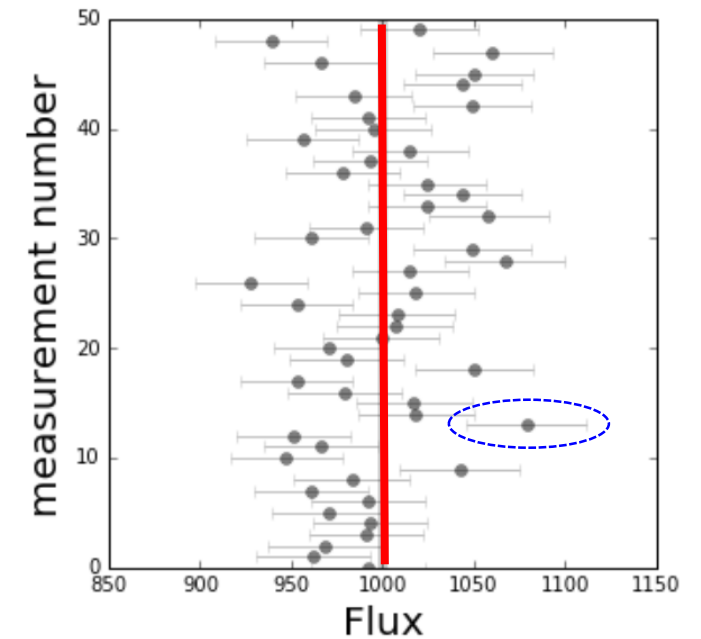
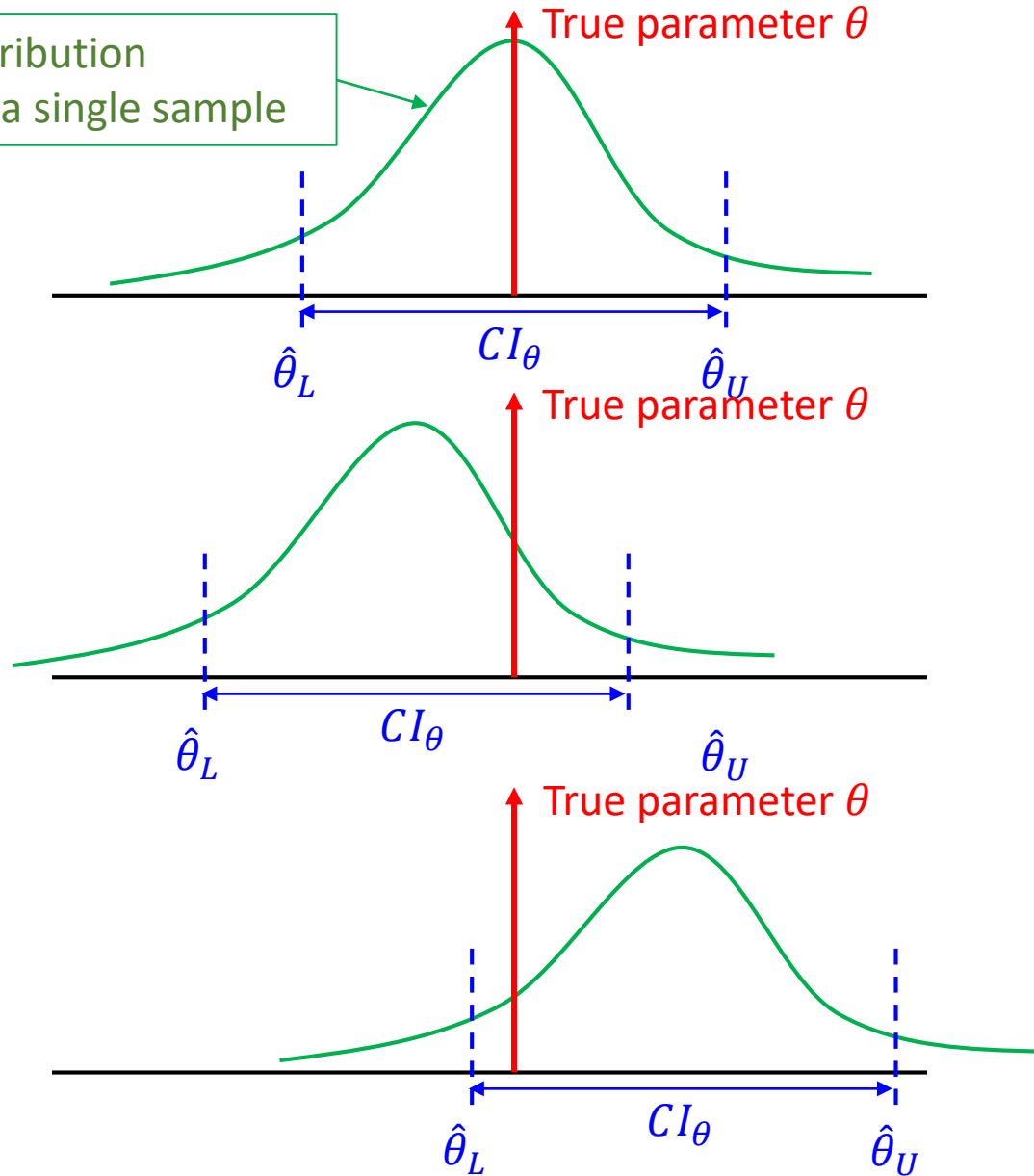
Consequently, we estimate the difference in mean wear to be 1300 miles, and we expect the error of estimation to be less than  $2\sigma_{\bar{Y}_1 - \bar{Y}_2}$ , or 368.8 miles, with a probability of approximately 0.95.

### Motivation

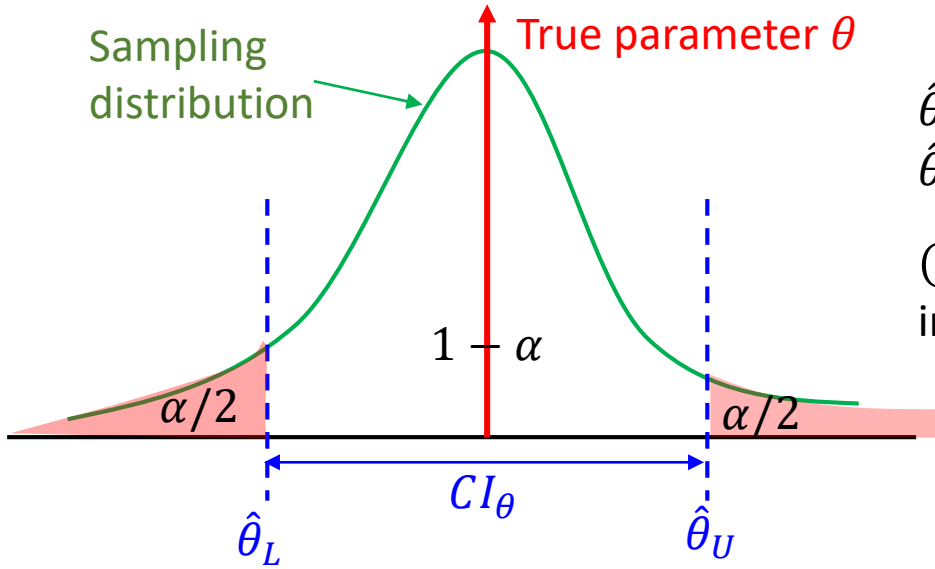
- An interval estimator is a rule specifying the method for using the sample measurements to calculate two numbers that form the endpoints of the interval.
- Ideally, the resulting interval will have two properties:
  - First, it will contain the target parameter  $\theta$ ;
  - Second, it will be relatively narrow.
- One or both of the end points of the interval, **being functions of the sample measurements (*Random variables*)**, **will vary randomly from sample to sample**.
  - Thus, the length and location of the interval are random quantities, and **we cannot be certain that the (fixed) target parameter  $\theta$  will fall between the endpoints of any single interval calculated from a single sample**.
  - This being the case, our objective is to find an interval estimator capable of generating **narrow intervals that have a high probability of enclosing  $\theta$** .

## Confidence Interval

Sampling distribution  
Estimated by a single sample



# Confidence Interval



$\hat{\theta}_L$ : *The lower confidence limit*, which is a random (function of a random samples)  
 $\hat{\theta}_U$ : *The upper confidence limit*, which is a random (function of a random samples)

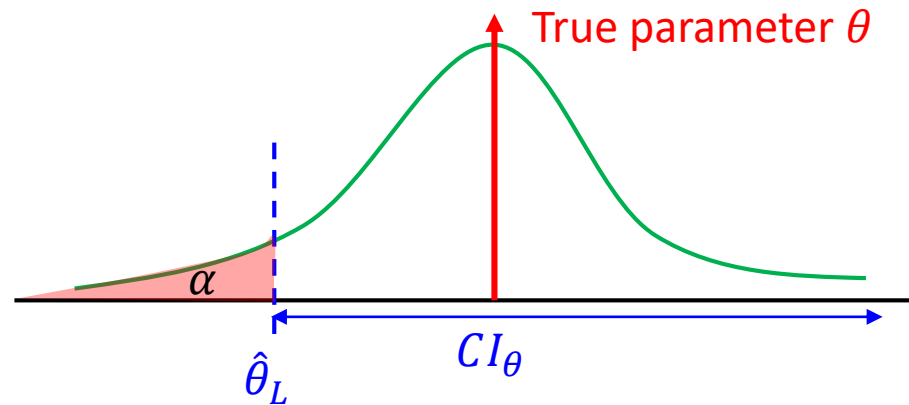
$(1 - \alpha)$  : *confidence coefficient*, the probability that a (random) confidence interval will enclose  $\theta$  (a fixed quantity) is called the confidence coefficient

$$P(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha$$

- “There is a  $(1 - \alpha)$  % probability that when I compute the confidence interval (CI) from *a current data sample*, the computed CI contains  $\theta$ 
  - From current data set, We can only say that  $\theta \in \text{CI}$  or  $\theta \notin \text{CI}$
- From a practical point of view, the confidence coefficient identifies the fraction of the time, **in repeated sampling**, that the intervals constructed will contain the target parameter  $\theta$ .
  - If the confidence coefficient is high, we can be highly confident that any confidence interval, **constructed by using the results from a single sample**, will enclose  $\theta$ .

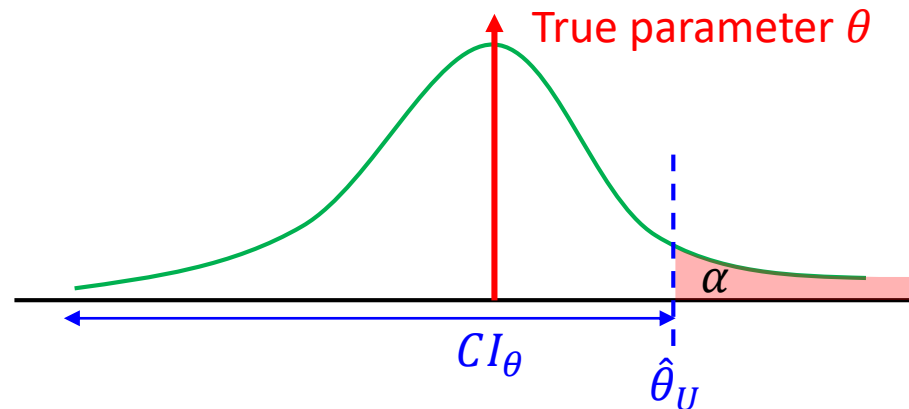
# Confidence Interval

- An lower one-sided confidence interval



$$P(\hat{\theta}_L \leq \theta) = 1 - \alpha$$

- An upper one-sided confidence interval



$$P(\theta \leq \hat{\theta}_U) = 1 - \alpha$$

## Pivotal Method for Estimating Confidence Interval

- This method depends on finding a pivotal quantity that possesses two characteristics
  - ✓ It is a function of the sample measurements and the unknown parameter  $\theta$ , where  $\theta$  is the only unknown quantity.
  - ✓ Its probability distribution does not depend on the parameter  $\theta$ .

$$P(a \leq Y \leq b) = .7$$

$$\Rightarrow P(ca \leq cY \leq cb) = .7$$

$$\Rightarrow P(a + d \leq Y + d \leq b + d) = .7$$

- That is, the probability of the event  $(a \leq Y \leq b)$  is unaffected by a change of scale or a translation of  $Y$ .
- Thus, if we know the probability distribution of a pivotal quantity, we may be able to use operations like these to form the desired interval estimator.



### Example

#### EXAMPLE 8.4

Suppose that we are to obtain a single observation  $Y$  from an exponential distribution with mean  $\theta$ . Use  $Y$  to form a confidence interval for  $\theta$  with confidence coefficient 0.90.

## Example

### SOLUTION 8.4

The probability density function for  $Y$  is given by

$$f(y) = \begin{cases} \left(\frac{1}{\theta}\right) e^{-y/\theta}, & y \geq 0 \\ 0, & \text{elsewhere} \end{cases}$$

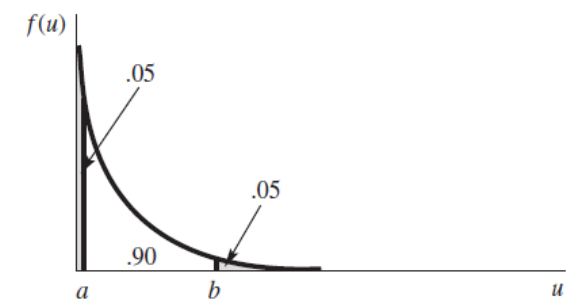
By the transformation method we can see that  $U = Y/\theta$  has the exponential density function given by

$$f_U(u) = \begin{cases} e^{-u}, & u > 0 \\ 0, & \text{elsewhere} \end{cases}$$

$U = Y/\theta$  is a function of  $Y$  (the sample measurement) and  $\theta$ , and the distribution of  $U$  does not depend on  $\theta$ . Thus, we can use  $U = Y/\theta$  as a pivotal quantity. Because we want an interval estimator with confidence coefficient equal to 0.90, we find two numbers  $a$  and  $b$  such that

$$P(a \leq U \leq b) = 0.90.$$

FIGURE 8.5  
Density function for  
 $U$ , Example 8.4



## Example

## SOLUTION 8.4 (Cont.)

One way to do this is to choose  $a$  and  $b$  to satisfy

$$p(U < a) = \int_0^a e^{-u} du = 0.05, \quad p(U > b) = \int_b^{\infty} e^{-u} du = 0.05$$

Thus,  $a = 0.051$ ,  $b = 2.996$ .

$$0.90 = P(0.051 \leq U \leq 2.996) = p(0.051 \leq \frac{Y}{\theta} \leq 2.996)$$

Because we seek an interval estimator for  $\theta$ , let us manipulate the inequalities describing the event to isolate  $\theta$  in the middle.  $Y$  has an exponential distribution, so  $p(Y > 0) = 1$ , and we maintain the direction of the inequalities if we divide through by  $Y$ . Thus, we conclude that

$$0.90 = P\left(\frac{Y}{2.996} \leq \theta \leq \frac{Y}{0.051}\right)$$

We know that limits of the form  $(Y/2.996, Y/0.051)$  will include the true (unknown) values of  $\theta$  for 90% of the values of  $Y$  we would obtain by repeatedly sampling from this exponential distribution.

## Example

**EXAMPLE 8.5**

Suppose that we take a sample of size  $n = 1$  from a uniform distribution defined on the interval  $[0, \theta]$ , where  $\theta$  is unknown. Find a 95% lower confidence bound for  $\theta$ .

## Example

## SOLUTION 8.5

Because  $Y$  is uniform on  $[0, \theta]$ , the methods of Chapter 6 can be used to show that  $U = Y/\theta$  is uniformly distributed over  $[0, 1]$ . That is,

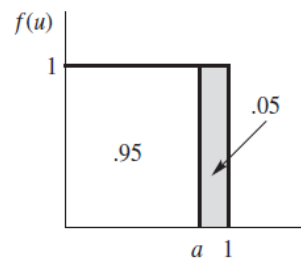
$$f_U(u) = \begin{cases} 1, & 0 \leq u \leq 1 \\ 0, & \text{elsewhere} \end{cases}$$

Again, we see that  $U$  satisfies the requirements of a pivotal quantity. Because we seek a 95% lower confidence limit

for  $\theta$ , let us determine the value for  $a$  so that  $P(U \leq a) = 0.95$ . Thus,  $a = 0.95$ .

$$p(U \leq 0.95) = p\left(\frac{Y}{\theta} \leq 0.95\right) = p\left(\frac{Y}{0.95} \leq \theta\right) = 0.95.$$

FIGURE 8.6  
Density function for  
 $U$ , Example 8.5



We see that  $Y/0.95$  is a lower confidence limit for  $\theta$ , with confidence coefficient 0.95. Because any observed  $Y$  must be less than  $\theta$ , it is intuitively reasonable to have the lower confidence limit for  $\theta$  slightly larger than the observed value of  $Y$ .

## Motivation

- First, the expected values and standard errors for  $Y$  and  $Y_1 - Y_2$  given in the table are valid **regardless of the distribution of the population(s)** from which the sample(s) is (are) taken.
- Second, **all four estimators possess probability distributions that are approximately normal for large samples.**

Table 8.1 Expected values and standard errors of some common point estimators

Target Parameter $\theta$	Sample Size(s)	Point Estimator $\hat{\theta}$	$E(\hat{\theta})$	Standard Error $\sigma_{\hat{\theta}}$
$\mu$	$n$	$\bar{Y}$	$\mu$	$\frac{\sigma}{\sqrt{n}}$
$p$	$n$	$\hat{p} = \frac{Y}{n}$	$p$	$\sqrt{\frac{pq}{n}}$
$\mu_1 - \mu_2$	$n_1$ and $n_2$	$\bar{Y}_1 - \bar{Y}_2$	$\mu_1 - \mu_2$	$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}^{*\dagger}$
$p_1 - p_2$	$n_1$ and $n_2$	$\hat{p}_1 - \hat{p}_2$	$p_1 - p_2$	$\sqrt{\frac{p_1q_1}{n_1} + \frac{p_2q_2}{n_2}}^{\dagger}$

\* $\sigma_1^2$  and  $\sigma_2^2$  are the variances of populations 1 and 2, respectively.

$\dagger$ The two samples are assumed to be independent.

### Motivation

- Second, **all four estimators possess probability distributions that are approximately normal for large samples.**
- That is, under the conditions of Section 8.3, if the target parameter  $\theta$  is,  $\mu$ ,  $p$ ,  $\mu_1 - \mu_2$ , or  $p_1 - p_2$ , then for large samples,

$$Z = \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}}$$

possesses approximately a standard normal distribution.

### Example

#### EXAMPLE 8.6

Let  $\hat{\theta}$  be a statistic that is normally distributed with mean  $\theta$  and standard error  $\sigma_{\hat{\theta}}$ . Find a confidence interval for  $\theta$  that possesses a confidence coefficient equal to  $(1 - \alpha)$ .



## Example

## SOLUTION 8.6

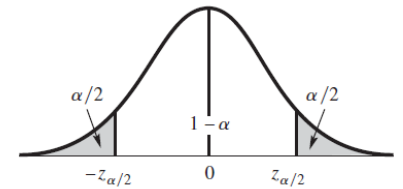
$$Z = \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}}$$

has a standard normal distribution. Now select two values in the tails of this distribution,  $z_{\alpha/2}$  and  $-z_{\alpha/2}$ , such that (see Figure 8.7)  $P(-z_{\alpha/2} \leq Z = \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}} \leq z_{\alpha/2}) = 1 - \alpha$ .

Thus,

$$P(\hat{\theta} - z_{\alpha/2}\sigma_{\hat{\theta}} \leq \theta \leq \hat{\theta} + z_{\alpha/2}\sigma_{\hat{\theta}}) = 1 - \alpha$$

FIGURE 8.7  
Location of  $z_{\alpha/2}$   
and  $-z_{\alpha/2}$



## Example

**EXAMPLE 8.7**

The shopping times of  $n = 64$  randomly selected customers at a local supermarket were recorded. The average and variance of the 64 shopping times were  $33 \text{ minutes}$  and  $256 \text{ minutes}^2$ , respectively. Estimate  $\mu$ , the true average shopping time per customer, with a confidence coefficient of  $1 - \alpha = 0.90$ .

## Example

## SOLUTION 8.7

In this case, we are interested in the parameter  $\theta = \mu$ . Thus,  $\hat{\theta} = y = 33$  and  $s^2 = 256$  for a sample of  $n = 64$  shopping times. The population variance  $\sigma^2$  is unknown, so (as in Section 8.3), we use  $s^2$  as its estimated value. The confidence interval  $\hat{\theta} \pm z_{\alpha/2} \sigma_{\hat{\theta}}$  has the form

$$\bar{y} \pm z_{\frac{\alpha}{2}} \left( \frac{\sigma}{\sqrt{n}} \right) \approx \bar{y} \pm z_{\frac{\alpha}{2}} \left( \frac{s}{\sqrt{n}} \right)$$

$$\bar{y} - z_{\frac{\alpha}{2}} \left( \frac{s}{\sqrt{n}} \right) = 33 - 1.645 \left( \frac{16}{8} \right) = 29.71,$$

$$\bar{y} + z_{\frac{\alpha}{2}} \left( \frac{s}{\sqrt{n}} \right) = 33 + 1.645 \left( \frac{16}{8} \right) = 36.29.$$

## Example

### EXAMPLE 8.8

Two brands of refrigerators, denoted A and B, are each guaranteed for 1 year. In a random sample of 50 refrigerators of brand A, 12 were observed to fail before the guarantee period ended. An independent random sample of 60 brand B refrigerators also revealed 12 failures during the guarantee period. Estimate the true difference  $(p_1 - p_2)$  between proportions of failures during the guarantee period, with confidence coefficient approximately 0.98.

## Example

## SOLUTION 8.8

The confidence interval

$$\hat{\theta} \pm z_{\alpha/2} \sigma_{\hat{\theta}}$$

has the form

$$(\hat{p}_1 - \hat{p}_2) \pm z_{\alpha/2} \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

In this case,

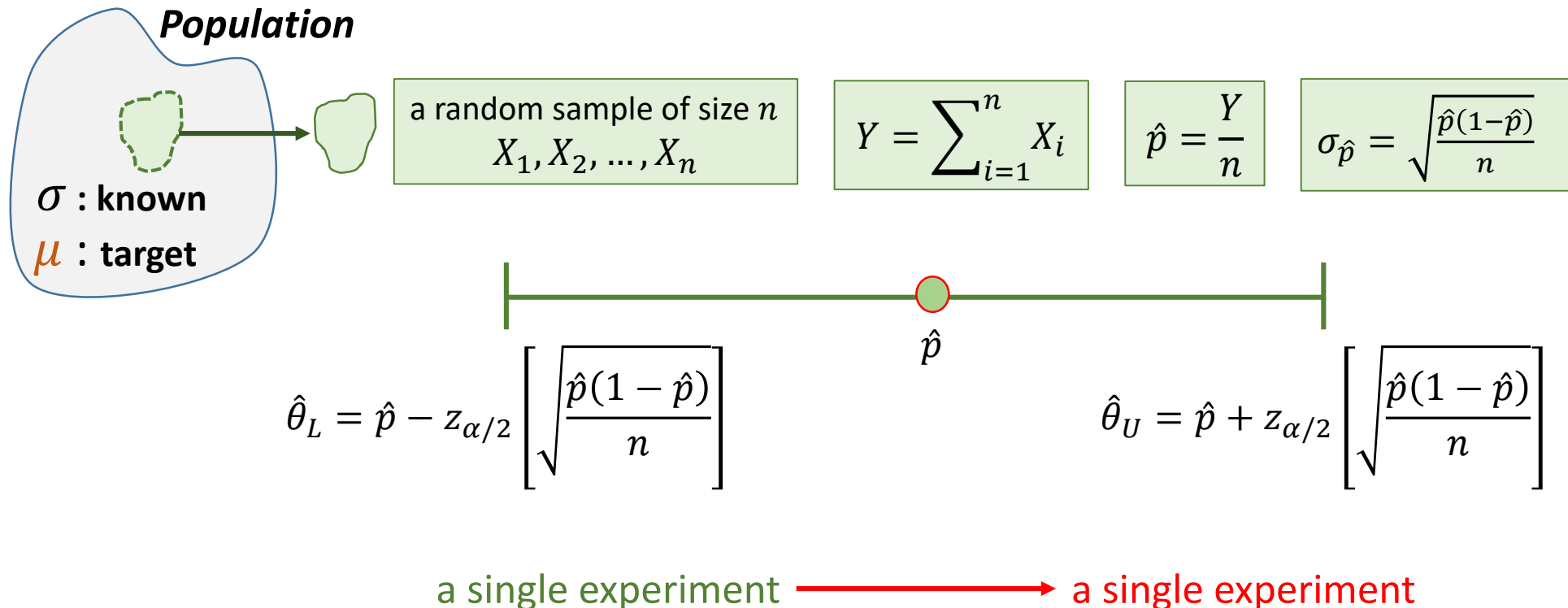
$$0.24 - 0.20 \pm 2.33 \sqrt{\frac{(0.24)(0.76)}{50} + \frac{(0.20)(0.80)}{60}} = [-0.1451, 0.2251]$$

Notice that this confidence interval contains zero. Thus, a zero value for the difference in proportions  $(p_1 - p_2)$  is “believable” (at approximately the 98% confidence level) on the basis of the observed data. However, the interval also includes the value 0.1. Thus, 0.1 represents another value of  $(p_1 - p_2)$  that is “believable” on the basis of the data that we have analyzed.

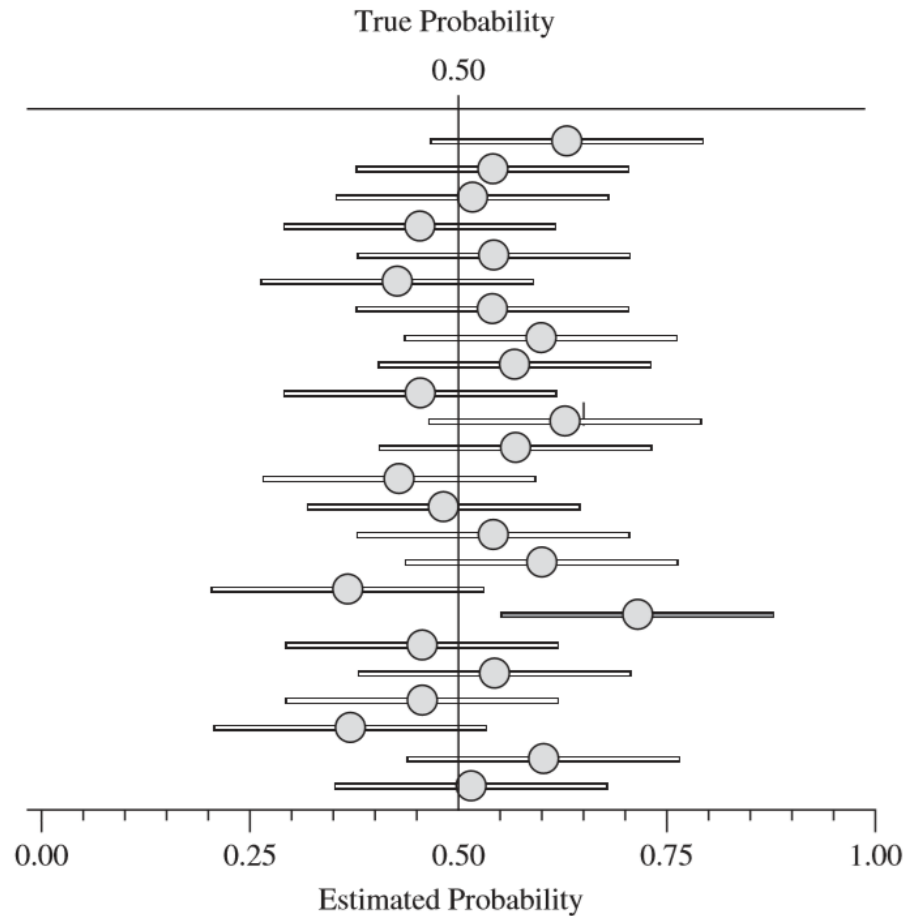
## Remarks

- Empirical investigation of the performance of the large-sample interval estimation procedure for a single population proportion  $p$ , based on  $Y$ , the number of successes observed during  $n$  trials in a binomial experiment.

- $\theta = p$
- $\hat{\theta} = \hat{p} = Y/n$
- $\sigma_{\hat{\theta}} = \sigma_{\hat{p}} = \sqrt{p(1-p)/n} = \sqrt{\hat{p}(1-\hat{p})/n} \quad \because \sqrt{\hat{p}(1-\hat{p})/n} \text{ provides a good approximation for } \sigma_{\hat{p}}$



## Remarks



- Figure shows the results of 24 independent binomial experiments, each based on 35 trials when the true value of  $p = 0.5$ .
- Notice that each individual interval either contains the true value of  $p$  or it does not.
- However, the true value of  $p$  is contained in 23 out of the 24 (95.8%) of intervals observed.

### Motivation

- The design of an experiment is essentially **a plan for purchasing a quantity of information**.
- Like any other commodity, **information may be acquired at varying prices** depending on the manner in which the data is obtained.
- Some measurements contain a large amount of information about the parameter of interest; others may contain little or none.
- Research, scientific or otherwise, is done in order to obtain information. Obviously, **we should seek to obtain information at minimum cost**.





### Motivation

- **The sampling procedure** (i.e., experimental design) affects the quantity of information per measurement.
- **The sample size  $n$**  controls the total amount of relevant information in a sample.

**At this point in our study, we will be concerned with the simplest sampling situation:  
random sampling from a relatively large population.**

- Indeed, one of the most frequent questions asked of the statistician is, How many measurements should be included in the sample?
- Unfortunately, the statistician cannot answer this question without knowing how much information the experimenter wishes to obtain.
- Referring specifically to estimation, **we would like to know how accurate the experimenter wishes the estimate to be.**
- The experimenter can indicate the desired accuracy by **specifying a bound on the error of estimation.**

### Motivation

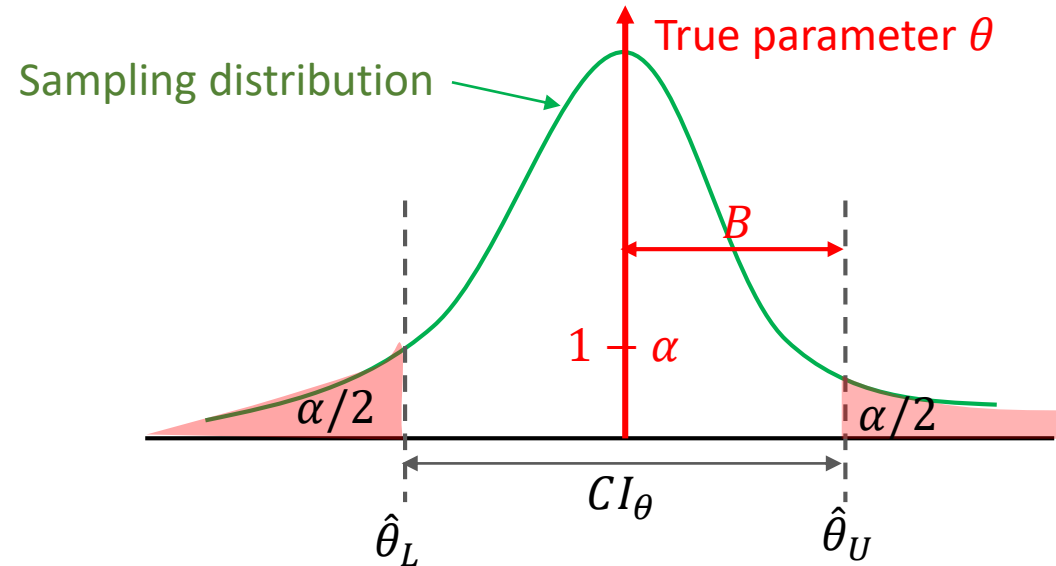
- For instance, suppose that we wish to estimate the average daily yield  $\mu$  of a chemical and we wish the error of estimation to be less than 5 tons with probability .95.
- Because approximately 95% of the sample means will lie within  $2\sigma_{\bar{Y}}$  of  $\mu$  in repeated sampling, we are asking that  $2\sigma_{\bar{Y}}$  equal 5 tons. Then

$$2\sigma_{\bar{Y}} = \frac{2\sigma}{\sqrt{(n)}} = 5 \quad \text{and} \quad n = \frac{4\sigma^2}{25}$$

- We cannot obtain an exact numerical value for  $n$  unless the population standard deviation  $\sigma$  is known.
  - This is exactly what we would expect because the variability associated with the estimator  $Y$  depends on the variability exhibited in the population from which the sample will be drawn.
- Lacking an exact value for  $\sigma$ , we use the best approximation available such as an estimate  $s$  obtained from a previous sample or knowledge of the range of the measurements in the population.
- Because the range is approximately equal to  $4\sigma$  (recall the empirical rule), one-fourth of the range provides an approximate value for  $\sigma$ . For example the range of the daily yields is 84 then  $\sigma \approx \frac{84}{4} = 21$

$$n = \frac{4\sigma^2}{25} = \frac{4(21)^2}{25} = 70.56 = 71$$

## Procedure for Selecting the Sample Size



- The experimenter **must specify a desired bound**  $B$  on the error of estimation
- The experimenter **must specify an associated confidence level**  $1 - \alpha$ .
  - For example, if the parameter is  $\theta$  and the desired bound is  $B$ , we equate

$$z_{\alpha/2} \sigma_{\hat{\theta}} = B$$

where  $P(Z > z_{\alpha/2}) = \alpha/2$

### Example

#### EXAMPLE 8.9

The reaction of an individual to a stimulus in a psychological experiment may take one of two forms, A or B. If an experimenter wishes to estimate the probability  $p$  that a person will react in manner A, how many people must be included in the experiment?

Assume that the experimenter will be satisfied if the error of estimation is less than 0.04 with probability equal to 0.90. Assume also that he expects  $p$  to lie somewhere in the neighborhood of 0.6.

## Example

## SOLUTION 8.9

Because we have specified that  $1 - \alpha = 0.90$ ,  $\alpha$  must equal 0.10 and  $\alpha/2 = 0.05$ . The  $z$  value corresponding to an area equal to 0.05 in the upper tail of the standard normal distribution is  $z_{\alpha/2} = 1.645$ . We then require that

$$1.645\sigma_{\hat{p}} = 1.645\sqrt{\frac{pq}{n}} = 0.04$$

Because the standard error of  $\hat{p}$  depends on  $p$ , which is unknown, we could use the guessed value of  $p = 0.6$  provided by the experimenter as an approximate value for  $n$ . Then  $n = 406$ .

In this example, we assumed that  $p \approx 0.60$ . How would we proceed if we had no idea about the true value of  $p$ ? In Exercise 7.76(a), we established that the *maximum* value for the variance of  $\hat{p} = Y/n$  occurs when  $p = 0.5$ . If we did not know that  $p \approx 0.6$ , we would use  $p = 0.5$ , which would yield the maximum possible value for  $n$  :  $n = 423$ .

No matter what the true value for  $p$ ,  $n = 423$  is large enough to provide an estimate that is within  $B = 0.04$  of  $p$  with probability 0.90.

### Example

#### EXAMPLE 8.10

An experimenter wishes to compare the effectiveness of two methods of training industrial employees to perform an assembly operation. The selected employees are to be divided into two groups of equal size, the first receiving training method 1 and the second receiving training method 2. After training, each employee will perform the assembly operation, and the length of assembly time will be recorded. The experimenter expects the measurements for both groups to have a range of approximately 8 minutes. If the estimate of the difference in mean assembly times is to be correct to within 1 minute with probability 0.95, how many workers must be included in each training group?

## Example

## SOLUTION 8.10

The manufacturer specified  $1 - \alpha = 0.95$ . Thus,  $\alpha = 0.05$  and  $z_{\alpha/2} = 1.96$ .

Equating  $1.96\sigma_{(\bar{Y}_1 - \bar{Y}_2)}$  to 1 minute, we obtain

$$1.96 \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = 1$$

Alternatively, because we desire  $n_1$  to equal  $n_2$ , we may let  $n_1 = n_2 = n$  and obtain the equation

$$1.96 \sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}} = 1$$

As noted earlier, the variability of each method of assembly is approximately the same; hence,  $\sigma_1^2 = \sigma_2^2 = \sigma^2$ . Because the range, 8 minutes, is approximately equal to  $4\sigma$ , we have  $4\sigma \approx 8$ , or equivalently,  $\sigma \approx 2$ .

Substituting this value for  $\sigma_1$  and  $\sigma_2$  in the earlier equation,

$$1.96 \sqrt{\frac{2^2}{n} + \frac{2^2}{n}} = 1$$

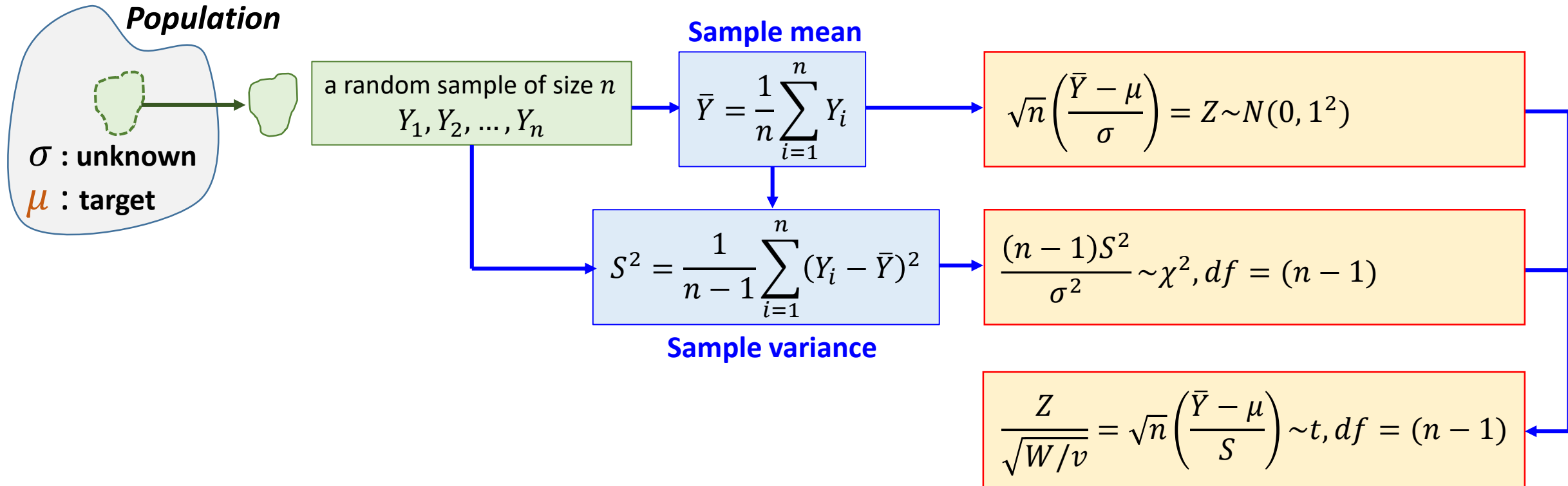
Solving, we obtain  $n = 30.73$ . Therefore, each group should contain  $n = 31$  members.

### Motivation

- The confidence intervals for a population mean  $\mu$  that we discuss in this section are based on the assumption that the experimenter's sample has been randomly selected from a normal population.
- The intervals are appropriate for samples of any size, and the confidence coefficients of the intervals are close to the specified values even when the population is not normal, as long as the departure from normality is not excessive.
- We rarely know the form of the population frequency distribution before we sample. Consequently, if an interval estimator is to be of any value, it must work reasonably well even when the population is not normal.
  - “Working well” means that the confidence coefficient should not be affected by modest departures from normality.
- For most mound-shaped population distributions, experimental studies indicate that these confidence intervals maintain confidence coefficients close to the nominal values used in their calculation.



# Estimation of the mean $\mu$ for the population distribution



- When  $Y_i \sim N(\mu, \sigma^2)$  with  $E(Y_i) = \mu$  and  $V(Y_i) = \sigma^2$
- When  $n$  is small (we cannot apply CLT)
- When  $\sigma$  is assumed to be **unknown**

$$\sqrt{n} \left( \frac{\bar{Y} - \mu}{S} \right) \sim t, df = (n-1)$$

### Estimation of the mean $\mu$ for the population distribution

- We assume that  $Y_1, Y_2, \dots, Y_n$  represent a random sample selected from a normal population, and we let  $\bar{Y}$  and  $S^2$  represent the sample mean and sample variance, respectively.
- We would like to construct a confidence interval for the population mean when  $V(Y_i) = \sigma^2$  is unknown and the sample size is too small to permit us to apply the large-sample techniques of the previous section.

$$T = \frac{\bar{Y} - \mu}{\sqrt{S^2/n}} = \frac{\bar{Y} - \mu}{S/\sqrt{n}} \sim t, df = (n - 1)$$

- The quantity  $T$  serves as the pivotal quantity that we will use to form a confidence interval for  $\mu$ .

## Estimation of the mean $\mu$ for the population distribution

- The quantity  $T$  serves as the pivotal quantity that we will use to form a confidence interval for  $\mu$ .

$$T = \frac{Z}{\sqrt{W/v}} = \sqrt{n} \left( \frac{\bar{Y} - \mu}{S} \right)$$

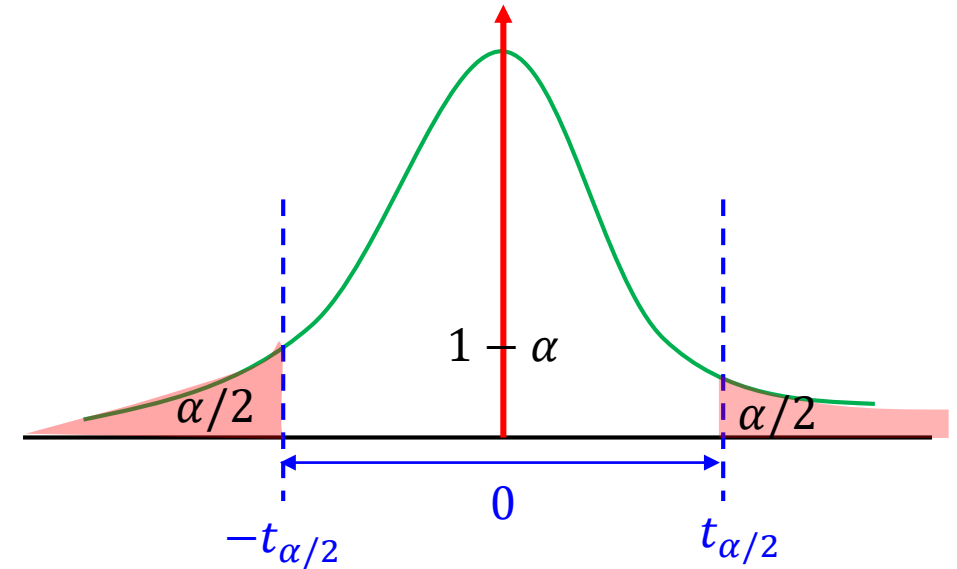
$$P(-t_{\alpha/2} \leq T \leq t_{\alpha/2}) = 1 - \alpha$$

$$P\left(-t_{\alpha/2} \leq \frac{\bar{Y} - \mu}{S/\sqrt{n}} \leq t_{\alpha/2}\right) = 1 - \alpha$$

$$P\left(\bar{Y} - t_{\alpha/2} \left(\frac{S}{\sqrt{n}}\right) \leq \mu \leq \bar{Y} + t_{\alpha/2} \left(\frac{S}{\sqrt{n}}\right)\right) = 1 - \alpha$$

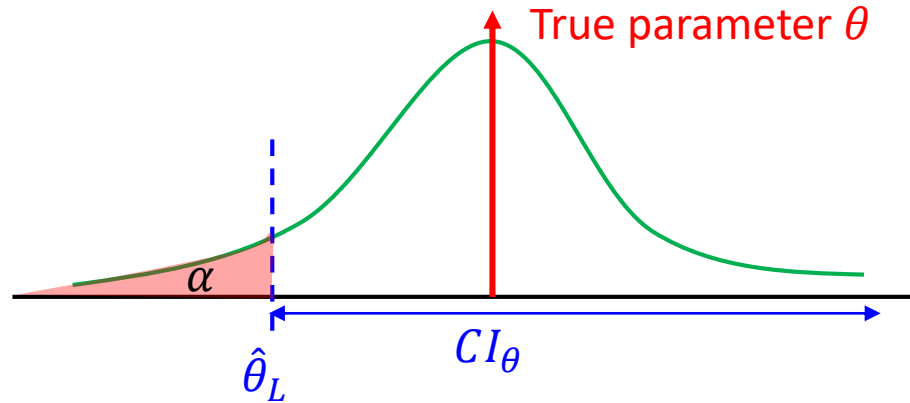
- The resulting confidence interval for  $\mu$

$$\bar{Y} \pm t_{\alpha/2} \left( \frac{S}{\sqrt{n}} \right)$$



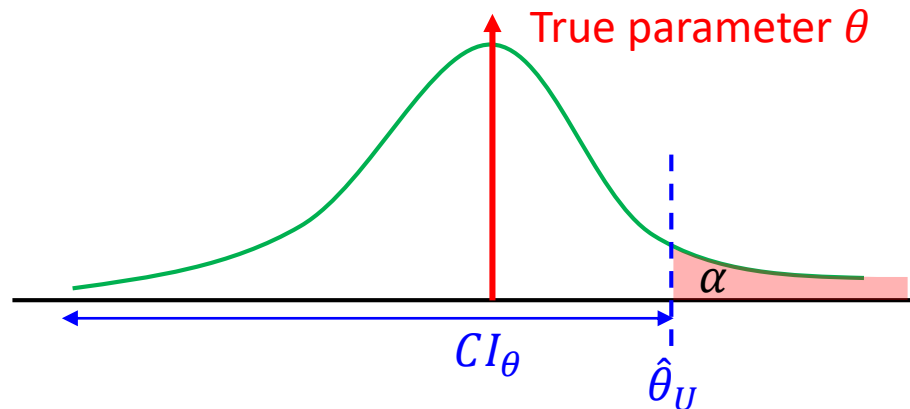
## Estimation of the mean $\mu$ for the population distribution

- $100(1 - \alpha)\%$  one-sided confidence limits for  $\mu$  (lower bound)

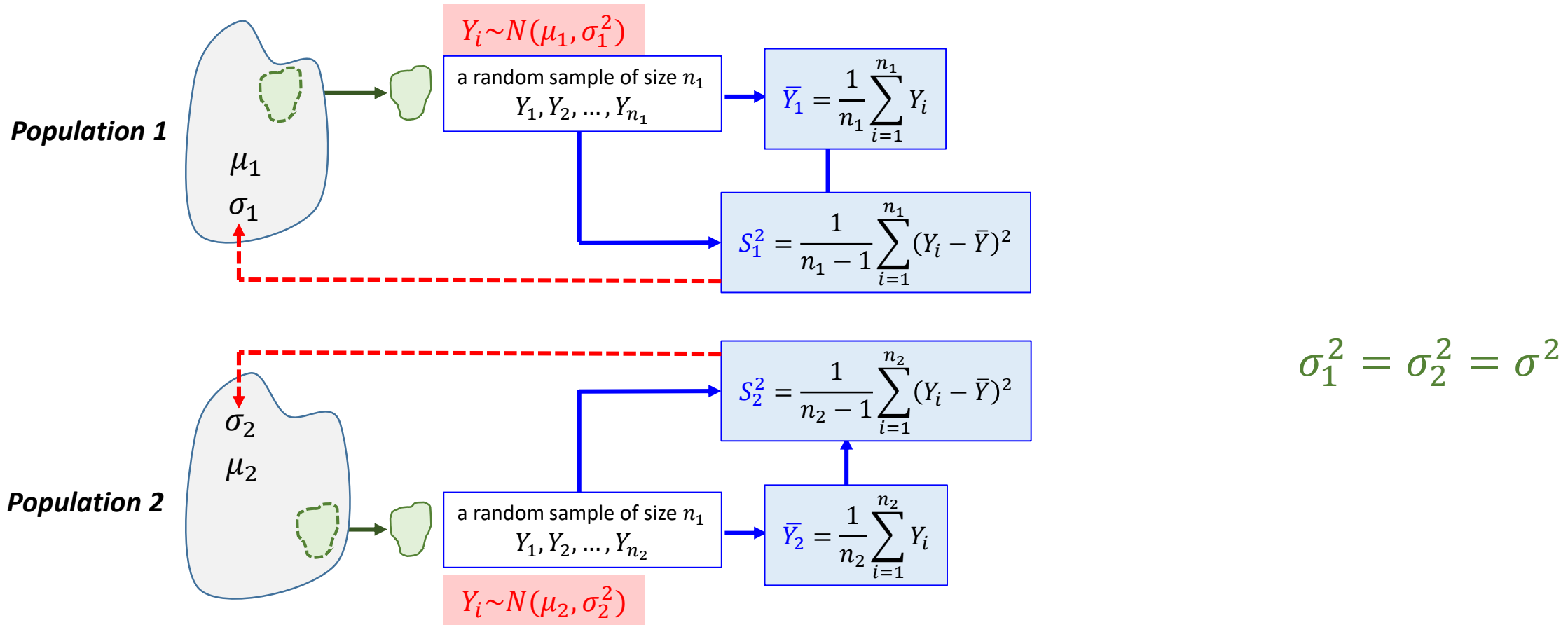


$$P\left(\bar{Y} - t_{\alpha/2}\left(\frac{S}{\sqrt{n}}\right) \leq \mu\right) = 1 - \alpha$$

- $100(1 - \alpha)\%$  one-sided confidence limits for  $\mu$  (upper bound)



$$P\left(\mu \leq \bar{Y} + t_{\alpha/2}\left(\frac{S}{\sqrt{n}}\right)\right) = 1 - \alpha$$

Estimation of  $\mu_1 - \mu_2$ 

- Suppose that we are interested in comparing the means of two normal populations,
- If the samples are independent, confidence intervals for  $\mu_1 - \mu_2$  based on a t-distributed random variable can be constructed if we assume that the two populations have a common but unknown variance  $\sigma_1^2 = \sigma_2^2 = \sigma^2$

### Estimation of $\mu_1 - \mu_2$

- If  $Y_1$  and  $Y_2$  are the respective sample means obtained from independent random samples from normal populations, the large-sample confidence interval for  $(\mu_1 - \mu_2)$  is developed by using as a pivotal quantity

$$Z = \frac{(\bar{Y}_1 - \bar{Y}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{(\bar{Y}_1 - \bar{Y}_2) - (\mu_1 - \mu_2)}{\sigma^2 \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

- Because  $\sigma^2$  is unknown, we need to find an estimator of the common variance  $\sigma^2$

$$S_P^2 = \frac{\sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_1)^2 + \sum_{i=1}^{n_2} (Y_{2i} - \bar{Y}_2)^2}{n_1 + n_2 - 2} = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

✓  $S_P^2$  is simply the average of  $S_1^2$  and  $S_2^2$ , with larger weight given to the sample variance associated with the larger sample size.

- $W$  is the sum of two independent  $\chi^2$ -distribution with  $(n_1 - 1)$  and  $(n_2 - 1)$  df, respectively

$$W = \frac{(n_1 + n_2 - 2)S_P^2}{\sigma^2} = \frac{\sum_{i=1}^{n_1} (Y_{1i} - \bar{Y}_1)^2}{\sigma^2} + \frac{\sum_{i=1}^{n_2} (Y_{2i} - \bar{Y}_2)^2}{\sigma^2}$$

✓  $W$  has a  $\chi^2$ -distribution with  $\nu = (n_1 - 1) + (n_2 - 1) = (n_1 + n_2 - 2)$

## Estimation of $\mu_1 - \mu_2$

- We now use the  $\chi^2$ -distributed variable  $W$  and the independent standard normal quantity  $Z$  defined in the previous paragraph to form a pivotal quantity:

$$T = \frac{Z}{\sqrt{W/v}} = \frac{\left[ \frac{(\bar{Y}_1 - \bar{Y}_2) - (\mu_1 - \mu_2)}{\sigma^2 \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right]}{\left[ \sqrt{\frac{(n_1 + n_2 - 2)S_P^2}{\sigma^2(n_1 + n_2 - 2)}} \right]} = \frac{(\bar{Y}_1 - \bar{Y}_2) - (\mu_1 - \mu_2)}{S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

a quantity that by construction has a  $t$  distribution with  $(n_1 + n_2 - 2)$  df.

- Proceeding as we did earlier in this section, we see that the confidence interval for  $(\mu_1 - \mu_2)$  has the form

$$(\bar{Y}_1 - \bar{Y}_2) \pm t_{\alpha/2} S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

## Estimation of $\mu_1 - \mu_2$

### Summary of Small-Sample Confidence Intervals for Means of Normal Distributions with Unknown Variance(s)

<i>Parameter</i>	<i>Confidence Interval (<math>v = df</math>)</i>
$\mu$	$\bar{Y} \pm t_{\alpha/2} \left( \frac{S}{\sqrt{n}} \right), \quad v = n - 1.$

$\mu_1 - \mu_2$	$(\bar{Y}_1 - \bar{Y}_2) \pm t_{\alpha/2} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}},$
-----------------	--

where  $v = n_1 + n_2 - 2$  and

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

(requires that the samples are independent and the assumption that  $\sigma_1^2 = \sigma_2^2$ ).

- As the sample size (or sizes) gets large, the number of degrees of freedom for the  $t$  distribution increases, and the  $t$  distribution can be approximated quite closely by the standard normal distribution
- The confidence intervals for a single mean and the difference in two means were developed **under the assumptions that the populations of interest are normally distributed**.
  - ✓ There is considerable empirical evidence that these intervals maintain their nominal confidence coefficient as long as the populations sampled have roughly mound-shaped distributions



## Example

### EXAMPLE 8.11

A manufacturer of gunpowder has developed a new powder, which was tested in eight shells. The resulting muzzle velocities, in feet per second, were as follows:

3005 2925 2935 2965  
2995 3005 2937 2905

Find a 95% confidence interval for the true average velocity  $\mu$  for shells of this type. Assume that muzzle velocities are approximately normally distributed.

## Example

### SOLUTION 8.11

If we assume that the velocities  $Y_i$  are normally distributed, the confidence interval for  $\mu$  is

$$\bar{Y} \pm t_{\alpha/2} \frac{S}{\sqrt{n}}$$

where  $t_{\alpha/2}$  is determined for  $n - 1$  df. For the given data,  $y = 2959$  and  $s = 39.1$ .

In this example, we have  $n - 1 = 7$  df and, using Table 5, Appendix 3,  $t_{\alpha/2} = 2.365$ .

Thus, we obtain

$$2959 \pm 2.365 \frac{39.1}{\sqrt{8}} \text{ or } 2959 \pm 32.7$$

## Example

### EXAMPLE 8.12

Table 8.3 Data for Example 8.12

Procedure	Measurements								
Standard	32	37	35	28	41	44	35	31	34
New	35	31	29	25	34	40	27	32	31

To reach maximum efficiency in performing an assembly operation in a manufacturing plant, new employees require approximately a 1-month training period. A new method of training was suggested, and a test was conducted to compare the new method with the standard procedure. Two groups of nine new employees each were trained for a period of 3 weeks, one group using the new method and the other following the standard training procedure. The length of time (in minutes) required for each employee to assemble the device was recorded at the end of the 3-week period. The resulting measurements are as shown in Table 8.3. Estimate the true mean difference ( $\mu_1 - \mu_2$ ) with confidence coefficient 0.95.

Assume that the assembly times are approximately normally distributed, that the variances of the assembly times are approximately equal for the two methods, and that the samples are independent.

## Example

### SOLUTION 8.12

For the data in Table 8.3, with sample 1 denoting the standard procedure, we have

$$\begin{aligned}\bar{y}_1 &= 35.22, & \bar{y}_2 &= 31.56, \\ \sum_{i=1}^9 (y_{1i} - \bar{y}_1)^2 &= 195.56, & \sum_{i=1}^9 (y_{2i} - \bar{y}_2)^2 &= 160.22, \\ s_1^2 &= 24.445, & s_2^2 &= 20.027.\end{aligned}$$

Hence,

$$s_p^2 = \frac{8(24.445) + 8(20.027)}{9 + 9 - 2} = \frac{195.56 + 160.22}{16} = 22.236$$

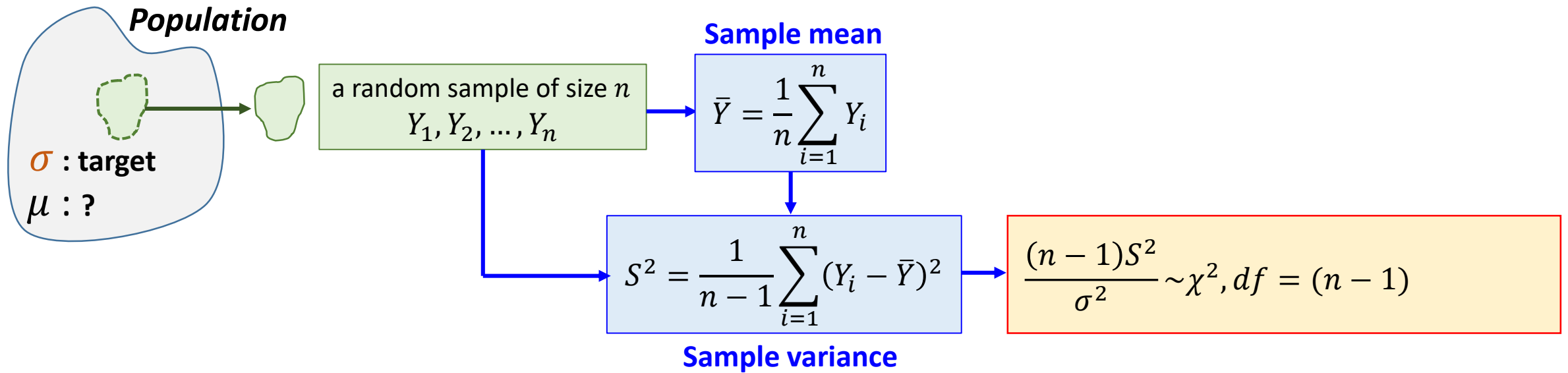
$$s_p = 4.716$$

Notice that, because  $n_1 = n_2 = 9$ ,  $s_p^2$  is the simple average of  $s_1^2$  and  $s_2^2$ . Also,  $t_{0.025} = 2.120$  for  $(n_1 + n_2 - 2) = 16$  df. The observed confidence interval is therefore

$$(\bar{y}_1 - \bar{y}_2) \pm t_{\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \quad \text{or} \quad 3.66 \pm 4.71 = [-1.05, 8.37]$$

The interval is fairly wide and includes both positive and negative values. If  $\mu_1 - \mu_2$  is positive,  $\mu_1 > \mu_2$  and the standard procedure has a larger expected assembly time than the new procedure. If  $\mu_1 - \mu_2$  is really negative, the reverse is true. Because the interval contains both positive and negative values, neither training method can be said to produce a mean assembly time that differs from the other.

# Motivation



- When  $Y_i \sim N(\mu, \sigma^2)$
- When  $n$  is small or large

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2, df = (n-1)$$

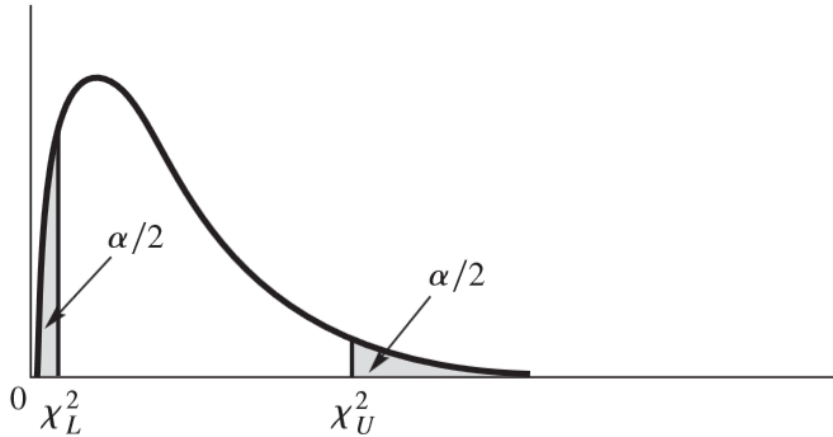
## Motivation

- We can then proceed by the pivotal method to find two numbers  $\chi_L^2$  and  $\chi_U^2$  such that

$$P \left[ \chi_L^2 \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi_U^2 \right] = 1 - \alpha$$

- We compromise by choosing points that cut off equal tail areas

$$P \left[ \chi_{1-\alpha/2}^2 \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi_{\alpha/2}^2 \right] = 1 - \alpha \Rightarrow P \left[ \frac{(n-1)S^2}{\chi_{\alpha/2}^2} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi_{1-\alpha/2}^2} \right] = 1 - \alpha$$



A  $100(1 - \alpha)$  Confidence interval for  $\sigma^2$

$$\left( \frac{(n-1)S^2}{\chi_{\alpha/2}^2}, \frac{(n-1)S^2}{\chi_{1-\alpha/2}^2} \right)$$

- The intervals for  $\sigma^2$  presented in this section can have confidence coefficients that differ markedly from the nominal level if the sampled population is not normally distributed.

## Example

### EXAMPLE 8.13

An experimenter wanted to check the variability of measurements obtained by using equipment designed to measure the volume of an audio source. Three independent measurements recorded by this equipment for the same sound were 4.1, 5.2, and 10.2. Estimate  $\sigma^2$  with confidence coefficient 0.90.

## Example

## SOLUTION 8.13

If normality of the measurements recorded by this equipment can be assumed, the confidence interval just developed applies. For the data given,  $s^2 = 10.57$ . With  $\alpha/2 = 0.05$  and  $(n - 1) = 2$  df, Table 6, Appendix 3, gives  $\chi^2_{0.95} = 0.103$  and  $\chi^2_{0.05} = 5.991$ . Thus, the 90% confidence interval for  $\sigma^2$  is

$$\left( \frac{(n-1)s^2}{\chi^2_{0.05}}, \frac{(n-1)s^2}{\chi^2_{0.95}} \right) \text{ or } \left( \frac{(2)(10.57)}{5.991}, \frac{(2)(10.57)}{0.103} \right) = (3.53, 205.24)$$

Notice that this interval for  $\sigma^2$  is very wide, primarily because  $n$  is quite small.