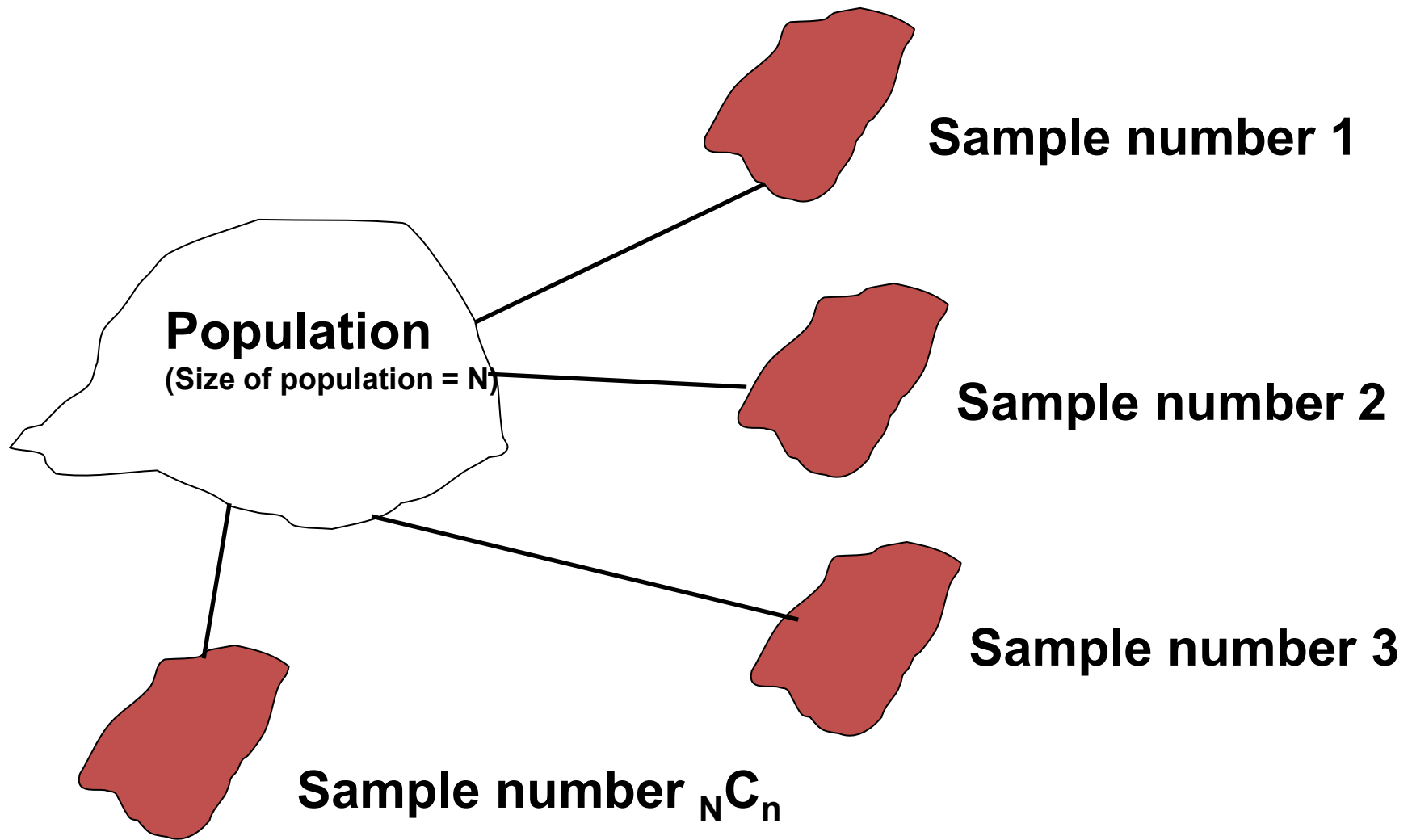


Lecture 7

Large and Small-Sample Estimation

Large sample estimation



Each sample size = n

INTRODUCTION

- ✗ Populations are described by their probability distributions and parameters.
 - + For quantitative populations, the location and shape are described by μ and σ .
 - + For a binomial populations, the location and shape are determined by p .
- ✗ If the values of parameters are unknown, we make inferences about them using sample information.

Types of Inference

- **Estimation:**
 - Estimating or predicting the value of the parameter
 - “What is (are) the most likely values of m or p ?”
- **Hypothesis Testing:**
 - Deciding about the value of a parameter based on some preconceived idea.
 - “Did the sample come from a population with $\mu = 5$ or $p = 0.2$?”

Types of Inference

- **Examples:**

- A consumer wants to estimate the average price of similar homes in her city before putting her home on the market.

Estimation: Estimate μ , the average home price.

- A manufacturer wants to know if a new type of steel is more resistant to high temperatures than an old type was.

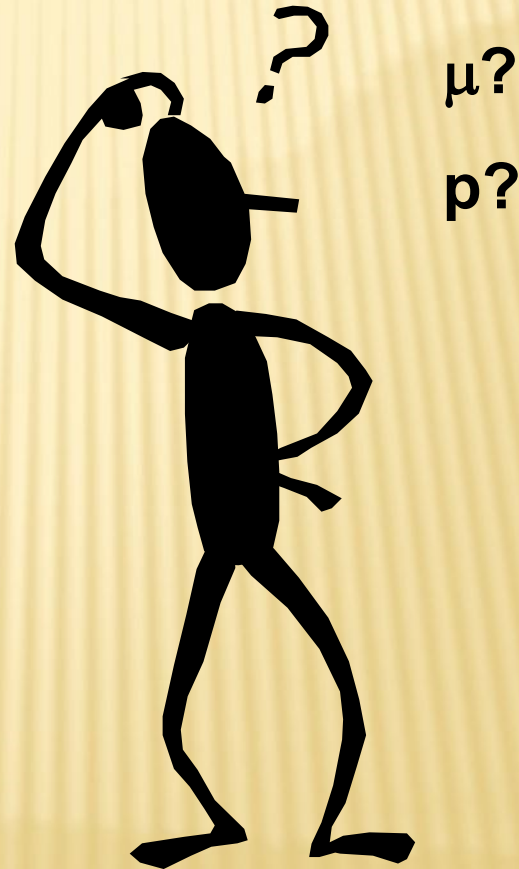
Hypothesis test: Is the new average resistance, μ_N equal to the old average resistance, μ_O ?

Types of Inference

- Whether you are estimating parameters or testing hypotheses, statistical methods are important because they provide:
 - Methods for making the inference
 - A numerical measure of the goodness or reliability of the inference

WHAT DO WE FREQUENTLY NEED TO ESTIMATE?

- ✗ An unknown population proportion p
- ✗ An unknown population mean μ



DEFINITIONS

- ✖ An **estimator** is a rule, usually a formula, that tells you how to calculate the estimate based on the sample.
 - + **Point estimation:** A single number is calculated to estimate the parameter.
 - + **Interval estimation:** Two numbers are calculated to create an interval within which the parameter is expected to lie.

A. Point Estimators

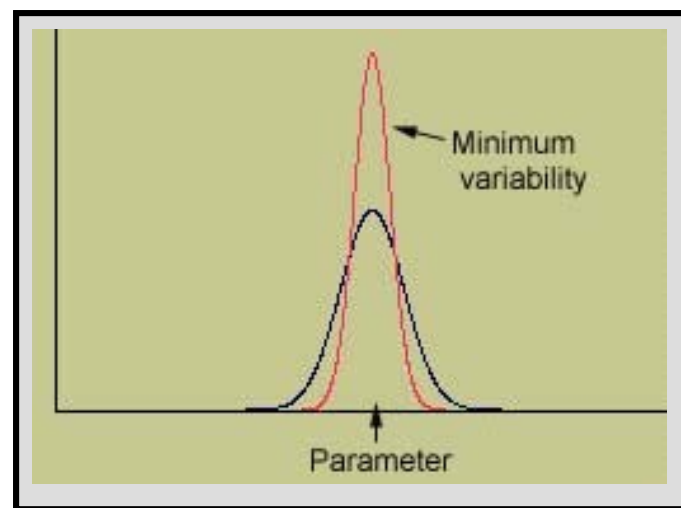
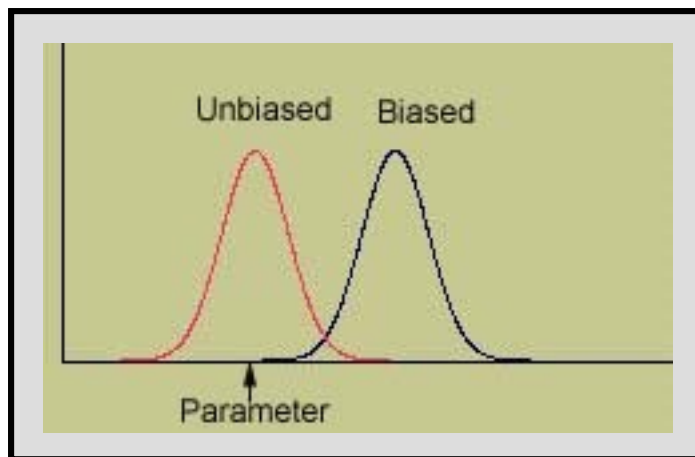
Properties

- ▶ Since an estimator is calculated from sample values, it varies from sample to sample according to its **sampling distribution**.
- ▶ An **estimator** is **unbiased** if the mean of its sampling distribution equals the parameter of interest.
- ▶ It does not systematically overestimate or underestimate the target parameter.



Properties

- ▶ Of all the **unbiased** estimators, we prefer the estimator whose sampling distribution has the **smallest spread** or **variability**.



Measuring the Goodness of an Estimator



- The distance between an estimate and the true value of the parameter is the **error of estimation**.

The distance between the bullet and the bull's-eye.

- When the sample sizes are large, our *unbiased* estimators will have **normal** distributions.

Because of the Central Limit Theorem.



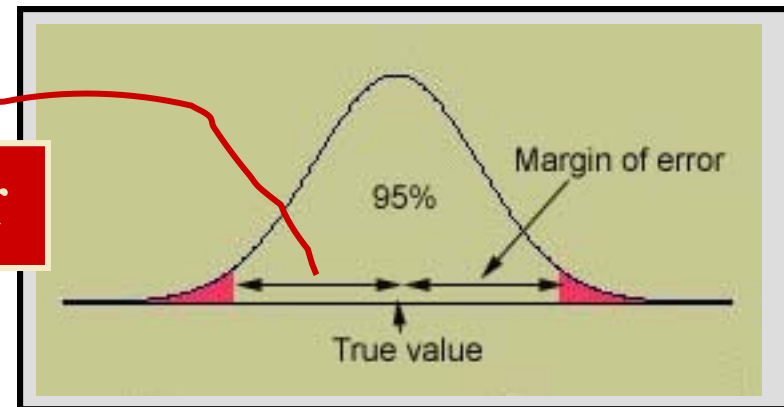
The Margin of Error



➤ **Margin of error:** The maximum error of estimation, is the maximum likely difference observed between sample mean \bar{x} and true population mean μ , calculated as :

1.645 1.96 2.33 2.575

$z_{\alpha/2} \times \text{std error of the estimator}$

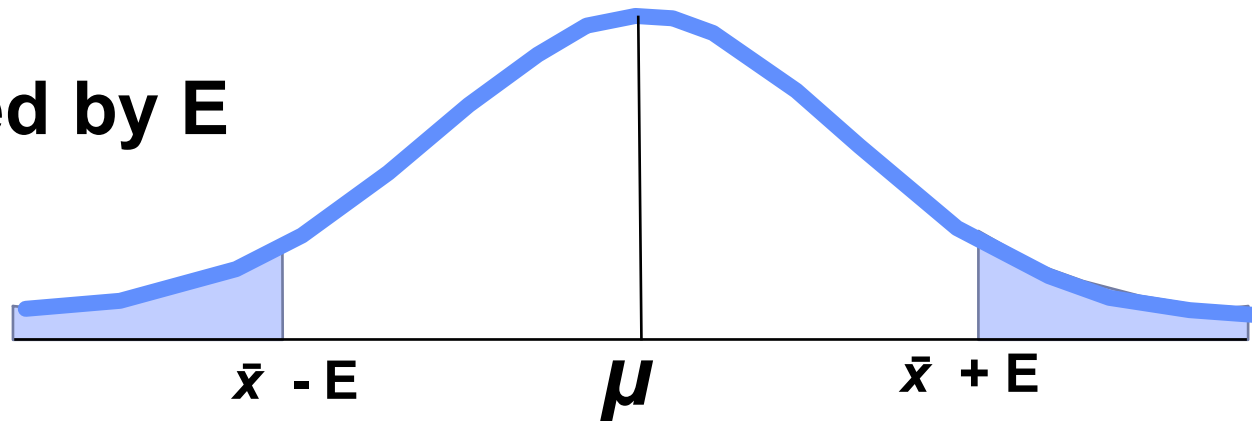


Definition

Margin of Error

is the maximum likely difference observed between sample mean \bar{x} and true population mean μ .

denoted by E



$$\bar{x} - E < \mu < \bar{x} + E$$

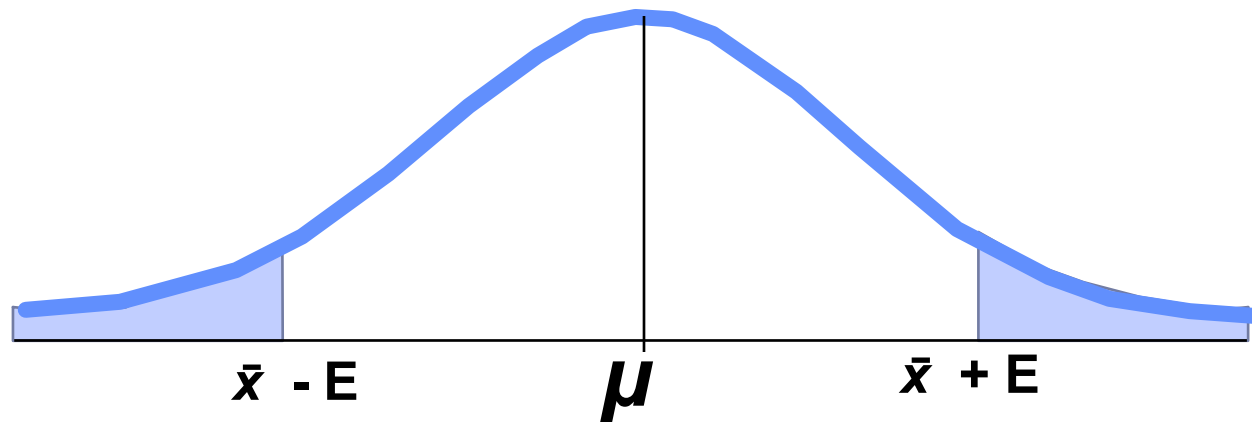
lower limit

upper limit

Definition

Margin of Error

$$E = z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$



also called the maximum error of the estimate

Estimating Means and Proportions

- For a quantitative population,

Point estimator of population mean $\mu : \bar{x}$

Margin of error ($n > 30$) : $\pm z_{\alpha/2} \frac{s}{\sqrt{n}}$

- For a binomial population,

Point estimator of population proportion $p : \hat{p} = x/n$

Margin of error ($n > 30$) : $\pm z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$

SE

1.645
1.96
2.33
2.575

Example



- A homeowner randomly samples 64 homes similar to her own and finds that the average selling price is \$252,000 with a standard deviation of \$15,000. Estimate the average selling price for all similar homes in the city.

Point estimator of μ : $\bar{x} = 252,000$

$$\text{Margin of error : } \pm 1.96 \frac{s}{\sqrt{n}} = \pm 1.96 \frac{15,000}{\sqrt{64}} = \pm 3675$$

Example



A quality control technician wants to estimate the proportion of soda cans that are underfilled. He randomly samples 200 cans of soda and finds 10 underfilled cans.

$n = 200$ $p =$ proportion of underfilled cans

Point estimator of p : $\hat{p} = x/n = 10/200 = .05$

Margin of error : $\pm 1.96 \sqrt{\frac{\hat{p}\hat{q}}{n}} = \pm 1.96 \sqrt{\frac{(.05)(.95)}{200}} = \pm .03$

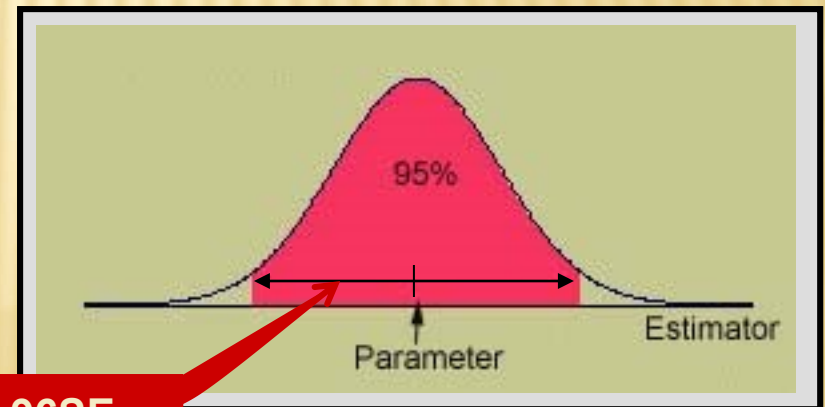
B. INTERVAL ESTIMATION



- Create an interval (a, b) so that you are fairly sure that the parameter lies between these two values.
- “Fairly sure” means “with high probability”, measured using the **confidence coefficient, $1-\alpha$** .

Usually, $1-\alpha = .90, .95, .98, .99$

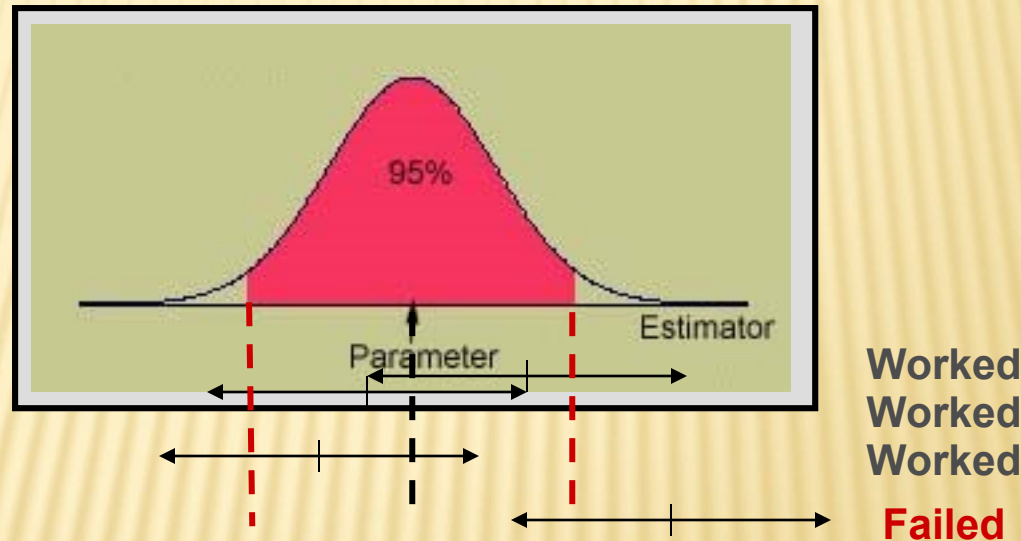
- Suppose $1-\alpha = .95$ and that the estimator has a normal distribution.



Parameter $\pm 1.96SE$

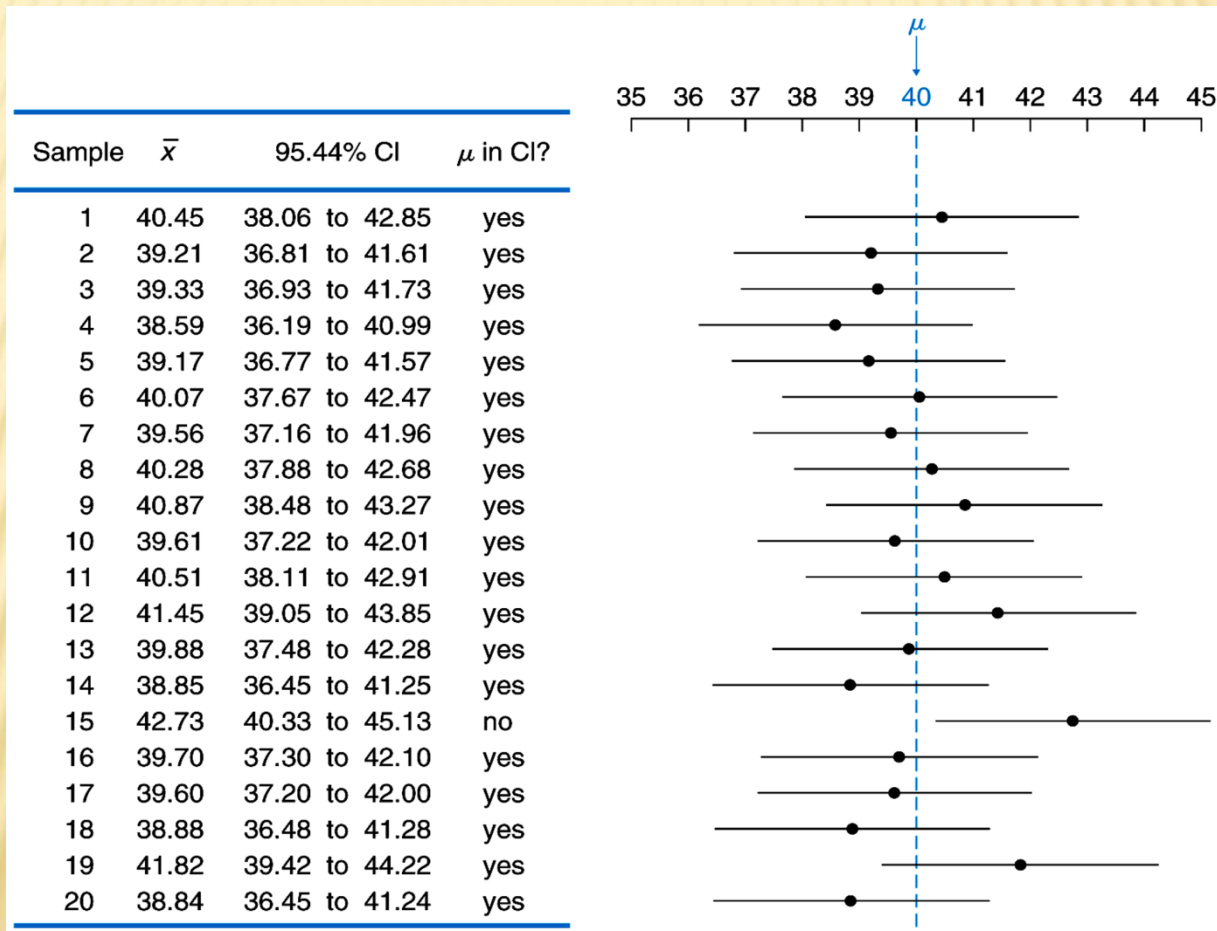
INTERVAL ESTIMATION (CONT'D)

- Since we don't know the value of the parameter, consider **Estimator ± 1.96 SE** which has a variable center.



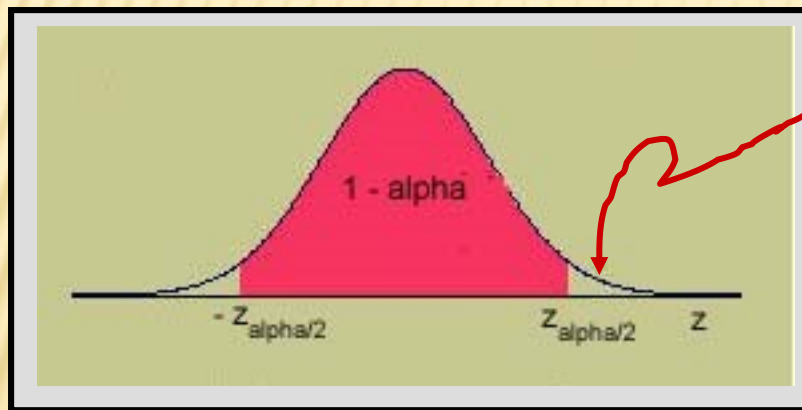
- Only if the estimator falls in the tail areas will the interval fail to enclose the parameter. This happens only 5% of the time.

INTERVAL ESTIMATION/CONFIDENCE INTERVALS



TO CHANGE THE CONFIDENCE LEVEL

- To change to a general confidence level, $1-\alpha$, pick a value of z that puts area $1-\alpha$ in the center of the z distribution.



Tail area	$z_{\alpha/2}$
.05	1.645
.025	1.96
.01	2.33
.005	2.575

$100(1-\alpha)\%$ Confidence Interval: Estimator $\pm z_{\alpha/2}SE$

1. CONFIDENCE INTERVALS FOR MEANS AND PROPORTIONS

- For a quantitative population,

Confidence interval for a population mean μ :

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$$

- For a binomial population,

Confidence interval for a population proportion p :

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}\hat{q}}{n}}$$

1.96

The diagram shows a box containing the value 1.96. Two arrows originate from this box: one points to the $z_{\alpha/2}$ term in the formula for the population mean confidence interval, and the other points to the $z_{\alpha/2}$ term in the formula for the population proportion confidence interval.

EXAMPLE



- A random sample of $n = 50$ males showed a mean average daily intake of dairy products equal to 756 grams with a standard deviation of 35 grams. Find a 95% confidence interval for the population average μ .

1.96

$$\bar{x} \pm z_{0.05/2} \frac{s}{\sqrt{n}} \Rightarrow 756 \pm 1.96 \frac{35}{\sqrt{50}} \Rightarrow 756 \pm 9.70$$

or $746.30 < \mu < 765.70$ grams.

EXAMPLE



- Find a 99% confidence interval for μ , the population average daily intake of dairy products for men.

2.575

$$\bar{x} \pm z_{0.01/2} \frac{s}{\sqrt{n}} \Rightarrow 756 \pm 2.58 \frac{35}{\sqrt{50}} \Rightarrow 756 \pm 12.75$$

or $743.25 < \mu < 768.75$ grams.

or $746.30 < \mu < 765.70$ grams.

The interval must be wider to provide for the increased confidence that it does indeed enclose the true value of μ .

EXAMPLE



- Of a random sample of $n = 150$ college students, 104 of the students said that they had played on a soccer team during their K-12 years. Estimate the proportion of college students who played soccer in their youth with a 98% confidence interval.

2.33

$$\hat{p} \pm z_{0.02/2} \sqrt{\frac{\hat{p}\hat{q}}{n}} \Rightarrow \frac{104}{150} \pm 2.33 \sqrt{\frac{.69(.31)}{150}}$$
$$\Rightarrow .69 \pm .09 \quad \text{or } .60 < p < .78.$$

2. ESTIMATING THE DIFFERENCE BETWEEN TWO MEANS

- Sometimes we are interested in comparing the means of two populations.
 - The average growth of plants fed using two different nutrients.
 - The average scores for students taught with two different teaching methods.
- To make this comparison,

A random sample of size n_1 drawn from population 1 with

A random sample of size n_2 drawn from population 2 with mean μ_2 and variance σ_2^2 .

ESTIMATING THE DIFFERENCE BETWEEN TWO MEANS (*CONT'D*)

- We compare the two averages by making inferences about $\mu_1 - \mu_2$, the difference in the two population averages.
 - If the two population averages are the same, then $\mu_1 - \mu_2 = 0$.
 - The best estimate of $\mu_1 - \mu_2$ is the difference in the two sample means,

$$\bar{x}_1 - \bar{x}_2$$



THE SAMPLING DISTRIBUTION OF $\bar{x}_1 - \bar{x}_2$

- Properties of the Sampling Distribution of $\bar{x}_1 - \bar{x}_2$

- Expected Value

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

- Standard Deviation/Standard Error

$$\sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

where: σ_1 = standard deviation of population 1

σ_2 = standard deviation of population 2

n_1 = sample size from population 1

n_2 = sample size from population 2



INTERVAL ESTIMATE OF $\mu_1 - \mu_2$:

LARGE-SAMPLE CASE ($n_1 > 30$ AND $n_2 > 30$)

Interval Estimate with σ_1 and σ_2 Known

$$\bar{x}_1 - \bar{x}_2 \pm z_{\alpha/2} \sigma_{\bar{x}_1 - \bar{x}_2}$$

SE

where:

$1 - \alpha$ is the confidence coefficient

□ Interval Estimate with σ_1 and σ_2 Unknown

$$\bar{x}_1 - \bar{x}_2 \pm z_{\alpha/2} s_{\bar{x}_1 - \bar{x}_2}$$

where:

$$s_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$\sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

EXAMPLE



Avg Daily Intakes	Men	Women
Sample size	50	50
Sample mean	756	762
Sample Std Dev	35	30

1.96

- Compare the average daily intake of dairy products of men and women using a 95% confidence interval.

$$(\bar{x}_1 - \bar{x}_2) \pm z_{0.05/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$
$$\Rightarrow (756 - 762) \pm 1.96 \sqrt{\frac{35^2}{50} + \frac{30^2}{50}} \Rightarrow -6 \pm 12.78$$
$$\text{or } -18.78 < \mu_1 - \mu_2 < 6.78.$$

EXAMPLE (CONT'D)

$$-18.78 < \mu_1 - \mu_2 < 6.78$$

- Could you conclude, based on this confidence interval, that there is a difference in the average daily intake of dairy products for men and women?
- The confidence interval contains the value $\mu_1 - \mu_2 = 0$. Therefore, it is possible that $\mu_1 = \mu_2$. You would not want to conclude that there is a difference in average daily intake of dairy products for men and women.



3. Estimating the Difference between Two Proportions

- Sometimes we are interested in comparing the proportion of “successes” in two binomial populations.
 - The germination rates of untreated seeds and seeds treated with a fungicide.
 - The proportion of male and female voters who favor a particular candidate for governor.
- To make this comparison,

A random sample of size n_1 drawn from

binomial population 1

A random sample of size n_2 drawn from

binomial population 2 with parameter p_2 .



Estimating the Difference between Two Proportions (*cont'd*)

- We compare the two proportions by making inferences about $p_1 - p_2$, the difference in the two population proportions.
 - If the two population proportions are the same, then $p_1 - p_2 = 0$.
 - The best estimate of $p_1 - p_2$ is the difference in the two sample proportions,

$$\hat{p}_1 - \hat{p}_2 = \frac{x_1}{n_1} - \frac{x_2}{n_2}$$

The Sampling Distribution of $\hat{p}_1 - \hat{p}_2$

- Expected Value/mean

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

- Standard Deviation/Standard Error

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

- Distribution Form

If the sample sizes are large ($n_1 p_1, n_1 q_1, n_2 p_2, n_2 q_2$) are all greater than to 5), the sampling distribution of $\hat{p}_1 - \hat{p}_2$ can be approximated by a normal probability distribution.

Interval Estimate of $p_1 - p_2$: Large-Sample Case

□ Interval Estimate with p_1 and p_2 Known

$$\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sigma_{\hat{p}_1 - \hat{p}_2}$$

where:

$1 - \alpha$ is the confidence coefficient

□ Interval Estimate with p_1 and p_2 Unknown

$$\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sigma_{\hat{p}_1 - \hat{p}_2}$$

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

where: $\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$

Example



Youth Soccer	Male	Female
Sample size	80	70
Played soccer	65	39

- Compare the proportion of male and female college students who said that they had played on a soccer team during their K-12 years using a 99% confidence interval.

2.575

$$(\hat{p}_1 - \hat{p}_2) \pm z_{0.01/2} \sqrt{\frac{\hat{p}_1 \hat{q}_1}{n_1} + \frac{\hat{p}_2 \hat{q}_2}{n_2}}$$

$$\Rightarrow \left(\frac{65}{80} - \frac{39}{70} \right) \pm 2.575 \sqrt{\frac{.81(.19)}{80} + \frac{.56(.44)}{70}} \Rightarrow 0.26 \pm 0.19$$

$$\text{or } 0.07 < p_1 - p_2 < 0.45$$

Example (cont'd)



$$0.07 < p_1 - p_2 < 0.45$$

- Could you conclude, based on this confidence interval, that there is a difference in the proportion of male and female college students who said that they had played on a soccer team during their K-12 years?
- The confidence interval does not contain the value $p_1 - p_2 = 0$. Therefore, it is not likely that $p_1 = p_2$. You would conclude that there is a difference in the proportions for males and females.

A higher proportion of males than females played soccer in their youth.

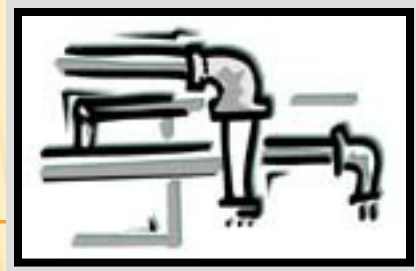
ONE SIDED CONFIDENCE BOUNDS

- ✗ Confidence intervals are by their nature **two-sided** since they produce upper and lower bounds for the parameter.
- ✗ **One-sided bounds** can be constructed simply by using a value of z that puts α rather than $\alpha/2$ in the tail of the z distribution.

LCB : Estimator $- z_{\alpha} \times (\text{Std Error of Estimator})$

UCB : Estimator $+ z_{\alpha} \times (\text{Std Error of Estimator})$

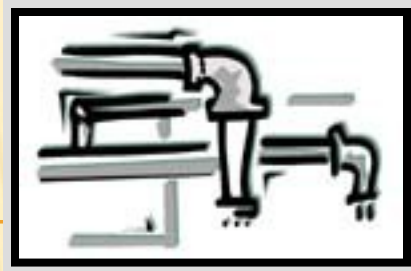
CHOOSING THE SAMPLE SIZE



- ✗ The total amount of relevant information in a sample is controlled by two factors:
 - The **sampling plan** or **experimental design**: the procedure for collecting the information
 - The **sample size n** : the amount of information you collect.
- ✗ In a statistical estimation problem, the accuracy of the estimation is measured by the **margin of error** or the **width of the confidence interval**.

CHOOSING THE SAMPLE SIZE (CONT'D)

1. Determine the size of the **margin of error, E** , that you are willing to tolerate.
2. Choose the sample size by solving for n or $n = n_1 = n_2$ in the inequality: **$1.96 SE \leq E$** , where SE is a function of the sample size n .
3. For quantitative populations, estimate the population standard deviation using a previously calculated value of **s** or the range approximation **$\sigma \approx \text{Range} / 4$** .
4. For binomial populations, use the conservative approach and approximate p using the value **$p = .5$** .



EXAMPLE

A producer of PVC pipe wants to survey wholesalers who buy his product in order to estimate the proportion who plan to increase their purchases next year. What sample size is required if he wants his estimate to be within 0.04 of the actual proportion with probability equal to 0.95?

$$1.96\sqrt{\frac{pq}{n}} \leq .04 \Rightarrow 1.96\sqrt{\frac{0.5(0.5)}{n}} \leq 0.04$$

$$\Rightarrow \sqrt{n} \geq \frac{1.96\sqrt{0.5(0.5)}}{0.04} = 24.5 \Rightarrow n \geq 24.5^2 = 600.25$$

He should survey at least 600 wholesalers.

4. Estimating the Variance

The **sample variance** is defined by

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$
$$= \frac{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}}{n-1}$$

$$\begin{aligned} \text{Var}(X) &= E(X - E(X))^2 \\ &= E(X - \mu)^2 = \sigma^2 \end{aligned}$$

$$\begin{aligned} E\left(\sum_{i=1}^n (x_i - \bar{x})^2\right) &= E\left(\sum_{i=1}^n (x_i - \mu)(\mu - \bar{x})\right)^2 \\ &= \sum E[(x_i - \mu)^2] - nE[(\bar{x} - \mu)^2] \\ &= n\sigma^2 - \frac{n\sigma^2}{n} \end{aligned}$$

$$E\left(\sum_{i=1}^n (x_i - \mu)^2\right) = (n-1)\sigma^2$$

$$E\left(\frac{\sum_{i=1}^n (x_i - \mu)^2}{n-1}\right) = \sigma^2 \Rightarrow E(s^2) = \sigma^2$$

Analysis of Sample Variance

If s^2 is the variance of a random sample size n from a normal population, a $100(1-\alpha)\%$ confidence interval for σ^2 is

$$\frac{(n-1)s^2}{\chi_{\alpha/2}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{1-\alpha/2}^2}$$

Where $\chi_{\alpha/2}^2$ and $\chi_{1-\alpha/2}^2$ are χ^2 values with $(n-1)$ degrees of freedom.

