

# Five

---

## Estimation

---

We now learn about a powerful use of statistics:

### STATISTICAL INFERENCE

about POPULATION PARAMETERS

using SAMPLE DATA.

In case you wonder about the relevance of learning about probability and sampling distribution, this is why:

- Statistical inference methods use probability calculations that assume that the data were gathered with a random sample or a randomized experiment.
- The probability calculations refer to a sampling distribution of a statistic, which is often approximately a normal distribution.

There are two types of statistical inference methods

- estimation of population parameters; and
- testing hypotheses about the parameter values.

This chapter discusses the first — estimating population parameters.

### TWO TYPES OF ESTIMATIONS

#### Point estimation

Based on sample data, a single number is calculated to estimate the population parameter. The rule or formula that describes this calculation is called the **point estimator**. The resulting number is called a **point estimate**.

#### Interval estimation

Based on sample data, two numbers are calculated to form an interval within which the parameter is expected to lie.

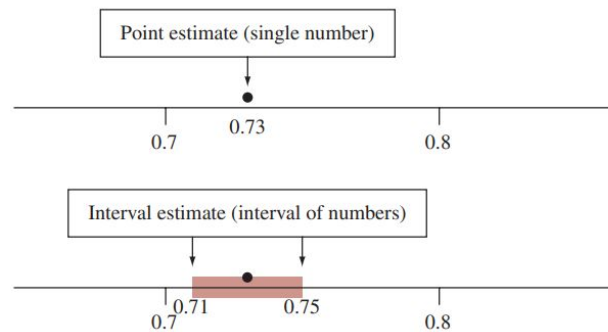
**EXAMPLE 5.1**

One survey asked, "Do you believe in hell?"

From **sample** data, the **point estimate** for the proportion of adult (in the **population**) who would respond "yes" is 0.73. The adjective "point" refers to using a single number as the parameter estimate.

An **interval estimate** predicts that the proportion of adult (in the **population**) who believe in hell falls between 0.71 and 0.75.

The next figure illustrates the difference between **point estimate** and **interval estimate** for the previous example.



## 1 POINT ESTIMATION

Suppose we are interested to estimate the parameter  $\mu$ , the population mean. Assume that we have the following data, a random sample consisting

$$X_1, X_2, \dots, X_n.$$

**DEFINITION 1 (ESTIMATOR)**

An **estimator** is a rule, usually expressed as a formula, that tells us how to calculate an **estimate** based on information in the sample.

**EXAMPLE 5.2 (POINT ESTIMATOR)**

We want to estimate the average waiting time for a bus ( $\mu$ ) for students attending ST2334. The lecturer asked 4 students their waiting times  $X_1, \dots, X_4$  for a bus. The (observed) results are

$$x_1 = 6, x_2 = 1, x_3 = 4, x_4 = 9.$$

We can use  $\bar{X} = \frac{1}{4}(X_1 + \dots + X_4)$  to estimate  $\mu$ . In this case,  $\bar{X}$  is the **estimator** (for  $\mu$ ), and the computed value  $\bar{x} = 5$  is the **estimate**.

**QUESTIONS**

- How good is the estimator?
- What would be a criteria for a “good” estimator?

**Unbiased Estimator**

One of the reasons we think  $\bar{X}$  is a good estimator of  $\mu$  is because  $E(\bar{X}) = \mu$ . That is, “on average”, the estimator is right.

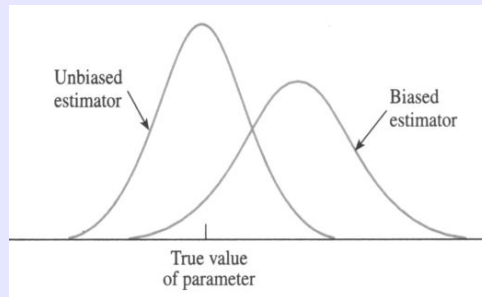
In general, we represent the parameter of interest by  $\theta$ . For example,  $\theta$  can be  $p, \mu$ , or  $\sigma$ .

**DEFINITION 2 (UNBIASED ESTIMATOR)**

Let  $\hat{\Theta}$  be an estimator of  $\theta$ . Then  $\hat{\Theta}$  is a random variable based on the sample. If  $E(\hat{\Theta}) = \theta$ , we call  $\hat{\Theta}$  an **unbiased estimator** of  $\theta$ .

**REMARK**

An unbiased estimator has mean value equals to the true value of the parameter.

**EXAMPLE 5.3 (UNBIASED ESTIMATOR)**

Let  $X_1, X_2, \dots, X_n$  be a random sample from the same population with mean  $\mu$  and variance  $\sigma^2$ . Then

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

is an unbiased estimator of  $\sigma^2$  since  $E(S^2) = \sigma^2$ .

**L-EXAMPLE 5.1 (UNBIASED ESTIMATOR)**

A bus arrives at the bus stop according to a  $U(0, \theta)$  distribution. The lecturer wants to estimate  $\theta$ . So this morning he randomly selected 4 students and asked their waiting time for a bus. The following values are obtained:

$$X_1, \dots, X_4.$$

- (a) Is  $\bar{X}$  an unbiased estimator of  $\theta$ ?
- (b) Using (a), construct an unbiased estimator of  $\theta$ .
- (c) Is there another unbiased estimator of  $\theta$ ?

**Solution:**

We know that  $X_1, \dots, X_4$  have the same distribution as  $X \sim U(0, \theta)$ .

- (a) No. This is because

$$E(\bar{X}) = E(X) = \int_0^\theta \frac{x}{\theta} dx = \frac{\theta}{2} \neq \theta.$$

- (b)  $2\bar{X}$  is one. This is because  $E(2\bar{X}) = 2E(\bar{X}) = \theta$ .
- (c) Yes.  $2X_1$  is another since  $E(2X_1) = E(2X) = \theta$ .

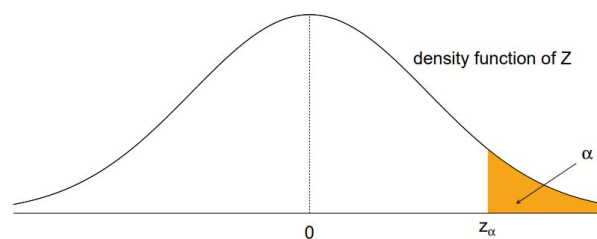
**Maximum Error of Estimate**

Typically  $\bar{X} \neq \mu$ , so  $\bar{X} - \mu$  measures the difference between the estimator and the true value of the parameter.

Recall that if the population is normal or if  $n$  is large,  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  follows a standard normal or an approximately standard normal distribution.

**DEFINITION 3 ( $z_\alpha$ )**

Define  $z_\alpha$  to be the number with an upper-tail probability of  $\alpha$  for the standard normal distribution  $Z$ . That is,  $P(Z > z_\alpha) = \alpha$ .

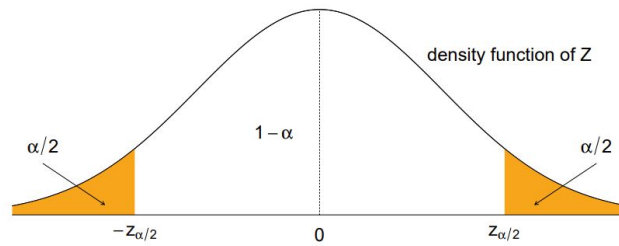


From the above definition, we then have

$$P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = 1 - \alpha.$$

In other words,

$$P\left(\frac{|\bar{X} - \mu|}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = P\left(|\bar{X} - \mu| \leq z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$



This means that, with probability  $1 - \alpha$ , the error  $|\bar{X} - \mu|$  is less than

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

#### DEFINITION 4 (MAXIMUM ERROR OF ESTIMATE)

The quantity

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

is called the *maximum error of estimate*.

#### EXAMPLE 5.4 (TV TIME FOR INTERNET USERS)

An investigator is interested in the amount of time internet users spend watching television per week.

Based on historical experience, he assumes that the standard deviation is  $\sigma = 3.5$  hours.

He proposes to select a random sample of  $n = 50$  internet users, poll them, and take the sample mean to estimate the population mean  $\mu$ .

What can he assert with probability 0.99 about the maximum error of estimate?

**Solution:**

As  $n = 50 \geq 30$  is large,  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$  is approximately normal.

So we can use the previous result, with  $\sigma = 3.5$ ,  $\alpha = 0.01$  and  $z_{\alpha/2} = z_{0.005} = 2.576$ .

With probability 0.99, the error is at most

$$E = 2.576 \times \frac{3.5}{\sqrt{50}} \approx 1.27.$$

#### REMARK

$z_{0.005}$  is the same as the 0.995 quantile of the standard normal. The value of 2.576 can be obtained from tables or software.

Use the command `qnorm(0.995)` or `qnorm(0.005, lower.tail=F)` to obtain the value via <https://rdr.io/snippets/>.

Alternatively, you may use Radian to get the same value as well.

#### Determination of Sample Size

We often want to know what the minimum sample size should be, so that with probability  $1 - \alpha$ , the error is at most  $E_0$ .

To answer this, consider the fact that we want

$$z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq E_0.$$

Solving for  $n$ , we have

$$n \geq \left( \frac{z_{\alpha/2} \cdot \sigma}{E_0} \right)^2.$$

#### L-EXAMPLE 5.2 (TV TIME FOR INTERNET USERS II)

What is the sample size  $n$  required such that the television investigator can assert with 99% probability that his estimation error is at most 0.5 hour?

**Solution:**

The sample size required is

$$n \geq \left( \frac{2.576 \times 3.5}{0.5} \right)^2 \approx 325.15.$$

**Different Cases**

We had previously understood the sampling distribution of  $\bar{X}$  for a variety of cases. Repeating the same arguments above, we have the following table.

**DIFFERENT CASES**

	Population	$\sigma$	$n$	Statistic	$E$	$n$ for desired $E_0$ and $\alpha$
I	Normal	known	any	$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$	$z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$	$\left( \frac{z_{\alpha/2} \cdot \sigma}{E_0} \right)^2$
II	any	known	large	$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$	$z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$	$\left( \frac{z_{\alpha/2} \cdot \sigma}{E_0} \right)^2$
III	Normal	unknown	small	$T = \frac{\bar{X} - \mu}{s/\sqrt{n}}$	$t_{n-1; \alpha/2} \cdot \frac{s}{\sqrt{n}}$	$\left( \frac{t_{n-1; \alpha/2} \cdot s}{E_0} \right)^2$
IV	any	unknown	large	$Z = \frac{\bar{X} - \mu}{s/\sqrt{n}}$	$z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$	$\left( \frac{z_{\alpha/2} \cdot s}{E_0} \right)^2$

**L-EXAMPLE 5.3 (TV TIME FOR INTERNET USERS III)**

Back to the case where the investigator polls  $n = 50$  internet users.

Suppose we do not trust the historical assumption that  $\sigma = 3.5$ . Instead, we estimate  $\sigma$  using the sample standard deviation  $s = 2.6$ .

With 99% confidence, what is  $E$ , the maximum error of our estimate?

**Solution:**

We are in Case IV. So  $E$  is given as

$$E = z_{\alpha/2} \cdot \frac{s}{\sqrt{n}} = 2.576 \cdot \frac{2.6}{\sqrt{50}} \approx 0.947.$$

## 2 CONFIDENCE INTERVALS FOR THE MEAN

Since a point estimate is almost never right, one might be interested in asking for an interval where the parameter lies in.

### DEFINITION 5 (CONFIDENCE INTERVAL)

An **interval estimator** is a rule for calculating, from the sample, an interval  $(a, b)$  in which you are fairly certain the parameter of interest lies in.

This “fairly certain” can be quantified by the **degree of confidence** also known as **confidence level**  $(1 - \alpha)$ , in the sense that

$$P(a < \mu < b) = 1 - \alpha.$$

$(a, b)$  is called the  $(1 - \alpha)$  **confidence interval**.

### Case I: $\sigma$ known, data normal

Consider the case where  $\sigma$  is known, and data comes from a normal population.

We learnt previously that

$$P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = 1 - \alpha.$$

Rearranging, we have

$$P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

So

$$\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = \left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$$

is a  $(1 - \alpha)$  confidence interval.



**EXAMPLE 5.5**

In order to set inventory levels, a computer company samples **demand during lead time** over 25 time periods:

235 374 309 499 253 421 361 514 462 369 394 439  
348 344 330 261 374 302 466 535 386 316 296 332 334

It is known that the (population) standard deviation of **demand over lead time** is 75 computers. Given that  $\bar{x} = 370.16$ , estimate the mean demand over lead time with 95% confidence. Assume a normal distribution for the population.

**Solution:**

Note that  $z_{\alpha/2} = z_{0.025} = 1.96$ . The 95% confidence interval is

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 370.16 \pm 1.96 \frac{75}{\sqrt{25}} = 370.16 \pm 29.4$$

or (340.76, 399.56).

**REMARK**

Notice that our  $(1 - \alpha)$  confidence interval can be written as  $\bar{X} \pm E$ .

This is not a coincidence: recall that there is  $(1 - \alpha)$  confidence that the error  $|\bar{X} - \mu|$  is within  $E$ .

For the other cases, based on our understanding of the sampling distribution of  $\bar{X}$ , we can construct our confidence intervals for the different cases  $\bar{X} \pm E$ , based on the conditions given.

**CONFIDENCE INTERVALS FOR THE MEAN**

The table below gives the  $(1 - \alpha)$  confidence interval (formulas) for the population mean.

Case	Population	$\sigma$	$n$	Confidence Interval
I	Normal	known	any	$\bar{x} \pm z_{\alpha/2} \cdot \sigma / \sqrt{n}$
II	any	known	large	$\bar{x} \pm z_{\alpha/2} \cdot \sigma / \sqrt{n}$
III	Normal	unknown	small	$\bar{x} \pm t_{n-1; \alpha/2} \cdot s / \sqrt{n}$
IV	any	unknown	large	$\bar{x} \pm z_{\alpha/2} \cdot s / \sqrt{n}$

Note that  $n$  is considered large when  $n \geq 30$ .

**EXAMPLE 5.6 (WHICH CASE?)**

The following data set collects  $n = 41$  randomly sampled waiting times of students from ST2334 to receive reply for their email from a survey in the day time.

2.50	23.28	19.34	4.74	7.03	21.85	2.72
17.73	21.55	9.71	30.24	0.37	31.26	35.24
7.81	16.69	66.54	1.88	14.14	46.59	28.17
0.06	9.32	0.03	10.75	6.97	56.86	2.89
7.67	30.16	0.33	0.44	3.77	25.07	7.05
0.08	10.64	13.10	7.92	112.77	11.93	

Given that  $\bar{x} = 17.736$  and  $s = 21.7$ , construct a 98% confidence interval for the mean waiting time of *all* ST2334 students.

**Solution:**

Note that  $\sigma$  is unknown, and  $n$  is large. So we are in Case IV.

Note that  $z_{\alpha/2} = z_{0.01} = 2.326$ . So our 98% confidence interval is

$$\begin{aligned}\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}} &= 17.736 \pm 2.326 \times \frac{21.7}{\sqrt{41}} \\ &= (9.85, 25.62).\end{aligned}$$

**EXAMPLE 5.7 (WHICH CASE AGAIN?)**

The contents of 7 similar containers of sulphuric acid (in litres) are

9.8	10.2	10.4	9.8	10.0	10.2	9.6
-----	------	------	-----	------	------	-----

It can shown that  $\bar{x} = 10$  and  $s^2 = 0.08$ . Find a 95% confidence interval for the mean content of all such containers, assuming an approximate normal distribution for container contents.

**Solution:**

We are in Case III.

Using software, we obtain  $t_{6;0.025} = 2.447$ .

Thus a 95% confidence interval for the mean content of all such containers is given as

$$\bar{x} \pm t_{n-1;\alpha/2} \cdot \frac{s}{\sqrt{n}} = 10 \pm 2.447 \cdot \frac{\sqrt{0.08}}{\sqrt{7}} = (9.738, 10.262).$$

**L-EXAMPLE 5.4 (AGAIN, WHICH CASE?)**

A major department store chain is interested in estimating the average amount its credit card customers spent on their first visit to the chain's new store in the mall.

Fifty credit card accounts were randomly sampled and analyzed with the following results:

$$\bar{x} = \$62.56 \quad \text{and} \quad s = \$20.$$

- Identify the population the department store chain is interested in learning about.
- Which population parameter does the chain wish to estimate?
- Construct a 90% confidence interval for the parameter identified in the previous part.

**Solution:**

- The population is all its credit card customers.
- We are interested in  $\mu$ , the average amount credit card customers spent on their first visit to the chain's new store in the mall.
- As  $n$  is large and  $\sigma$  is unknown, we are in Case IV.

Note that  $z_{0.05} = 1.645$ . Thus the 90% confidence interval for  $\mu$  is

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{s}{\sqrt{n}} = 62.56 \pm 1.645 \cdot \frac{20}{\sqrt{50}} = (57.907, 67.213).$$

**L-EXAMPLE 5.5 (BUYING PHONES ON EBAY)**

eBay is a popular Internet company for auctioning just about anything.

The following is a random sample of 11 completed auctions for an unlocked Apple iPhone 5s with 16GB storage in new condition (item not used, but original packaging might be missing), obtained from eBay in July 2014.

Closing Price (in \$):

570 620 610 590 540 590 565 590 580 570 595

We are given that  $\bar{x} = 583.64$ ,  $s = 22.15$ .

The retail price of this phone is \$649.

By constructing a 95% confidence interval for the mean closing price (of this phone) on eBay, find out how much can you save by buying items on eBay compared to their actual retail price. Assume an approximate normal distribution for the closing price.

**Solution:**

Let  $\mu$  denote the population mean for the closing price of the auction.

The point estimate of  $\mu$  is the sample mean  $\bar{x} = 583.64$ .

We are in Case III, and a 95% confidence interval for  $\mu$  is given as

$$\bar{x} \pm t_{10;0.025} \cdot \frac{s}{\sqrt{n}} = 583.64 \pm 2.228 \cdot \frac{22.15}{\sqrt{11}} = (568.8, 598.5).$$

The upper bound of this confidence interval is \$50 below the retail price of \$649, so some potential savings can be made when buying this phone from eBay.

**INTERPRETING CONFIDENCE INTERVALS I**

- We saw that  $\bar{X} \pm E$  has probability  $(1 - \alpha)$  of containing  $\mu$ .

This is a probability statement about the **procedure** by which we compute the interval — the **interval estimator**.

- Each time we take a sample, and go through this construction, we get a different confidence interval.
- Sometimes we get a confidence interval that **contains**  $\mu$ , and sometimes we get one that **does not contain**  $\mu$ .
- Once an interval is **computed**,  $\mu$  is either in it or not. There is no more randomness.

**INTERPRETING CONFIDENCE INTERVALS II**

- Since  $\mu$  is typically not known, there is no way to determine whether a particular confidence interval succeeded in capturing the population mean.
- However, if we repeat this procedure of taking a sample and computing a confidence interval many times, about  $(1 - \alpha)$  of the many confidence intervals that we get will contain the true parameter.

This is what “confidence” means — a **confidence in the method used**.

- The following R Shiny app allows us to explore this fact:  
<https://istats.shinyapps.io/ExploreCoverage/>

**3 COMPARING TWO POPULATIONS**

In real applications, it is quite common to compare the means of two populations.

Imagine that we have two populations

- Population 1 has mean  $\mu_1$ , variance  $\sigma_1^2$ .
- Population 2 has mean  $\mu_2$ , variance  $\sigma_2^2$ .

### Experimental Design

In order to compare two populations, a number of observations from each population need to be collected. Experimental design refers to the manner in which samples from populations are collected.

#### TWO BASIC DESIGNS FOR COMPARING TWO TREATMENTS

- Independent samples — complete randomization.
- Matched pairs samples — randomization between matched pairs.

#### EXAMPLE 5.8 (INDEPENDENT SAMPLES)

In order to compare the examination scores of male and female students attending ST2334,

- 10 scores of female students are randomly sampled — Sample I,
- 8 scores of male students are randomly sampled — Sample II.

Note that all observations are independent —

- Sample I and Sample II are independent;
- Individuals within Sample I are independent;
- Individuals within Sample II are independent.

#### EXAMPLE 5.9 (MATCHED PAIRS SAMPLES)

In order to study whether there exists income difference between male and female, 100 **married couples** are sampled, and their monthly incomes are collected.

In this example, the treatment groups are the female group and male group.

Note that observations are dependent in a special way —

- Within the pair, the observations are dependent (since they are married to one another);
- Between pairs, observations are independent.

#### 4 INDEPENDENT SAMPLES: UNEQUAL VARIANCES

Our interest is to make statistical inference on  $\mu_1 - \mu_2$ . Consider the following assumptions:

##### INDEPENDENT SAMPLES (KNOWN AND UNEQUAL VARIANCES)

1. A random sample of size  $n_1$  from population 1 with mean  $\mu_1$  and variance  $\sigma_1^2$ .
2. A random sample of size  $n_2$  from population 2 with mean  $\mu_2$  and variance  $\sigma_2^2$ .
3. The two samples are **independent**.
4. The population **variances are known** and **not the same**:  $\sigma_1^2 \neq \sigma_2^2$
5. Either one of the following conditions holds:
  - The two populations are **normal**; **OR**
  - Both samples are **large**:  $n_1 \geq 30, n_2 \geq 30$ .

Consider  $X_1, \dots, X_{n_1}$  and  $Y_1, \dots, Y_{n_2}$ , random samples from the two populations of interest. Let

$$\bar{X} = \frac{1}{n_1} \sum_{i=1}^{n_1} X_i, \text{ and } \bar{Y} = \frac{1}{n_2} \sum_{i=1}^{n_2} Y_i$$

be the means of random samples. Then,

$$E(\bar{X}) = \mu_1, \quad V(\bar{X}) = \frac{\sigma_1^2}{n_1}, \quad E(\bar{Y}) = \mu_2, \quad V(\bar{Y}) = \frac{\sigma_2^2}{n_2}.$$

Thus

$$E(\bar{X} - \bar{Y}) = \mu_1 - \mu_2 = \delta,$$

and, using the independence assumption,

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

When

- the two populations are normal, **OR**
- both samples are large,

we have

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \approx N(0, 1).$$

**Confidence Intervals for  $\mu_1 - \mu_2$** 

We are interested in the difference

$$\delta = \mu_1 - \mu_2,$$

with confidence  $100(1 - \alpha)\%$  for any  $0 < \alpha < 1$ .

If  $\sigma_1^2$  and  $\sigma_2^2$  are **known**, by the distributions above, we have

$$P\left(\left|\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right| < z_{\alpha/2}\right) = 1 - \alpha$$

or

$$P\left((\bar{X} - \bar{Y}) - z_{\alpha/2}\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < (\bar{X} - \bar{Y}) + z_{\alpha/2}\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right) = 1 - \alpha.$$

Thus the  $100(1 - \alpha)\%$  confidence interval for  $\mu_1 - \mu_2$  is

$$\left((\bar{X} - \bar{Y}) - z_{\alpha/2}\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X} - \bar{Y}) + z_{\alpha/2}\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right)$$

or

$$(\bar{X} - \bar{Y}) \pm z_{\alpha/2}\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}.$$

**CONFIDENCE INTERVALS: KNOWN AND UNEQUAL VARIANCES**

Suppose we have **independent** populations with **known and unequal variances**, and that either one of the following conditions holds:

- The two populations are **normal**; **OR**
- Both samples are **large**:  $n_1 \geq 30, n_2 \geq 30$ .

The  $100(1 - \alpha)\%$  confidence interval for  $\mu_1 - \mu_2$ , is then given as

$$(\bar{x} - \bar{y}) \pm z_{\alpha/2}\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}.$$

**EXAMPLE 5.10**

A study was conducted to compare two types of engines, A and B.

Gas mileage, in miles per gallon, was measured. 50 experiments were conducted using engine A. 75 experiments were done for engine type B. The gasoline used and other conditions were held constant.

- The average gas mileage for 50 experiments using engine A was 36 miles per gallon and
- The average gas mileage for the 75 experiments using machine B was 42 miles per gallon.

Find a 96% confidence interval on  $\mu_B - \mu_A$ , where  $\mu_A$  and  $\mu_B$  are the population mean gas mileage for machine types A and B, respectively.

Assume that the population standard deviations are 6 and 8 for machine types A and B, respectively.

**Solution:**

For a 96% confidence interval,  $\alpha = 0.04$  and  $z_{0.02} = 2.05$ . We are also given that

$$\begin{aligned} n_1 &= 50, \bar{x}_A = 36, \sigma_1^2 = 6^2 \\ n_2 &= 75, \bar{x}_B = 42, \sigma_2^2 = 8^2 \end{aligned}$$

The sample sizes are large, so a 96% confidence interval for  $\mu_B - \mu_A$  is

$$\begin{aligned} &(\bar{x}_B - \bar{x}_A) \pm z_{\alpha/2} \sqrt{\sigma_2^2/n_2 + \sigma_1^2/n_1} \\ &= (42 - 36) \pm 2.05 \cdot \sqrt{8^2/75 + 6^2/50} \\ &= (3.428, 8.571). \end{aligned}$$

We next consider the following assumptions/case:

**INDEPENDENT SAMPLES (LARGE, WITH UNKNOWN VARIANCES)**

1. A random sample of size  $n_1$  from population 1 with mean  $\mu_1$  and variance  $\sigma_1^2$ .
2. A random sample of size  $n_2$  from population 2 with mean  $\mu_2$  and variance  $\sigma_2^2$ .
3. The two samples are **independent**.
4. The population **variances are unknown** and **not the same**:  $\sigma_1^2 \neq \sigma_2^2$
5. Both samples are **large**:  $n_1 \geq 30, n_2 \geq 30$ .

Since  $\sigma_1$  and  $\sigma_2$  are unknown, let

$$S_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (X_i - \bar{X})^2, \text{ and } S_2^2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (Y_i - \bar{Y})^2$$



and use

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \approx N(0, 1).$$

If  $\sigma_1^2$  and  $\sigma_2^2$  are **unknown**, the  $100(1 - \alpha)\%$  CI for  $\mu_1 - \mu_2$  is

$$\left( (\bar{X} - \bar{Y}) - z_{\alpha/2} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, (\bar{X} - \bar{Y}) + z_{\alpha/2} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right)$$

or

$$(\bar{X} - \bar{Y}) \pm z_{\alpha/2} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}.$$

#### CONFIDENCE INTERVALS: LARGE, WITH UNKNOWN VARIANCES

Suppose we have **independent** populations with **unknown and unequal variances**, and that both samples are **large**:  $n_1 \geq 30, n_2 \geq 30$ .

The  $100(1 - \alpha)\%$  confidence interval for  $\mu_1 - \mu_2$ , is then given as

$$(\bar{x} - \bar{y}) \pm z_{\alpha/2} \sqrt{s_1^2/n_1 + s_2^2/n_2}.$$

## 5 INDEPENDENT SAMPLES: EQUAL VARIANCES

Consider the following assumptions:

#### INDEPENDENT SAMPLES: SMALL, WITH EQUAL VARIANCES

1. A random sample of size  $n_1$  from population 1 with mean  $\mu_1$  and variance  $\sigma_1^2$ .
2. A random sample of size  $n_2$  from population 2 with mean  $\mu_2$  and variance  $\sigma_2^2$ .
3. The two samples are **independent**.
4. The population **variances are unknown** and **the same**:  $\sigma_1^2 = \sigma_2^2 = \sigma^2$ .
5. Both samples are **small**:  $n_1 < 30, n_2 < 30$
6. Both populations are **normally distributed**.

#### THE EQUAL VARIANCE ASSUMPTION

*In real applications, the equal variance assumption is usually unknown and needs to be checked.*

Based upon the normal distribution and equal variance assumptions

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim N(0, 1).$$

Since  $\sigma$  is unknown, we shall estimate it.

Note that  $S_1^2$  and  $S_2^2$  are both unbiased estimators of  $\sigma^2$  under the equal variance assumption.

We can use the **pooled estimator** to estimate  $\sigma^2$  better.

**DEFINITION 6 (THE POOLED ESTIMATOR:  $S_p^2$ )**

$\sigma^2$  can be estimated by the **pooled sample variance**

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

with  $S_1^2$  and  $S_2^2$  being the sample variances of the first and second samples respectively.

When we estimate  $\sigma^2$  using  $S_p^2$ , the resulting statistic

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}$$

follows a  $t$ -distribution with degrees of freedom  $n_1 + n_2 - 2$ .

We then have

$$P \left( -t_{n_1 + n_2 - 2; \alpha/2} < \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} < t_{n_1 + n_2 - 2; \alpha/2} \right) = 1 - \alpha.$$

**CONFIDENCE INTERVALS: SMALL, WITH EQUAL VARIANCES**

Suppose we have **independent, normal** populations with **unknown and equal variances**, and that both samples are **small**:  $n_1 < 30, n_2 < 30$ .

A  $100(1 - \alpha)\%$  confidence interval for  $\mu_1 - \mu_2$  is given as

$$(\bar{X} - \bar{Y}) \pm t_{n_1 + n_2 - 2; \alpha/2} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}.$$

**EXAMPLE 5.11**

A course in mathematics is taught to 12 students by the conventional classroom procedure. A second group of 10 students was given the same course by means of programmed materials.

At the end of the semester the same examination was given to each group.

- The 12 students meeting in the classroom made an average grade of 85 with standard deviation of 4.
- The 10 students using programmed materials made an average of 81 with a standard deviation of 5.

Find a 90% confidence interval for the difference between the population means, assuming the populations are approximately normally distributed with equal variances.

**Solution:**

Let  $\mu_1$  and  $\mu_2$  represent the average grades of all students who might take this course by the classroom and programmed presentations respectively.

So  $\bar{x} - \bar{y} = 85 - 81 = 4$  is the point estimate for  $\mu_1 - \mu_2$ .

As we assume equal population variance, we estimate it by the pooled variance

$$s_p^2 = \frac{(12-1) \times 4^2 + (10-1) \times 5^2}{12+10-2} = 20.05.$$

In this case,  $t_{n_1+n_2-2; \alpha/2} = t_{20; 0.05} = 1.7247$ . Thus a 90% confidence interval for  $\mu_1 - \mu_2$  is given as

$$\begin{aligned} & (\bar{x} - \bar{y}) \pm t_{n_1+n_2-2; \alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \\ &= (85 - 81) \pm 1.7247 \times \sqrt{20.05} \times \sqrt{\frac{1}{12} + \frac{1}{10}} \\ &= (0.693, 7.307). \end{aligned}$$

**Independent Large Samples with Equal Variance**

Note that for large samples such that  $n_1 \geq 30, n_2 \geq 30$ , we can replace  $t_{n_1+n_2-2; \alpha/2}$  by  $z_{\alpha/2}$  in the previous formula.

**CONFIDENCE INTERVALS: LARGE, WITH EQUAL VARIANCES**

Suppose we have **independent** populations with **unknown and equal variances**, and that both samples are **large**:  $n_1 \geq 30, n_2 \geq 30$ .

A  $100(1 - \alpha)\%$  confidence interval for  $\mu_1 - \mu_2$  is given as

$$(\bar{X} - \bar{Y}) \pm z_{\alpha/2} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}.$$

**L-EXAMPLE 5.6 (ELECTRICAL USAGE)**

As a baseline for a study on the effects of changing electrical pricing for electricity during peak hours, July usage during peak hours was obtained for  $n_1 = 45$  homes with air-conditioning and  $n_2 = 55$  homes without. The summarized results are provided below

sample	Size	Mean	Variance
With	45	204.4	13,825.3
Without	55	130.0	8,632.0

Construct a 95% confidence interval for  $\delta = \mu_1 - \mu_2$ .

**Solution:**

For a 95% confidence interval,  $\alpha = 0.05$  and  $z_{0.025} = 1.96$ .

If we assume that  $\sigma_1^2 \neq \sigma_2^2$ , the 95% confidence interval is then constructed via the formula:

$$\begin{aligned} & (\bar{x} - \bar{y}) \pm z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \\ &= (204.4 - 130.0) \pm 1.96 \sqrt{\frac{13825.3}{45} + \frac{8632.0}{55}} \\ &= (32.1724, 116.6276). \end{aligned}$$

If we assume that  $\sigma_1^2 = \sigma_2^2$ , the 95% confidence interval is then constructed via the formula:

$$\begin{aligned} & (\bar{x} - \bar{y}) \pm z_{\alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \\ &= (204.4 - 130) \pm 1.96 \cdot \sqrt{10963.69} \sqrt{\frac{1}{45} + \frac{1}{55}} \\ &= (33.1478, 115.6522). \end{aligned}$$

Here

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} = 10963.69.$$

**REMARK**

We can roughly assume equal variance if

$$1/2 \leq S_1/S_2 \leq 2.$$

This is because the statistic is not overly sensitive to small difference between the population variances.

**6 PAIRED DATA**

Some times, like in the couple income example, it makes sense to take matched pairs instead of independent samples.

Because of dependence in the sample, the methods discussed previously are not applicable.

Consider the assumptions that follows.

**PAIRED DATA**

1.  $(X_1, Y_1), \dots, (X_n, Y_n)$  are matched pairs, where  $X_1, \dots, X_n$  is a random sample from population 1,  $Y_1, \dots, Y_n$  is a random sample from population 2.
2.  $X_i$  and  $Y_i$  are dependent.
3.  $(X_i, Y_i)$  and  $(X_j, Y_j)$  are independent for any  $i \neq j$ .
4. For matched pairs, define  $D_i = X_i - Y_i$ ,  $\mu_D = \mu_1 - \mu_2$ .
5. Now we can treat  $D_1, D_2, \dots, D_n$  as a random sample from a single population with mean  $\mu_D$  and variance  $\sigma_D^2$ .

All techniques derived for a single population can now be employed.

- We consider the statistic

$$T = \frac{\bar{D} - \mu_D}{S_D / \sqrt{n}}, \quad \text{where} \quad \bar{D} = \frac{\sum_{i=1}^n D_i}{n}, \quad S_D^2 = \frac{\sum_{i=1}^n (D_i - \bar{D})^2}{n-1}.$$

- If  $n < 30$  and the population is normally distributed then

$$T \sim t_{n-1}.$$

- If  $n \geq 30$ , then

$$T \sim N(0, 1).$$

**CONFIDENCE INTERVALS: PAIRED DATA**

For **paired data**, if  $n$  is **small** ( $n < 30$ ) and the population is **normally distributed**, a  $(1 - \alpha)100\%$  confidence interval for  $\mu_D$  is

$$\bar{d} \pm t_{n-1; \alpha/2} \cdot \frac{s_D}{\sqrt{n}}.$$

If  $n$  is **large** ( $n \geq 30$ ), a  $(1 - \alpha)100\%$  confidence interval for  $\mu_D$  is

$$\bar{d} \pm z_{\alpha/2} \cdot \frac{s_D}{\sqrt{n}}.$$

**EXAMPLE 5.12**

Twenty students were divided into 10 pairs, each member of the pair having approximately the same IQ.

One of each pair was selected at random and assigned to a mathematics section using programmed materials only. The other member of each pair was assigned to a section in which the professor lectured.

At the end of the semester each group was given the same examination and the following results were recorded.

Pair	1	2	3	4	5	6	7	8	9	10
P.M.	76	60	85	58	91	75	82	64	79	88
Lecture	81	52	87	70	86	77	90	63	85	83
d	-5	8	-2	-12	5	-2	-8	1	-6	5

Given that  $\bar{d} = -1.6$  and  $s_D^2 = 40.71$ , compute a 98% confidence interval for the true difference in the two learning procedures.

**Solution:**

Since  $\alpha = 0.02$ , we have  $t_{n-1; \alpha/2} = t_{9, 0.01} = 2.821$ . Thus a 98% confidence interval for the true difference  $\mu_D$  is given as

$$\bar{d} \pm t_{n-1; \alpha/2} \cdot \frac{s_D}{\sqrt{n}} = -1.6 \pm 2.821 \times \sqrt{\frac{40.71}{10}} = (-7.292, 4.092).$$

**L-EXAMPLE 5.7 (WATER TREATMENT)**

A state law requires municipal waste water treatment plants to monitor their discharges into rivers and streams. A treatment plant could choose to send its samples to a commercial laboratory of its choosing.

Concern over this self-monitoring led a civil engineer to design a matched pairs experiment. Exactly the same bottle of effluent cannot

be sent to two different laboratories. To match “identical” as closely as possible, she would take a sample of effluent in a large sample bottle and pour it back and forth over two open specimen bottles.

When they were filled and capped, a coin was flipped to see if the one on the right was sent to commercial laboratory or the state laboratory.

This process was repeated 11 times. The results, for the response suspended solids (SS) are

Sample	1	2	3	4	5	6	7	8	9	10	11
Commercial lab	27	23	64	44	30	75	26	124	54	30	14
State lab	15	13	22	29	31	64	30	64	56	20	21
Difference $X_i - Y_i$	12	10	42	15	-1	11	-4	60	-2	10	-7

Obtain a 95% confidence interval for the difference in SS from the two labs.

**Solution:**

We assume a normal distribution for the population. From the data, we compute the following

$$n = 11, \bar{d} = 13.27, s_D^2 = 418.61.$$

Further, since  $\alpha = 0.05$ , we have  $t_{n-1;\alpha/2} = t_{10,0.025} = 2.228$ .

The 95% confidence interval is given as

$$13.27 \pm 2.228 \sqrt{\frac{418.61}{11}} = (-0.47, 27.01).$$

