

# Chapter 4: Interval estimation

STAT2602A Probability and statistics II  
(2024-2025 1st Semester)

# Contents

4.1 Basic concepts

4.2 Confidence intervals for means

4.3 Confidence intervals for variances

4.4 Confidence intervals: Large samples

## 4.1 Basic concepts

- ▶ **Motivation:** A point estimator for  $\theta$  does not provide much information about the accuracy of the estimator. It is desirable to generate a narrow interval that will cover the unknown parameter  $\theta$  with a large probability (confidence).
- ▶ **Definition 4.1.** (**Interval estimator**) An interval estimator of  $\theta$  is a random interval  $[L(\mathbf{X}), U(\mathbf{X})]$ , where  $L(\mathbf{X}) := L(X_1, \dots, X_n)$  and  $U(\mathbf{X}) := U(X_1, \dots, X_n)$  are two statistics such that  $L(\mathbf{X}) \leq U(\mathbf{X})$  with probability one.
- ▶ **Definition 4.2.** (**Interval estimate**) If  $\mathbf{X} = \mathbf{x}$  is observed,  $[L(\mathbf{x}), U(\mathbf{x})]$  is the interval estimate of  $\theta$ .
- ▶ **Remark:** *it will sometimes be more natural to use an open interval  $(L(\mathbf{X}), U(\mathbf{X}))$ , a half-open and half-closed interval  $(L(\mathbf{X}), U(\mathbf{X})]$  (or  $[L(\mathbf{X}), U(\mathbf{X}))$ ), or an one-sided interval  $(\infty, U(\mathbf{X})]$  (or  $[L(\mathbf{X}), \infty)$ ).*

## 4.1 Basic concepts

**Example 4.1.** For an independent random sample  $X_1, X_2, X_3, X_4$  from  $N(\mu, 1)$ , consider an interval estimator of  $\mu$  by  $[\bar{X} - 1, \bar{X} + 1]$ . Then, the probability that  $\mu$  is covered by the interval  $[\bar{X} - 1, \bar{X} + 1]$  can be calculated by

$$\begin{aligned}P(\mu \in [\bar{X} - 1, \bar{X} + 1]) &= P(\bar{X} - 1 \leq \mu \leq \bar{X} + 1) \\&= P\left(-2 \leq (\bar{X} - \mu)/\sqrt{1/4} \leq 2\right) \\&= P(-2 \leq Z \leq 2) \\&\approx 0.9544,\end{aligned}$$

where  $Z \sim N(0, 1)$  and with the fact that  $\bar{X} \sim N(\mu, 1/4)$ . Thus, we have over a 95% chance of covering the unknown parameter with our interval estimator. Note that for any **point estimator**  $\hat{\mu}$  of  $\mu$ , we have  $P(\hat{\mu} = \mu) = 0$ . **Sacrificing some precision in the interval estimator**, in moving from a point to an interval, has resulted in increased confidence that our assertion about  $\mu$  is correct. □

## 4.1 Basic concepts

- **Definition 4.3.** (*Confidence coefficient*) For an interval estimator  $[L(\mathbf{X}), U(\mathbf{X})]$  of  $\theta$ , the confidence coefficient of  $[L(\mathbf{X}), U(\mathbf{X})]$ , denoted by  $(1 - \alpha)$ , is

$$1 - \alpha = P(\theta \in [L(\mathbf{X}), U(\mathbf{X})]),$$

where  $P(\theta \in [L(\mathbf{X}), U(\mathbf{X})])$  is the coverage probability of  $[L(\mathbf{X}), U(\mathbf{X})]$ .

- **Remark:** In some situations, the coverage probability  $P(\theta \in [L(\mathbf{X}), U(\mathbf{X})])$  may depend on  $\theta$ , and then the confidence coefficient is defined as

$$1 - \alpha = \inf_{\theta} P(\theta \in [L(\mathbf{X}), U(\mathbf{X})]).$$

- Interval estimator, together with a measure of confidence (say, the confidence coefficient), is sometimes known as **confidence interval**. A confidence interval with confidence coefficient equal to  $1 - \alpha$ , is called a  **$1 - \alpha$  confidence interval**.

## 4.1 Basic concepts

- ▶ **Definition 4.4.** (*Pivotal Quantity*) A random variable  $Q(\mathbf{X}, \theta) = Q(X_1, \dots, X_n, \theta)$  is a pivotal quantity if the distribution of  $Q(\mathbf{X}, \theta)$  is free of  $\theta$ . That is, regardless of the distribution of  $\mathbf{X}$ ,  $Q(\mathbf{X}, \theta)$  has the same distribution for all values of  $\theta$ .
- ▶ **Remark:** Logically, when  $Q(\mathbf{X}, \theta)$  is a pivotal quantity, we can easily construct a  $1 - \alpha$  confidence interval for  $Q(\mathbf{X}, \theta)$  by

$$1 - \alpha = P \left( \tilde{L}_\alpha \leq Q(\mathbf{X}, \theta) \leq \tilde{U}_\alpha \right), \quad (1.1)$$

where  $\tilde{L}_\alpha$  and  $\tilde{U}_\alpha$  do not depend on  $\theta$ . Suppose that the inequalities  $\tilde{L}_\alpha \leq Q(\mathbf{X}, \theta) \leq \tilde{U}_\alpha$  in (1.1) are equivalent to the inequalities  $L(\mathbf{X}) \leq \theta \leq U(\mathbf{X})$ . Then, from (1.1), a  $1 - \alpha$  confidence interval of  $\theta$  is  $[L(\mathbf{X}), U(\mathbf{X})]$ .

## 4.2 Confidence intervals for means - One-sample case

**Formula derivation** Let  $\mathbf{X} = \{X_1, \dots, X_n\}$  be an independent random sample from the population  $N(\mu, \sigma^2)$ . We first consider the interval estimator of  $\mu$  when  $\sigma^2$  is known. Note that

$$\bar{X} \sim N(\mu, \sigma^2/n). \quad (2.1)$$

Hence, when  $\sigma^2$  is known,  $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$  is a pivotal quantity involving  $\mu$ . Let

$$\begin{aligned} 1 - \alpha &= P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) \\ &= P\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) \\ &= P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right), \end{aligned}$$

where  $z_\alpha$  satisfies

$$P(Z \geq z_\alpha) = \alpha$$

for  $Z \sim N(0, 1)$

## 4.2 Confidence intervals for means - One-sample case

**Formula derivation (con't)** Usually, we call  $z_\alpha$  the upper percentile of  $N(0, 1)$  at the level  $\alpha$ . When  $\sigma^2$  is known, a  $1 - \alpha$  confidence interval of  $\mu$  is

$$\left[ \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]. \quad (2.2)$$

Given the observed value of  $\bar{X} = \bar{x}$  and the value of  $z_{\alpha/2}$ , we can calculate the interval estimate of  $\mu$  by

$$\left[ \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right].$$

As the point estimator, the  $1 - \alpha$  confidence interval is also not unique. Ideally, we should choose it as narrow as possible in some sense, but in practice, we usually choose the equal-tail confidence interval as in (2.2) for convenience, since tables for selecting equal probabilities in the two tails are readily available.



## 4.2 Confidence intervals for means - One-sample case

*Example 4.2.* A publishing company has just published a new college textbook. Before the company decides the price of the book, it wants to know the average price of all such textbooks in the market. The research department at the company took a sample of 36 such textbooks and collected information on their prices. This information produced a mean price of \$48.40 for this sample. It is known that the standard deviation of the prices of all such textbooks is \$4.50. Construct a 90% confidence interval for the mean price of all such college textbooks assuming that the underlying population is normal.

*Solution.* From the given information,  $n = 36$ ,  $\bar{x} = 48.40$  and  $\sigma = 4.50$ . Now,  $1 - \alpha = 0.9$ , i.e.,  $\alpha = 0.1$ , and by (2.2), the 90% confidence interval for the mean price of all such college textbooks is given by

$$\begin{aligned} \left[ \bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right] &= \left[ 48.40 - z_{0.05} \frac{4.50}{\sqrt{36}}, 48.40 + z_{0.05} \frac{4.50}{\sqrt{36}} \right] \\ &\approx [47.1662, 49.6338]. \end{aligned}$$

## 4.2 Confidence intervals for means - One-sample case

**Example 4.3.** Suppose the bureau of the census and statistics of a city wants to estimate the mean family annual income  $\mu$  for all families in the city. It is known that the standard deviation  $\sigma$  for the family annual income is 60 thousand dollars. How large a sample should the bureau select so that it can assert with probability 0.99 that the sample mean will differ from  $\mu$  by no more than 5 thousand dollars?

**Solution.** From the construction of a confidence interval, we have

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

where  $1 - \alpha = 0.99$  and  $\sigma = 60$  thousand dollars. It suffices to have  $z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq 5$  or

$$n \geq (60z_{\alpha/2}/5)^2 = (60 \times 2.576/5)^2 \approx 955.5517.$$

Thus, the sample size should be at least 956. (Note that we have to round 955.5517 up to the next higher integer. This is always the case when determining the sample size.)

## 4.2 Confidence intervals for means - One-sample case

**Case: consider the interval estimator of  $\mu$  when  $\sigma^2$  is unknown.**

- ▶ **Property 4.1.** (i)  $\bar{X}$  and  $S^2$  are independent;
- (ii)  $\frac{nS^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}$  is  $\chi_{n-1}^2$ , where  $\chi_k^2$  is a chi-square distribution with  $k$  degrees of freedom;
- (iii)  $T = \frac{\bar{X} - \mu}{S/\sqrt{n-1}}$  is  $t_{n-1}$ , where  $t_k$  is a t distribution with  $k$  degrees of freedom.
- ▶ **Property 4.2.** (i) If  $Z_1, \dots, Z_k$  are  $k$  independent  $N(0, 1)$  random variables, then  $\sum_{i=1}^k Z_i^2$  is  $\chi_k^2$ ;
- (ii) If  $Z$  is  $N(0, 1)$ ,  $U$  is  $\chi_k^2$ , and  $Z$  and  $U$  are independent, then  $T = \frac{Z}{\sqrt{U/k}}$  is  $t_k$ .

*The proof of Property 4.1. could be found in the lecture notes.*

## 4.2 Confidence intervals for means - One-sample case

**Derivation of the interval estimator of  $\mu$  when  $\sigma^2$  is unknown.**

From Property 4.1(iii), we know that  $T$  is a pivotal quantity of  $\mu$ .

Let

$$\begin{aligned}1 - \alpha &= P\left(-t_{\alpha/2, df=n-1} \leq T \leq t_{\alpha/2, df=n-1}\right) \\&= P\left(-t_{\alpha/2, df=n-1} \leq \frac{\bar{X} - \mu}{S/\sqrt{n-1}} \leq t_{\alpha/2, df=n-1}\right) \\&= P\left(\bar{X} - t_{\alpha/2, df=n-1} \frac{S}{\sqrt{n-1}} \leq \mu \leq \bar{X} + t_{\alpha/2, df=n-1} \frac{S}{\sqrt{n-1}}\right)\end{aligned}$$

where  $t_{\alpha, df=k}$  satisfies

$$P(T \geq t_{\alpha, df=k}) = \alpha$$

for a random variable  $T \sim t_k$ .

## 4.2 Confidence intervals for means - One-sample case

### Derivation con't

Therefore, when  $\sigma^2$  is unknown, a  $1 - \alpha$  confidence interval of  $\mu$  is

$$\left[ \bar{X} - t_{\alpha/2, df=n-1} \frac{S}{\sqrt{n-1}}, \bar{X} + t_{\alpha/2, df=n-1} \frac{S}{\sqrt{n-1}} \right]. \quad (2.3)$$

Given the observed value of  $\bar{X} = \bar{x}$ ,  $S = s$ , and the value of  $t_{\alpha/2, df=n-1}$ , we can calculate the interval estimate of  $\mu$  by

$$\left[ \bar{x} - t_{\alpha/2, df=n-1} \frac{s}{\sqrt{n-1}}, \bar{x} + t_{\alpha/2, df=n-1} \frac{s}{\sqrt{n-1}} \right].$$

- **Remark:** Usually there is a row with  $\infty$  degrees of freedom in a  $t$ -distribution table, which actually shows values of  $z_{\alpha}$ . In fact, when  $n \rightarrow \infty$ , the distribution function of  $t_n$  tends to that of  $N(0, 1)$ . That is, in tests or exams, if  $n$  is so large that the value of  $t_{\alpha, df=n}$  cannot be found, you may use  $z_{\alpha}$  instead.

## 4.2 Confidence intervals for means - One-sample case

**Example 4.4** A paint manufacturer wants to determine the average drying time of a new brand of interior wall paint. If for 12 test areas of equal size he obtained a mean drying time of 66.3 minutes and a standard deviation of 8.4 minutes, construct a 95% confidence interval for the true population mean assuming normality.

**Solution.** As  $n = 12$ ,  $\bar{x} = 66.3$ ,  $s = 8.4$ ,  $\alpha = 1 - 0.95 = 0.05$  and  $t_{\alpha/2, df=n-1} = t_{0.025, 11} \approx 2.201$ , the 95% confidence interval for  $\mu$  is

$$\left[ 66.3 - 2.201 \times \frac{8.4}{\sqrt{12-1}}, 66.3 + 2.201 \times \frac{8.4}{\sqrt{12-1}} \right],$$

that is, [61.1722, 71.4278].



## 4.2 Confidence intervals for means - One-sample case

*Example 4.5* Construct a 95% confidence interval for the mean hourly wage of apprentice geologists employed by the top 5 oil companies. For a sample of 50 apprentice geologists,  $\bar{x} = 14.75$  and  $s = 3.0$  (in dollars).

*Solution.* As  $n = 50$ ,  $\alpha = 1 - 0.95 = 0.05$ ,  $t_{0.025, df=49} \approx 2.010$ , we have

$$t_{\alpha/2, df=n-1} \frac{s}{\sqrt{n-1}} \approx 2.010 \times \frac{3.0}{\sqrt{50-1}} = 0.8614.$$

Thus, the 95% confidence interval is  $[14.75 - 0.86, 14.75 + 0.86]$ , or  $[13.89, 15.59]$ . □

## 4.2 Confidence intervals for means - Two-sample case

In the next part, we shall consider the problem of constructing confidence intervals for the difference of the means of two normal distributions when the variances are unknown.

### Formula derivation

Let  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$  and  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_m\}$  be random samples from independent distributions  $N(\mu_X, \sigma_X^2)$  and  $N(\mu_Y, \sigma_Y^2)$ , respectively. We are of interest to construct the confidence interval for  $\mu_X - \mu_Y$  when  $\sigma_X^2 = \sigma_Y^2 = \sigma^2$ .

First, we can show that

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\sigma^2/n + \sigma^2/m}}$$

is  $N(0, 1)$ . Also, by the independence of  $\mathbf{X}$  and  $\mathbf{Y}$ , from Property 4.1(ii), we know that

$$U = \frac{nS_X^2}{\sigma^2} + \frac{mS_Y^2}{\sigma^2}$$

is  $\chi_{n+m-2}^2$ .



## 4.2 Confidence intervals for means - Two-sample case

### Formula derivation con't

Moreover, by Property 4.1(i),  $Z$  and  $U$  are independent. Hence,

$$\begin{aligned} T &= \frac{Z}{\sqrt{U/(n+m-2)}} \\ &= \frac{[(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)] / \sqrt{\sigma^2/n + \sigma^2/m}}{\sqrt{(nS_X^2 + mS_Y^2)/[\sigma^2(n+m-2)]}} \\ &= \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{R} \end{aligned}$$

is  $t_{n+m-2}$ , where

$$R = \sqrt{\frac{nS_X^2 + mS_Y^2}{n+m-2} \left( \frac{1}{n} + \frac{1}{m} \right)}.$$

That is,  $T$  is a pivotal quantity of  $\mu_X - \mu_Y$ .

## 4.2 Confidence intervals for means - Two-sample case

### Formula derivation con't

Let

$$\begin{aligned}1 - \alpha &= P\left(-t_{\alpha/2, df=n+m-2} \leq T \leq t_{\alpha/2, df=n+m-2}\right) \\&= P\left(-t \leq \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{R} \leq t\right) \\&= P\left((\bar{X} - \bar{Y}) - tR \leq \mu_X - \mu_Y \leq (\bar{X} - \bar{Y}) + tR\right),\end{aligned}$$

So, when  $\sigma_X^2 = \sigma_Y^2 = \sigma^2$  is unknown, a  $1 - \alpha$  confidence interval of  $\mu_X - \mu_Y$  is

$$\left[(\bar{X} - \bar{Y}) - t_{\alpha/2, df=n+m-2}R, (\bar{X} - \bar{Y}) + t_{\alpha/2, df=n+m-2}R\right]. \quad (2.4)$$

## 4.2 Confidence intervals for means - Two-sample case

### Formula derivation con't

Given the observed value of  $\bar{X} = \bar{x}$ ,  $\bar{Y} = \bar{y}$ ,  $S_X = s_X$ ,  $S_Y = s_Y$ , and the value of  $t_{\alpha/2, df=n+m-2}$ , we can calculate the interval estimate of  $\mu_X - \mu_Y$  by

$$\left[ (\bar{x} - \bar{y}) - t_{\alpha/2, df=n+m-2} r, (\bar{x} - \bar{y}) + t_{\alpha/2, df=n+m-2} r \right].$$

where

$$r = \sqrt{\frac{ns_X^2 + ms_Y^2}{n+m-2} \left( \frac{1}{n} + \frac{1}{m} \right)}.$$

## 4.2 Confidence intervals for means - Two-sample case

**Example 4.6** Suppose that scores on a standardized test in mathematics taken by students from large and small high schools are  $N(\mu_X, \sigma^2)$  and  $N(\mu_Y, \sigma^2)$ , respectively, where  $\sigma^2$  is unknown. If a random sample of  $n = 9$  students from large high schools yielded  $\bar{x} = 81.31$ ,  $s_X^2 = 60.76$  and a random sample of  $m = 15$  students from small high schools yielded  $\bar{y} = 78.61$ ,  $s_Y^2 = 48.24$ , the endpoints for a 95% confidence interval for  $\mu_X - \mu_Y$  are given by

$$81.31 - 78.61 \pm 2.074 \sqrt{\frac{9 \times 60.76 + 15 \times 48.24}{22} \left( \frac{1}{9} + \frac{1}{15} \right)},$$

since  $P(T \leq 2.074) = 0.975$ . So, the 95% confidence interval is  $[-3.95, 9.35]$ . □

## 4.3 Confidence intervals for variances - One-sample case

**Formula derivation** First, we consider the one-sample case. By Property 4.1(ii),

$$\frac{nS^2}{\sigma^2} \sim \chi_{n-1}^2$$

is a pivotal quantity involving  $\sigma^2$ . Let

$$\begin{aligned} 1 - \alpha &= P \left( \chi_{1-\alpha/2, df=n-1}^2 \leq \frac{nS^2}{\sigma^2} \leq \chi_{\alpha/2, df=n-1}^2 \right) \\ &= P \left( \frac{nS^2}{\chi_{\alpha/2, df=n-1}^2} \leq \sigma^2 \leq \frac{nS^2}{\chi_{1-\alpha/2, df=n-1}^2} \right), \end{aligned}$$

where  $\chi_{\alpha, df=n}^2$  satisfies

$$P(T \geq \chi_{\alpha, df=n}^2) = \alpha$$

for a random variable  $T \sim \chi_n^2$ .

## 4.3 Confidence intervals for variances - One-sample case

### Formula derivation con't

Therefore, a  $1 - \alpha$  confidence interval of  $\sigma^2$  is

$$\left[ \frac{nS^2}{\chi_{\alpha/2, df=n-1}^2}, \frac{nS^2}{\chi_{1-\alpha/2, df=n-1}^2} \right]. \quad (3.1)$$

Given the observed value of  $S = s$  and the values of  $\chi_{\alpha/2, df=n-1}^2$  and  $\chi_{1-\alpha/2, df=n-1}^2$ , we can calculate the interval estimate of  $\sigma^2$  by

$$\left[ \frac{ns^2}{\chi_{\alpha/2, df=n-1}^2}, \frac{ns^2}{\chi_{1-\alpha/2, df=n-1}^2} \right].$$

## 4.3 Confidence intervals for variances - One-sample case

**Example 4.7.** A machine is set up to fill packages of cookies. A recently taken random sample of the weights of 25 packages from the production line gave a variance of  $2.9 \text{ g}^2$ . Construct a 95% confidence interval for the standard deviation of the weight of a randomly selected package from the production line.

*Solution.* As  $n = 25$ ,  $s^2 = 2.9$ ,  $\alpha = 0.05$ ,

$$\frac{ns^2}{\chi_{\alpha/2, df=n-1}^2} = \frac{25(2.9)}{\chi_{0.025, df=24}^2} \approx \frac{25(2.9)}{39.36} \approx 1.8420,$$
$$\frac{ns^2}{\chi_{1-\alpha/2, df=n-1}^2} = \frac{25(2.9)}{\chi_{0.975, df=24}^2} \approx \frac{25(2.9)}{12.40} \approx 5.8468,$$

the 95% confidence interval for the population variance is  $(1.8420, 5.8468)$ . Taking positive square roots, we obtain the 95% confidence interval for the population standard deviation to be  $(1.3572, 2.4180)$ . □

## 4.3 Confidence intervals for variances - Two-sample case

- **Basic settings:** Consider the two-sample case. Let

$$\mathbf{X} = \{X_1, X_2, \dots, X_n\} \text{ and } \mathbf{Y} = \{Y_1, Y_2, \dots, Y_m\}$$

be random samples from independent distributions  $N(\mu_X, \sigma_X^2)$  and  $N(\mu_Y, \sigma_Y^2)$ , respectively. We are of interest to construct the confidence interval for  $\sigma_X^2/\sigma_Y^2$ .

- **Property 4.3.** Suppose that  $U \sim \chi_{r_1}^2$  and  $V \sim \chi_{r_2}^2$  are independent. Then,

$$F_{r_1, r_2} = \frac{U/r_1}{V/r_2}$$

has an  $F_{r_1, r_2}$  distribution with  $r_1$  and  $r_2$  degrees of freedom.



## 4.3 Confidence intervals for variances - Two-sample case

### Formula derivation

By Property 4.1(ii),

$$\frac{nS_X^2}{\sigma_X^2} \sim \chi_{n-1}^2 \quad \text{and} \quad \frac{mS_Y^2}{\sigma_Y^2} \sim \chi_{m-1}^2.$$

Then, by Property 4.3, it follows that

$$\left[ \frac{mS_Y^2}{\sigma_Y^2(m-1)} \right] / \left[ \frac{nS_X^2}{\sigma_X^2(n-1)} \right] \sim F_{m-1, n-1},$$

which is a pivotal quantity involving  $\sigma_X^2/\sigma_Y^2$ .

## 4.3 Confidence intervals for variances - Two-sample case

### Formula derivation con't

Let

$$\begin{aligned} 1 - \alpha &= P \left( F_{1-\alpha/2, df=(m-1, n-1)} \leq \frac{\left[ \frac{mS_Y^2}{\sigma_Y^2(m-1)} \right]}{\left[ \frac{nS_X^2}{\sigma_X^2(n-1)} \right]} \leq F_{\alpha/2, df=(m-1, n-1)} \right) \\ &= P \left( \frac{n(m-1)S_X^2}{m(n-1)S_Y^2} F \leq \frac{\sigma_X^2}{\sigma_Y^2} \leq \frac{n(m-1)S_X^2}{m(n-1)S_Y^2} F \right), \end{aligned}$$

where  $F_{\alpha, df=(m,n)}$  satisfies

$$P(T \geq F_{\alpha, df=(m,n)}) = \alpha$$

for a random variable  $T \sim F_{m,n}$ .

## 4.3 Confidence intervals for variances - Two-sample case

### Formula derivation con't

Therefore, a  $1 - \alpha$  confidence interval of  $\sigma_X^2/\sigma_Y^2$  is

$$\left[ \frac{n(m-1)S_X^2}{m(n-1)S_Y^2} F_{1-\alpha/2, df=(m-1, n-1)}, \frac{n(m-1)S_X^2}{m(n-1)S_Y^2} F_{\alpha/2, df=(m-1, n-1)} \right]. \quad (3.2)$$

Given the observed value of  $S_X = s_X$ ,  $S_Y = s_Y$ , and the values of  $F_{\alpha/2, df=(m-1, n-1)}$  and  $F_{1-\alpha/2, df=(m-1, n-1)}$ , we can calculate the interval estimate of  $\sigma_X^2/\sigma_Y^2$  by

$$\left[ \frac{n(m-1)s_X^2}{m(n-1)s_Y^2} F_{1-\alpha/2, df=(m-1, n-1)}, \frac{n(m-1)s_X^2}{m(n-1)s_Y^2} F_{\alpha/2, df=(m-1, n-1)} \right].$$

## 4.4 Confidence intervals: Large samples

- ▶ **Motivation:** When the population is not normal, we can make use of the CLT to propose confidence intervals, which has an approximated confidence coefficient  $1 - \alpha$  for large  $n$ .
- ▶ **Theorem 4.1.** (Slutsky's theorem) If  $X_n \rightarrow_d X$  and  $Y_n \rightarrow_p C$  (a constant), then
  - (i)  $X_n + Y_n \rightarrow_d X + C$ ;
  - (ii)  $X_n \cdot Y_n \rightarrow_d X \cdot C$ ;
  - (iii)  $X_n/Y_n \rightarrow_d X/C$  provided that  $C \neq 0$ .

*Remark:* Note that Theorem 4.1 fails if  $C$  is not a constant.

## 4.4 Confidence intervals: Large samples

### Formula derivation

Let  $\mathbf{X}$  be an independent random sample from a population, which has the mean  $\mu$  and the variance  $\sigma^2 < \infty$ . According to Theorem 4.1, CLT and the fact that  $S \rightarrow_p \sigma$ , we have

$$\frac{\sqrt{n}(\bar{X} - \mu)}{S} = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \cdot \frac{S}{\sigma} \rightarrow_d N(0, 1), \quad (4.1)$$

for large  $n$ . Hence, by (4.1), it follows that for large  $n$ ,

$$\begin{aligned} 1 - \alpha &\approx P\left(-z_{\alpha/2} \leq \frac{\sqrt{n}(\bar{X} - \mu)}{S} \leq z_{\alpha/2}\right) \\ &= P\left(\bar{X} - z_{\alpha/2} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{S}{\sqrt{n}}\right). \end{aligned}$$

## 4.4 Confidence intervals: Large samples

### Formula derivation con't

So, an approximated  $1 - \alpha$  confidence interval of  $\mu$  is

$$\left[ \bar{X} - z_{\alpha/2} \frac{S}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{S}{\sqrt{n}} \right]. \quad (4.2)$$

Given the observed value of  $\bar{x}$  and  $s$ , we can calculate the interval estimate of  $\mu$  by

$$\left[ \bar{x} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right].$$

Note that the confidence interval in (4.2) only requires a large  $n$  but not the normal population assumption. Clearly, the similar idea can be applied to the two-sample case.

## 4.4 Confidence intervals: Large samples

### End of the Chapter

To end this chapter, we consider the interval estimator for percentage  $p$ , where

$$p = P(X \in (a, b)).$$

Define  $\xi = I(a < X < b)$ . Then,  $E(\xi) = p$ . This indicates that  $p$  is the theoretical mean of  $\xi$ . Hence, by (4.2), an approximated  $1 - \alpha$  confidence interval of  $p$  is

$$\left[ \bar{\xi} - z_{\alpha/2} \frac{S_{\xi}}{\sqrt{n}}, \bar{\xi} + z_{\alpha/2} \frac{S_{\xi}}{\sqrt{n}} \right], \quad (4.3)$$

where  $\bar{\xi} = \frac{1}{n} \sum_{i=1}^n \xi_i$  and  $S_{\xi}^2 = \frac{1}{n} \sum_{i=1}^n (\xi_i - \bar{\xi})^2 = \bar{\xi}(1 - \bar{\xi})$  with  $\xi_i = I(a < X_i < b)$ .

In general, we can treat the interval  $(a, b)$  as “success”,  $p = P(\text{“success”})$ , and  $\bar{\xi}$  = relative frequency of “success”.

## 4.4 Confidence intervals: Large samples

*Example 4.8.* In a certain political campaign, one candidate has a poll taken at random among the voting population. The results are  $n = 112$  and  $y = 59$  (for “Yes”). Should the candidate feel very confident of winning?

*Solution.* Let  $p = P(\text{“the candidate wins the campaign”})$ . Then,  $\bar{\xi} = 59/112 \approx 0.527$ . According to (4.3), since  $z_{0.025} \approx 1.96$ , an approximated 95% confident interval estimate for  $p$  is

$$\left[ 0.527 - z \sqrt{\frac{0.527 * (1 - 0.527)}{112}}, 0.527 + z \sqrt{\frac{0.527 * (1 - 0.527)}{112}} \right] \\ \approx [0.435, 0.619].$$

There has certain possibility that  $p$  is less than 50%, and the candidate should take this into account in campaigning. □