

Mathematical Statistics II

Ch.7 Interval Estimation

Jungsoon Choi

jungsoonchoi@hanyang.ac.kr

Table of Contents

- Confidence Intervals for Means
- Confidence Intervals for the Difference of Two Means
- Confidence Intervals for Proportions

Introduction - Interval Estimation

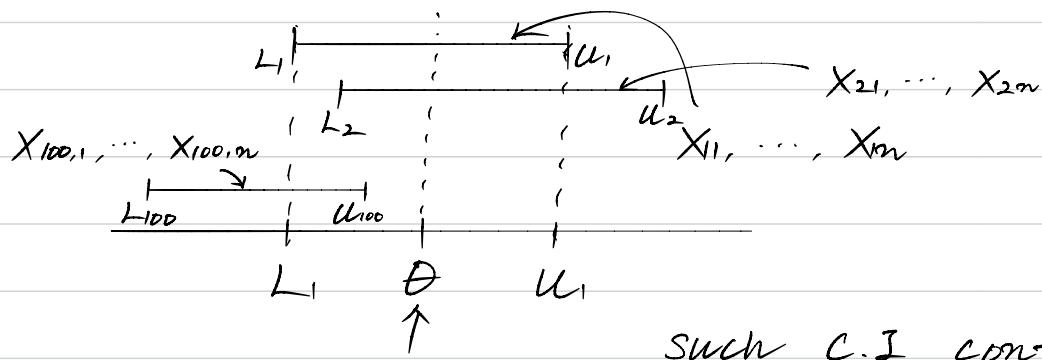
- **Confidence interval:** an interval of values that is likely to contain the true value of the parameter.
- A confidence interval for θ is an interval $[L, U]$ computed from the sample observations X_1, \dots, X_n .
- A $(1 - \alpha) \times 100\%$ confidence interval for θ is given by

$$P(L \leq \theta \leq U) = 1 - \alpha \quad \begin{matrix} \text{def of C.I} \\ \text{in Statistics} \end{matrix}$$

- L : lower confidence limit
- U : upper confidence limit
- $1 - \alpha$: confidence coefficient

X_1, \dots, X_{2n}
 X_1, \dots, X_m
 $(1-\alpha) \times 100\% \text{ C.I. for } \theta$
 POP
 θ
 unknown fixed value $P(L \leq \theta \leq U) = 1 - \alpha$
 $[L, U]$

95% C.I. for θ



such C.I. contain parameter θ or not.

Prob. of C.I. contain $\theta \neq 95\%$.
 확률은 95%가 아님

95% of C.I. contain θ

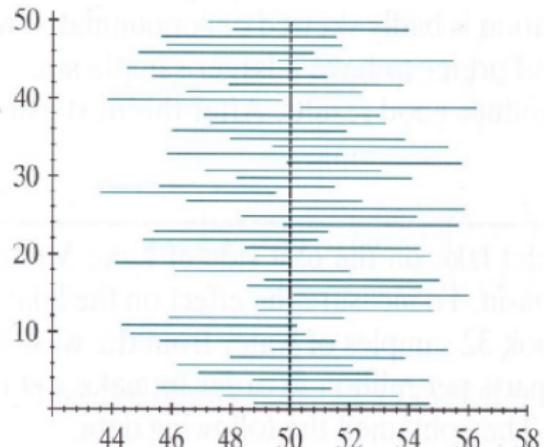
5% " not contain θ .

75% of C.I.

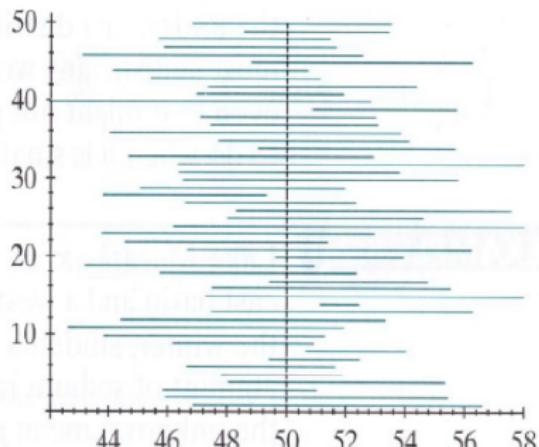
In Bayesian statistic,

θ not fixed, can move. C.I.

thus 95% C.I. means that contain θ with 95% P.



(a) 90% confidence intervals, σ known



(b) 90% confidence intervals, σ unknown

Confidence intervals using z and t

Ch7.1 Confidence Intervals for Means

Confidence Intervals for the population mean μ

Case 1: Normal with the population variance σ^2 known

Suppose that $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where σ^2 is known. The sample mean \bar{X} has a normal distribution, $N\left(\mu, \frac{\sigma^2}{n}\right)$. Then, the standardized random variable $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$.

$$P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha$$

POP
 $N(\mu, \sigma^2)$

$$X_i \sim N(\mu, \sigma^2) \quad \sigma^2 \text{ known} \quad i=1, \dots, n, \text{iid}$$

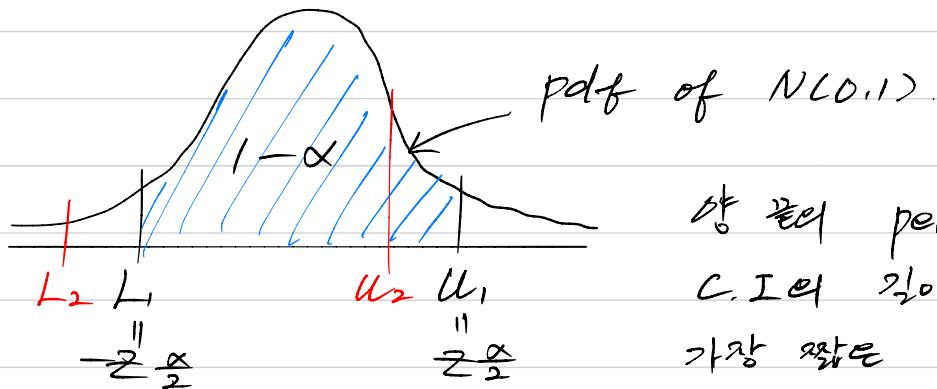
$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \quad \text{cf.) } (1-\alpha) \times 100\% \text{ C.I. for } \theta$$

$(1-\alpha) \times 100\% \text{ C.I. for } \mu$

$$P(L \leq \theta \leq U) = 1 - \alpha$$

$$P(-z_{\frac{\alpha}{2}} \leq Z \leq z_{\frac{\alpha}{2}}) = 1 - \alpha$$

$$P(-z_{\frac{\alpha}{2}} \leq \bar{X} - \mu \leq z_{\frac{\alpha}{2}}) = 1 - \alpha$$



각 절단 percentile의 차이를 통해
C.I.의 폭이 가장 좁아짐.

가장 좁은 C.I.가 더 좋음.

$$P(-z_{\frac{\alpha}{2}} \leq Z \leq z_{\frac{\alpha}{2}}) = 1 - \alpha$$

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

$$\Leftrightarrow P(-z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \bar{X} - \mu \leq z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

$$\Leftrightarrow P(-\bar{X} - z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

$$\Leftrightarrow P(\bar{X} - z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

시작

시작 대답해주세요

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

$(1 - \alpha) \times 100\%$ confidence interval for μ with known σ^2

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

Case 2: Normal with the population variance σ^2 unknown

Suppose that $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ from the observations X_1, \dots, X_n .

Since $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{(n-1)}^2$,

$$T = \frac{\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{(n-1)S^2}{\sigma^2}/(n-1)}} = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{(n-1)}$$

A $(1 - \alpha) \times 100\%$ confidence interval for μ with unknown σ^2 is given by

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

POP

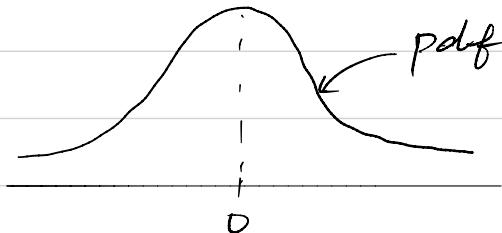
$$X_i \sim N(\mu, \sigma^2) \quad \sigma^2: \text{unknown}$$

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

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

Unknown

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1).$$



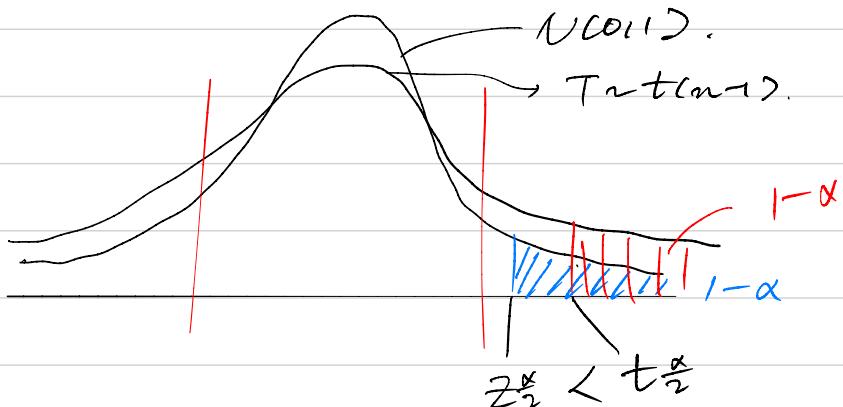
(1- α) × 100% C.I. for μ

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

$$\Leftrightarrow P(-t_{\frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq t_{\frac{\alpha}{2}}) = 1-\alpha.$$

$$\Leftrightarrow P(\bar{X} - t_{\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{\frac{\alpha}{2}} \cdot \frac{S}{\sqrt{n}}) = 1-\alpha$$

2. $Z \frac{\sigma}{\sqrt{n}}$ ~~vs~~ $2 \cdot t_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}$



Note

We know that $t_{\alpha/2, (n-1)} > z_{\alpha/2}$ in the tail parts. If $\sigma = s$, then the length of the confidence interval $\bar{x} \pm z_{\alpha/2}\sigma/\sqrt{n}$ is shorter than the length of the confidence interval $\bar{x} \pm t_{\alpha/2, (n-1)}s/\sqrt{n}$.

Case 3: the population variance σ^2 known and large n

Suppose that $X_1, \dots, X_n \sim (\mu, \sigma^2)$, where σ^2 is known and n is large enough ($n > 30$). Then, the standardized random variable $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ is approximate to $N(0, 1)$.

A $(1 - \alpha) \times 100\%$ confidence interval for μ with known σ^2 is given by

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

Case 4: the population variance σ^2 unknown and large n

Suppose that $X_1, \dots, X_n \sim (\mu, \sigma^2)$, where n is large enough. Then, the standardized random variable $Z = \frac{\bar{X} - \mu}{S/\sqrt{n}}$ (S^2 : sample variance) is approximate to $N(0, 1)$.

A $(1 - \alpha) \times 100\%$ confidence interval for μ with unknown σ^2 is given by

$$\begin{aligned} & \left[\bar{X} - t_{\alpha/2, (n-1)} \frac{S}{\sqrt{n}}, \bar{X} + t_{\alpha/2, (n-1)} \frac{S}{\sqrt{n}} \right] \\ & \approx \left[\bar{X} - z_{\alpha/2} \frac{S}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{S}{\sqrt{n}} \right] \end{aligned}$$

POP (μ, σ^2) $\rightarrow X_i \sim (\mu, \sigma^2)$. σ^2 : known.
 By CLT n : large.

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

$(1-\alpha) \times 100\%$ approximate for μ
 $(\bar{X} - Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}})$

$X_i \sim (\mu, \sigma^2)$. σ^2 : unknown
 n : large.

By CLT, approximation & replacement

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \longrightarrow N(0, 1).$$

$$P(-Z_{\frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{S/\sqrt{n}} \leq Z_{\frac{\alpha}{2}}) = 1 - \alpha$$

One-sided confidence interval

Suppose that $\bar{X} \sim N(\mu, \sigma^2/n)$ where σ^2 is known.

A $(1 - \alpha) \times 100\%$ one-sided lower confidence interval for μ is

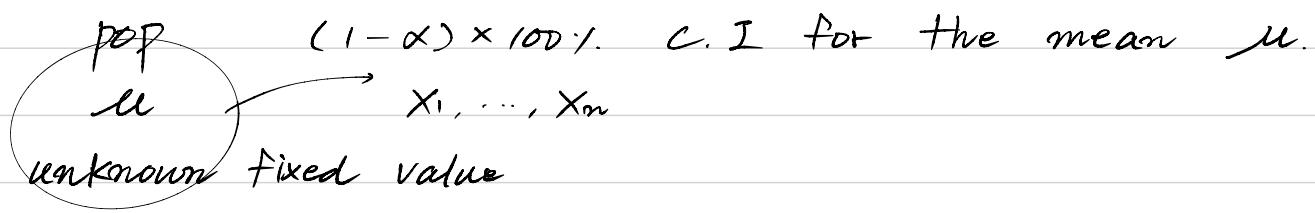
$$P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_\alpha\right) = 1 - \alpha$$

$$\begin{aligned} P\left(\bar{X} - z_\alpha \frac{\sigma}{\sqrt{n}} \leq \mu\right) &= 1 - \alpha \\ \left[\bar{x} - z_\alpha \frac{\sigma}{\sqrt{n}}, \infty\right) \end{aligned}$$

Similarly, a $(1 - \alpha) \times 100\%$ one-sided upper confidence interval for μ is

$$\left(-\infty, \bar{x} + z_\alpha \frac{\sigma}{\sqrt{n}}\right]$$





1) $X_i \sim N(\mu, \sigma^2)$ σ^2 : known.

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \Rightarrow Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

$$P(-Z_{\frac{\alpha}{2}} \leq Z \leq Z_{\frac{\alpha}{2}}) = 1 - \alpha$$

$$\Leftrightarrow P(\bar{X} - Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

2) $X_i \sim N(\mu, \sigma^2)$ σ^2 : unknown.

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$$

$$P(-t_{\frac{\alpha}{2}, (n-1)} \leq T \leq t_{\frac{\alpha}{2}, (n-1)}) = 1 - \alpha$$

$$\Leftrightarrow P(\bar{X} - t_{\frac{\alpha}{2}, (n-1)} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{\frac{\alpha}{2}, (n-1)} \frac{S}{\sqrt{n}}) = 1 - \alpha$$

3) $X_i \sim (\mu, \sigma^2)$ σ^2 : known. large n

By CLT. $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ $(\bar{X} - Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}})$

4) $X_i \sim (\mu, \sigma^2)$, σ^2 : unknown. large n

$[\bar{X} - Z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}, \bar{X} + Z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}]$ $(\bar{X} - Z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}}, \bar{X} + Z_{\frac{\alpha}{2}} \frac{S}{\sqrt{n}})$

$(\bar{X} - t_{\frac{\alpha}{2}, (n-1)} \frac{S}{\sqrt{n}}, \bar{X} + t_{\frac{\alpha}{2}, (n-1)} \frac{S}{\sqrt{n}})$

Ch7.2 Confidence Intervals for the Difference of Two Means

C.I for difference of two population means μ_X, μ_Y

Case 1: Normal with the population variances σ_X^2, σ_Y^2 known

Suppose that $X_1, \dots, X_n \sim N(\mu_X, \sigma_X^2)$ and $Y_1, \dots, Y_m \sim N(\mu_Y, \sigma_Y^2)$, where σ_X^2 and σ_Y^2 are known. Then, the difference of the sample means $(\bar{X} - \bar{Y})$ has a normal distribution,

$$N\left(\mu_X - \mu_Y, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right).$$

The standardized random variable

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\sigma_X^2/n + \sigma_Y^2/m}} \sim N(0, 1).$$

$$1 - \alpha = P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2})$$

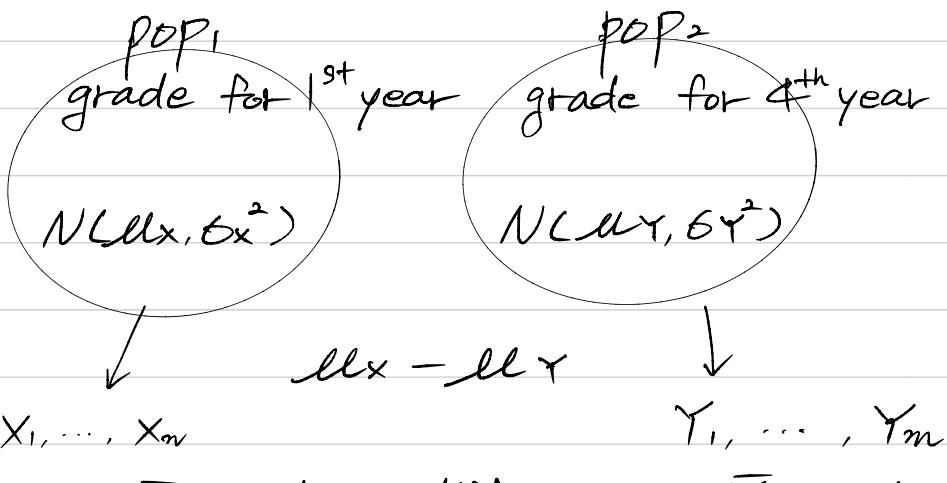
$$= P\left(-z_{\alpha/2} \leq \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\sigma_X^2/n + \sigma_Y^2/m}} \leq z_{\alpha/2}\right)$$

$$= P\left((\bar{X} - \bar{Y}) - z_{\alpha/2}\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}} \leq \mu_X - \mu_Y \leq (\bar{X} - \bar{Y}) + z_{\alpha/2}\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}\right)$$

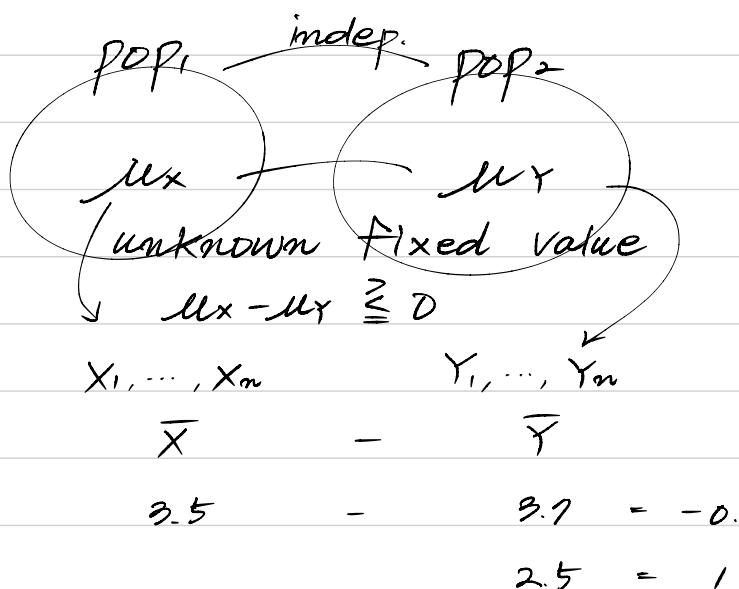
$(1 - \alpha) \times 100\%$ confidence interval for the difference $(\mu_X - \mu_Y)$

with known σ_X^2 , σ_Y^2

$$\left[(\bar{X} - \bar{Y}) - z_{\alpha/2}\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}, (\bar{X} - \bar{Y}) + z_{\alpha/2}\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}} \right]$$



(1 - α) · 100% C.I. for $\bar{x}_x - \bar{y}_Y$



1) $X_i \sim N(\mu_x, \sigma_x^2)$ $Y_i \sim N(\mu_Y, \sigma_Y^2)$ σ_x^2, σ_Y^2 : known.

$$\bar{X} \sim N\left(\mu_x, \frac{\sigma_x^2}{n}\right) \quad \bar{Y} \sim N\left(\mu_Y, \frac{\sigma_Y^2}{m}\right)$$

$$\bar{X} - \bar{Y} \sim N\left(\mu_x - \mu_Y, \frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}\right)$$

$$Z = \frac{(\bar{X} - \bar{Y}) - (\bar{x}_x - \bar{y}_Y)}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}} \sim N(0, 1)$$

$$P(-z_{\frac{\alpha}{2}} \leq Z \leq z_{\frac{\alpha}{2}}) = 1 - \alpha$$

length of C.I. depends on sample size as well as variance info

Case 2: Normal with unknown σ_X^2, σ_Y^2 but $\sigma_X^2 = \sigma_Y^2 = \sigma^2$

Suppose that $X_1, \dots, X_n \sim N(\mu_X, \sigma_X^2)$ and $Y_1, \dots, Y_m \sim N(\mu_Y, \sigma_Y^2)$, where σ_X^2 and σ_Y^2 are unknown but $\sigma_X^2 = \sigma_Y^2 = \sigma^2$.

Then, the difference of the sample means $(\bar{X} - \bar{Y})$ has a normal distribution, $N\left(\mu_X - \mu_Y, \sigma^2\left(\frac{1}{n} + \frac{1}{m}\right)\right)$.

The standardized random variable

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sigma\sqrt{1/n + 1/m}} \sim N(0, 1).$$

Also, since both random variables are independent,

$$U = \frac{(n-1)S_X^2}{\sigma^2} + \frac{(m-1)S_Y^2}{\sigma^2} \sim \chi_{(n+m-2)}^2.$$

$$2). X_i \sim N(\mu_x, \sigma_x^2) \quad Y_i \sim N(\mu_y, \sigma_y^2)$$

$$\sigma_x^2 = \sigma_y^2 = \sigma^2 \text{ unknown.}$$

Testing. $\sigma_x^2 = \sigma_y^2$ vs not.
 $\Rightarrow \sigma_x^2 = \sigma_y^2$

at this moment, just assume.

$$\bar{X} \sim N(\mu_x, \frac{\sigma^2}{n}), \quad Y \sim N(\mu_y, \frac{\sigma^2}{m})$$

$$\bar{X} - \bar{Y} \sim N(\mu_x - \mu_y, \sigma^2(\frac{1}{n} + \frac{1}{m}))$$

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0, 1)$$

$$\frac{(n-1)s_x^2}{\sigma^2} \sim \chi^2_{(n-1)}, \quad \frac{(m-1)s_y^2}{\sigma^2} \sim \chi^2_{(m-1)}$$

$$U = \frac{1}{\sigma^2} [(n-1)s_x^2 + (m-1)s_y^2] \sim \chi^2_{(n+m-2)}$$

$$T = \frac{Z}{\sqrt{U/(n+m-2)}}$$

$$= \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sqrt{\frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2} \cdot \left(\frac{1}{n} + \frac{1}{m}\right)}} \sim t(n+m-2).$$

$$S_p^2 = \frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^m (Y_i - \bar{Y})^2}{n+m-2}$$

point estimator for σ^2

$$(†). S_x^2 = \frac{\sum (X_i - \bar{X})^2}{n-1}$$

Since the sample means and sample variances are independent,

$$\begin{aligned} T &= \frac{Z}{\sqrt{U/(n+m-2)}} \sim t(n+m-2) \\ T &= \frac{\frac{(\bar{X}-\bar{Y})-(\mu_X-\mu_Y)}{\sigma\sqrt{1/n+1/m}}}{\sqrt{\left[\frac{(n-1)S_X^2}{\sigma^2} + \frac{(m-1)S_Y^2}{\sigma^2}\right]/(n+m-2)}} \\ &= \frac{(\bar{X}-\bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\left[\frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}\right] \left[\frac{1}{n} + \frac{1}{m}\right]}} \sim t(n+m-2) \end{aligned}$$

where S_p^2 is called the pooled estimator of σ^2 , defined by

$$S_p^2 = \frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}.$$

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t(n + m - 2)$$

$$1 - \alpha = P(-t_0 \leq T \leq t_0)$$

$$= P\left((\bar{X} - \bar{Y}) - t_0 S_p \sqrt{\frac{1}{n} + \frac{1}{m}} \leq \mu_X - \mu_Y \leq (\bar{X} - \bar{Y}) + t_0 S_p \sqrt{\frac{1}{n} + \frac{1}{m}}\right)$$

where $t_0 = t_{\alpha/2}(n + m - 2)$.

$(1 - \alpha) \times 100\%$ confidence interval for the difference $(\mu_X - \mu_Y)$

with known $\sigma_X^2 = \sigma_Y^2 = \sigma^2$

$$\left[(\bar{X} - \bar{Y}) - t_0 S_p \sqrt{\frac{1}{n} + \frac{1}{m}}, (\bar{X} - \bar{Y}) + t_0 S_p \sqrt{\frac{1}{n} + \frac{1}{m}} \right]$$

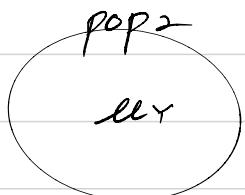
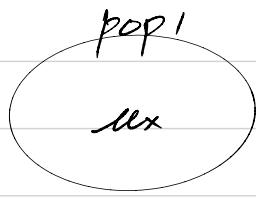


Case 3: Unknown σ_X^2, σ_Y^2 but large n, m

Suppose that $X_1, \dots, X_n \sim (\mu_X, \sigma_X^2)$ and $Y_1, \dots, Y_m \sim (\mu_Y, \sigma_Y^2)$, where σ_X^2 and σ_Y^2 are unknown but n, m are large. Then, by CLT, the difference of the sample means $(\bar{X} - \bar{Y})$ is approximate to a normal distribution, $N\left(\mu_X - \mu_Y, \frac{S_X^2}{n} + \frac{S_Y^2}{m}\right)$.

$(1 - \alpha) \times 100\%$ confidence interval for the difference $(\mu_X - \mu_Y)$ with unknown σ_X^2, σ_Y^2 but large n, m

$$\left[(\bar{X} - \bar{Y}) - z_{\alpha/2} \sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}, (\bar{X} - \bar{Y}) + z_{\alpha/2} \sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}} \right]$$



$$X_i \sim (\mu_x, \sigma_x^2) \quad Y_j \sim (\mu_Y, \sigma_Y^2) \quad \sigma_x^2, \sigma_Y^2: \text{unknown.}$$

$$i = 1, \dots, n \quad j = 1, \dots, m$$

n & m : large.

$(1-\alpha) \times 100\%$ C.I. for $\mu_x - \mu_Y$

By CLT,

$$\bar{X} \sim N(\mu_x, \frac{\sigma_x^2}{n}) \quad) \quad \bar{X} - \bar{Y} \sim N(\mu_x - \mu_Y, \frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m})$$

$$\bar{Y} \sim N(\mu_Y, \frac{\sigma_Y^2}{m}) \quad) \quad Z = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_Y)}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}} \sim N(0, 1)$$

$$P(-Z \leq Z \leq Z) = 1 - \alpha$$

$$\Leftrightarrow P((\bar{X} - \bar{Y}) - Z \frac{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}} \leq \mu_x - \mu_Y \leq (\bar{X} - \bar{Y}) + Z \frac{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}})$$

$$\Rightarrow P((\bar{X} - \bar{Y}) - Z \frac{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}} \leq \mu_x - \mu_Y \leq (\bar{X} - \bar{Y}) + Z \frac{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_Y^2}{m}}})$$

Case 4 (Welch's test): Normal with unknown σ_X^2, σ_Y^2 and small n, m

Suppose that $X_1, \dots, X_n \sim N(\mu_X, \sigma_X^2)$ and $Y_1, \dots, Y_m \sim N(\mu_Y, \sigma_Y^2)$, where σ_X^2 and σ_Y^2 are unknown and n, m are small. A $(1 - \alpha) \times 100\%$ confidence interval for the difference $(\mu_X - \mu_Y)$ is given by

$$\left[(\bar{X} - \bar{Y}) - t_{\alpha/2}([r]) \sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}, (\bar{X} - \bar{Y}) + t_{\alpha/2}([r]) \sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}} \right]$$

where $[r]$ is the largest integer not greater than r defined by

$$r = \frac{\left(\frac{S_X^2}{n} + \frac{S_Y^2}{m} \right)^2}{\frac{1}{n-1} \left(\frac{S_X^2}{n} \right)^2 + \frac{1}{m-1} \left(\frac{S_Y^2}{m} \right)^2}$$



Case 5: Paired datasets

Suppose that $(X_1, Y_1), \dots, (X_n, Y_n)$ are paired samples and $D_i = X_i - Y_i$ are from $N(\mu_D, \sigma_D^2)$, $i = 1, \dots, n$. The confidence interval of $\mu_D = \mu_X - \mu_Y$ is

$$T = \frac{\bar{D} - \mu_D}{S_D / \sqrt{n}} \sim t(n-1),$$

where \bar{D} and S_D^2 are the sample mean and variance of the differences.

A $(1 - \alpha) \times 100\%$ confidence interval for μ_D is given by

$$\left[\bar{D} - t_{\alpha/2}(n-1) \frac{S_d}{\sqrt{n}}, \bar{D} + t_{\alpha/2}(n-1) \frac{S_d}{\sqrt{n}} \right]$$

ID	POP1 dep	POP2	education program
1	Z_{11}		
2		Z_{21}	
:			
n		Z_{2n}	
	μ_x	μ_y	
			\Rightarrow
			$\begin{aligned} \text{POP}_1 - \text{POP}_2 \\ d_1 = Z_{11} - Z_{21} \\ d_2 = Z_{12} - Z_{22} \\ \vdots \\ \mu_D \end{aligned}$
			$N(\mu_D, \sigma_D^2)$
			unknown.
	$X_i \sim (\mu_x, \sigma_x^2)$	$Y_i \sim (\mu_y, \sigma_y^2)$	
	$i = 1, \dots, n.$		
	(X_i, Y_i) paired dataset		
	$i = 1, \dots, n.$		
			$D_i = X_i - Y_i \sim N(\mu_D, \sigma_D^2)$
			$i = 1, \dots, n.$

$$D. \bar{D} \sim N(\mu_D, \frac{\sigma_D^2}{n}) \Rightarrow Z = \frac{\bar{D} - \mu_D}{\sigma_D / \sqrt{n}} \sim N(0, 1).$$

$$U = \frac{(n-1)S_D^2}{\sigma_D^2} \sim \chi^2(n-1),$$

$$T = \frac{Z}{\sqrt{U/(n-1)}} = \frac{\bar{D} - \mu_D}{S_D / \sqrt{n}} \sim t(n-1)$$

$$P(-t_{\frac{\alpha}{2}(n-1)} \leq T \leq t_{\frac{\alpha}{2}(n-1)}) = 1 - \alpha.$$

Ch7.3 Confidence Intervals for Proportions

Confidence Intervals for Proportions

Case 1: Single population

Let X be the number of successes with n trials, $X \sim \text{Bin}(n, p)$. By CLT,

$$\frac{X/n - p}{\sqrt{p(1-p)/n}} = \frac{X - np}{\sqrt{np(1-p)}} \rightarrow N(0, 1), \quad \text{as } n \rightarrow \infty$$

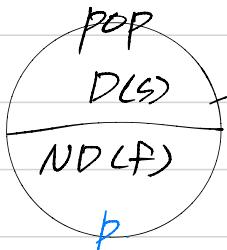
$X \sim (np, np(1-p))$
 $\frac{X}{n} \sim (p, \frac{p(1-p)}{n})$

$$1 - \alpha \approx P \left(-z_{\alpha/2} \leq \frac{X/n - p}{\sqrt{\frac{p(1-p)}{n}}} \leq z_{\alpha/2} \right)$$

$\sum X_i$

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

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



$X_1, \dots, X_n \sim \text{Ber}(p)$.

p : probability of S .

$$X_i = \begin{cases} 1 & \text{if } D(S) \\ 0 & \text{if } ND(f) \end{cases}$$

$$\hat{P} = \frac{1}{n} \cdot \sum_{i=1}^n X_i \quad n = 100 \quad \left[\begin{array}{l} 40 : P \\ 60 : ND \end{array} \right] \quad \hat{P}_{MLE} = 0.4$$

* $(1-\alpha) \times 100\%$. CI for p .

n large

$$\sum_{i=1}^n X_i \sim \text{Bin}(n, p)$$

By CLT

$$Z = \frac{\sum X_i - np}{\sqrt{np(1-p)}} \sim N(0, 1)$$

$$= \frac{\frac{\sum X_i}{n} - p}{\sqrt{\frac{p(1-p)}{n}}}$$

$$P(-z_{\frac{\alpha}{2}} \leq \frac{\frac{\sum X_i}{n} - p}{\sqrt{\frac{p(1-p)}{n}}} \leq z_{\frac{\alpha}{2}}) = 1 - \alpha$$

$$\Leftrightarrow P(\frac{\frac{\sum X_i}{n} - \hat{P}}{\sqrt{\frac{\hat{P}(1-\hat{P})}{n}}} \leq p \leq \frac{\frac{\sum X_i}{n} + z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}}{\sqrt{\frac{\hat{P}(1-\hat{P})}{n}}}) = 1 - \alpha$$

and approximately probability. plug in.

$$\Rightarrow P(\frac{\frac{\sum X_i}{n} - z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}}{\sqrt{\frac{\hat{P}(1-\hat{P})}{n}}} \leq p \leq \frac{\frac{\sum X_i}{n} + z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}}{\sqrt{\frac{\hat{P}(1-\hat{P})}{n}}}) = 1 - \alpha$$

$$\Leftrightarrow P\left(\left|\frac{\frac{\sum X_i}{n} - p}{\sqrt{\frac{p(1-p)}{n}}}\right| \leq z_{\frac{\alpha}{2}}\right) = 1 - \alpha$$

$$\Leftrightarrow P(-z_{\frac{\alpha}{2}} \leq p \leq z_{\frac{\alpha}{2}}) = 1 - \alpha$$

approximate $(1 - \alpha) \times 100\%$ confidence interval for p

$$\left[\frac{X}{n} - z_{\alpha/2} \sqrt{\frac{X/n(1-X/n)}{n}}, \frac{X}{n} + z_{\alpha/2} \sqrt{\frac{X/n(1-X/n)}{n}} \right]$$

$$X = \sum x_i$$

Alternative way is

$$\frac{|X/n - p|}{\sqrt{\frac{p(1-p)}{n}}} \leq z_{\alpha/2}$$

$$\left(\frac{X}{n} - p\right)^2 - \frac{z_{\alpha/2}^2 p(1-p)}{n} \leq 0$$

$$\left(1 + \frac{z_{\alpha/2}^2}{n}\right) p^2 - \left(2\hat{p} + \frac{z_{\alpha/2}^2}{n}\right) p + \hat{p}^2 \leq 0$$

where $\hat{p} = X/n$.

An approximate $(1 - \alpha) \times 100\%$ confidence interval for p is

$$\frac{\left[\hat{p} + z_{\alpha/2}^2/(2n)\right] \pm z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p})/n + z_{\alpha/2}^2/(4n^2)}}{1 + z_{\alpha/2}^2/n}$$

One-sided confidence interval

The one-sided upper confidence interval for p is

$$\left[0, \frac{X}{n} + z_\alpha \sqrt{\frac{X/n(1-X/n)}{n}} \right]$$

The one-sided lower confidence interval for p is

$$\left[\frac{X}{n} - z_\alpha \sqrt{\frac{X/n(1-X/n)}{n}}, 1 \right]$$

Case 2: $p_1 - p_2$ in two populations

Let X be the number of successes with n_1 trials for the first experiment. Let Y be the number of successes with n_2 trials for the second experiment. They are independent. Thus,

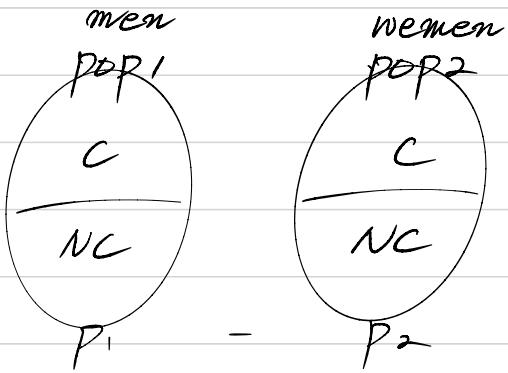
$X \sim \text{Bin}(n_1, p_1)$, $Y \sim \text{Bin}(n_2, p_2)$. As $n_1, n_2 \rightarrow \infty$,

$$\frac{X/n_1 - p_1}{\sqrt{p_1(1-p_1)/n_1}} \rightarrow N(0, 1), \quad \frac{Y/n_2 - p_2}{\sqrt{p_2(1-p_2)/n_2}} \rightarrow N(0, 1)$$

Thus the difference of both proportions is

$$\frac{X}{n_1} - \frac{Y}{n_2} \sim N\left(p_1 - p_2, \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}\right)$$

$$Z = \frac{\left(\frac{X}{n_1} - \frac{Y}{n_2}\right) - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}} \sim N(0, 1)$$



$$X_i \sim \text{Ber}(p) \quad i = 1, \dots, n \quad Y_j \sim \text{Ber}(p_2) \quad j = 1, \dots, m$$

$$\hat{P}_1 = \frac{\sum X_i}{n} \quad \hat{P}_2 = \frac{\sum Y_j}{m}$$

$$\hat{P}_1 = 0.01 \quad - \quad \hat{P}_2 = 0.008 = 0.002$$

* $(1-\alpha) \times 100\%$ CI for $P_1 - P_2$

n & m large.

$$Z_1 = \frac{\frac{\sum X_i}{n} - P_1}{\sqrt{\frac{P_1(1-P_1)}{n}}} \approx N(0,1) \quad Z_2 = \frac{\frac{\sum Y_j}{m} - P_2}{\sqrt{\frac{P_2(1-P_2)}{m}}} \approx N(0,1)$$

$$\Leftrightarrow \frac{\sum X_i}{n} = \hat{P}_1 \approx N(P_1, \frac{P_1(1-P_1)}{n}) \quad \frac{\sum Y_j}{m} = \hat{P}_2 \approx N(P_2, \frac{P_2(1-P_2)}{m})$$

$$\hat{P}_1 - \hat{P}_2 \approx N(P_1 - P_2, \frac{P_1(1-P_1)}{n} + \frac{P_2(1-P_2)}{m})$$

$$1 - \alpha \approx P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2})$$

$$= P\left(-z_{\alpha/2} \leq \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}} \leq z_{\alpha/2}\right)$$

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

where $q_1 = 1 - p_1$, $q_2 = 1 - p_2$, $\hat{p}_1 = X/n_1$, and $\hat{p}_2 = Y/n_2$.

approximate $(1 - \alpha) \times 100\%$ confidence interval for $p_1 - p_2$

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