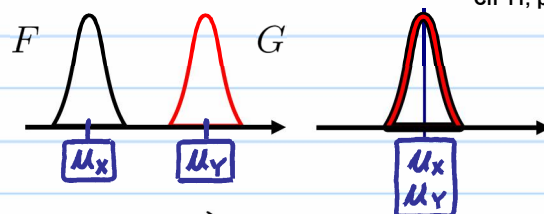


Methods based on normality assumptions

• variable: data
• parameter is fixed

- Assume that (1) F and G are normal, and (2) F and G have same variance.

- Thus, the statistical model is:



joint pdf

likelihood

cf

• variable: parameter
• data is fixed

(10/18)

$$\left. \begin{array}{l} \text{1st sample: } X_1, \dots, X_n \sim \text{i.i.d. } N(\mu_X, \sigma^2) \\ \text{2nd sample: } Y_1, \dots, Y_m \sim \text{i.i.d. } N(\mu_Y, \sigma^2) \end{array} \right\} \Leftarrow \text{independent } (*)$$

- This model contains three parameters: $\mu_X (\in \mathbb{R})$, $\mu_Y (\in \mathbb{R})$, $\sigma^2 (> 0)$.
- Under this model, the "difference" between F and G is simplified to be the difference between μ_X and μ_Y , i.e., $\Delta \equiv \mu_X - \mu_Y$ (\Leftarrow called "effect"), and $\mu_X - \mu_Y = 0 \Leftrightarrow$ no difference or no effect

a parameter

estimation testing

Review 1 (estimation of the parameters in one-sample normal model)

Consider $X_1, \dots, X_n \sim \text{i.i.d. } N(\mu, \sigma^2)$, and the statistics

2 parameters \rightarrow cf \rightarrow Data from S.R.S (LN, CH7, p11)

Textbook Sec. 6.3

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\text{unbiased}} \mu \quad \text{and} \quad s_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \xrightarrow{\text{unbiased}} \sigma^2$$

- distribution (exercise)

same as in survey sampling

not S_X^2

$(X_1, \dots, X_n) \in \mathbb{R}^n$
 $(X_1 - \bar{X}, \dots, X_n - \bar{X}) \in$ an $(n-1)$ -dim subspace of \mathbb{R}^n
 $\because U_1 + \dots + U_n = 0$

- \bar{X} and s_X^2 are independent \rightarrow joint = π marginals
- $\bar{X} \sim N(\mu, \sigma^2/n) \Rightarrow \sqrt{n}(\bar{X} - \mu)/\sigma \sim N(0, 1)$ (standardization)
- $(n-1)s_X^2 \sim \sigma^2 \chi_{n-1}^2 \Rightarrow (n-1)s_X^2/\sigma^2 \sim \chi_{n-1}^2$; $n-1$: degrees of freedom

r.v.
 $\sigma^2 \cdot Z$
 $Z \sim \chi_{n-1}^2$

- $(T_1 = \sum_{i=1}^n X_i, T_2 = \sum_{i=1}^n X_i^2)$ is a sufficient and complete statistic (exercise, Hint. 2-parameter exponential family)

MATH 2820, (統計學)

- Optimality

LN, CH8, p57

LN, CH8, p72-73

- $T_1/n = \bar{X}$ is the uniformly minimum variance unbiased estimator (UMVUE) of μ (exercise, Hint. Lehmann-Scheffe Thm)

- \bar{X} is the maximum likelihood estimator (MLE) of μ (exercise, Hint.

joint pdf:

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(X_i - \mu)^2}{2\sigma^2}}$$

$$\log\text{-likelihood} \propto -\frac{n}{2} \log(\sigma^2) - \frac{\sum_{i=1}^n (X_i - \mu)^2}{2\sigma^2}$$

LN, CH8, p20-21

$$\frac{1}{n-1} (T_2 - \frac{T_1^2}{n})$$

- s_X^2 is the UMVUE of σ^2 (exercise, Hint. Lehmann-Scheffe Thm)

- The MLE of σ^2 is $\frac{n-1}{n} s_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ (exercise)

Definition 1 (estimators of the parameters in the 2-sample normal model)

Under the two-sample normal model (*) in LNp.6, \rightarrow 3 parameters: μ_X, μ_Y, σ^2

- an intuitive estimator of μ_X is $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$, treated as one-sample
- an intuitive estimator of μ_Y is $\bar{Y} = \frac{1}{m} \sum_{j=1}^m Y_j$, $\because \bar{X} \xrightarrow{e} \mu_X, \bar{Y} \xrightarrow{e} \mu_Y$

- since $s_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ and $s_Y^2 = \frac{1}{m-1} \sum_{j=1}^m (Y_j - \bar{Y})^2$ estimate the same parameter σ^2 , we can pool them to get a better estimator:

n-1, m-1 d.f.s of S_X^2, S_Y^2

$$s_p^2 = \frac{(n-1)s_X^2 + (m-1)s_Y^2}{(n-1) + (m-1)} = \frac{(n-1)s_X^2 + (m-1)s_Y^2}{m+n-2}$$

Note 1 (Some notes about the estimator of σ^2)

- s_p^2 is called the **pooled sample variance**
- s_p^2 is a weighted average of the sample variances of the X_i 's and Y_j 's, where
 - the weights are proportional to the degrees of freedom, it is appropriate since if one sample is of much larger size than the other, the estimate of σ^2 from that sample is more reliable \Rightarrow it receives greater weight
 - since $E(s_X^2) = \sigma^2$ and $E(s_Y^2) = \sigma^2 \Rightarrow s_p^2$: an unbiased estimator of σ^2

from one-sample property in Review 1 (LNp.6)

Theorem 1 (distributions of the parameter estimators, 2-sample normal model)

- Since $(X_1, \dots, X_n), (Y_1, \dots, Y_m)$ are independent random variables

$\Rightarrow (\bar{X}, s_X^2, \bar{Y}, s_Y^2)$ are independent random variables

$\Rightarrow (\bar{X}, \bar{Y}, s_p^2)$ are independent random variables

$\because \bar{X}, \bar{Y}$ independent.
 $\text{Var}(\bar{X} - \bar{Y}) = \text{Var}(\bar{X}) + \text{Var}(\bar{Y}) = \sigma^2/n + \sigma^2/m$

• $\bar{X} \sim N(\mu_X, \sigma^2/n) \Rightarrow \sqrt{n}(\bar{X} - \mu_X)/\sigma \sim N(0, 1)$

• $\bar{Y} \sim N(\mu_Y, \sigma^2/m) \Rightarrow \sqrt{m}(\bar{Y} - \mu_Y)/\sigma \sim N(0, 1)$

• $\bar{X} - \bar{Y} \sim N\left(\mu_X - \mu_Y, \frac{\sigma^2}{n} + \frac{\sigma^2}{m}\right) \Rightarrow \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0, 1)$

standardization

$\Delta = \mu_X - \mu_Y \leftarrow e$

- Since (i) $(n-1)s_X^2/\sigma^2 \sim \chi_{n-1}^2$, (ii) $(m-1)s_Y^2/\sigma^2 \sim \chi_{m-1}^2$, and (iii) s_X^2 and s_Y^2 are independent,

$S_p^2 \sim \frac{\sigma^2}{m+n-2} \chi_{m+n-2}^2 \leftarrow$

$$\frac{(n-1)s_X^2 + (m-1)s_Y^2}{\sigma^2} = \frac{(m+n-2)s_p^2}{\sigma^2} \sim \chi_{m+n-2}^2$$

Theorem 2 (log-likelihood, 2-sample normal model)

testing \rightarrow likelihood ratio
estimation \rightarrow maximum likelihood

Under the two-sample normal model (*) in LNp.6, the log-likelihood is proportional to (exercise)

$$l(\mu_X, \mu_Y, \sigma^2) \propto -\frac{m+n}{2} \log(\sigma^2) - \frac{\sum_{i=1}^n (X_i - \mu_X)^2}{2\sigma^2} - \frac{\sum_{j=1}^m (Y_j - \mu_Y)^2}{2\sigma^2}$$

\square : data
 \square : parameter

$$= -\frac{1}{2\sigma^2} \left(\sum_{i=1}^n X_i^2 + \sum_{j=1}^m Y_j^2 \right) + \frac{\mu_X}{\sigma^2} \left(\sum_{i=1}^n X_i \right) + \frac{\mu_Y}{\sigma^2} \left(\sum_{j=1}^m Y_j \right)$$

$(-\frac{1}{2\sigma^2}, \frac{\mu_X}{\sigma^2}, \frac{\mu_Y}{\sigma^2})$
 $\in (-\infty, 0) \times \mathbb{R} \in$ 3-parameter exponential family

joint pdf:
 $\left[\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i - \mu_X)^2}{2\sigma^2}} \right] \times \left[\prod_{j=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(Y_j - \mu_Y)^2}{2\sigma^2}} \right]$

$\ast \mathbb{R}$

From the log-likelihood, we have

• $\frac{\partial l}{\partial \mu_X} = \frac{1}{\sigma^2} [(\sum_{i=1}^n X_i) - n \times \mu_X] = 0 \Rightarrow \hat{\mu}_{X, \text{MLE}} = \bar{X}$

• $\frac{\partial l}{\partial \mu_Y} = \frac{1}{\sigma^2} [(\sum_{j=1}^m Y_j) - m \times \mu_Y] = 0 \Rightarrow \hat{\mu}_{Y, \text{MLE}} = \bar{Y}$

• $\frac{\partial l}{\partial \sigma^2} = -\frac{m+n}{2\sigma^2} + \frac{\sum_{i=1}^n (X_i - \mu_X)^2}{2\sigma^4} + \frac{\sum_{j=1}^m (Y_j - \mu_Y)^2}{2\sigma^4} = 0 \Rightarrow \hat{\sigma}_{\text{MLE}}^2 = \frac{n+m-2}{n+m} S_p^2$

Theorem 3 (UMVUE and MLE of the parameters in the 2-sample normal model)

- $(R_1 = \sum_{i=1}^n X_i, R_2 = \sum_{j=1}^m Y_j, R_3 = \sum_{i=1}^n X_i^2 + \sum_{j=1}^m Y_j^2)$ is a sufficient and complete statistic (**Hint**. 3-parameter exponential family)
- \bar{X} ($= R_1/n$) is the UMVUE (by Lehmann-Scheffe Thm) and MLE of μ_X
- \bar{Y} ($= R_2/m$) is the UMVUE (by Lehmann-Scheffe Thm) and MLE of μ_Y

- (by Lehmann-Scheffe Thm) The pooled sample variance s_p^2 is the UMVUE of σ^2 , since (i) s_p^2 is unbiased, and (ii)

$$\begin{aligned} (m+n-2)s_p^2 &= (n-1)s_X^2 + (m-1)s_Y^2 \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{j=1}^m (Y_j - \bar{Y})^2 \\ &= \left(\sum_{i=1}^n X_i^2 \right) - n\bar{X}^2 + \left(\sum_{j=1}^m Y_j^2 \right) - m\bar{Y}^2 = R_3 - \left(\frac{R_1^2}{n} \right) - \left(\frac{R_2^2}{m} \right) \end{aligned}$$

$n\bar{X}^2 = n \left(\frac{\sum X_i}{n} \right)^2 = n \left(\frac{R_1}{n} \right)^2$

not unbiased

- The MLE of σ^2 is $\frac{m+n-2}{m+n} s_p^2 = \frac{(n-1)s_X^2 + (m-1)s_Y^2}{m+n}$ cf. $m+n-2$

Question 2 (how to claim $\Delta=0$ or $\Delta \neq 0$?)

Under the two-sample normal model (*) in LNp.6, consider the parameter

$$\Delta = \mu_X - \mu_Y.$$

Notice that

estimation of Δ

$$\Delta = 0 \Leftrightarrow \text{no difference in the two samples} \Rightarrow \text{MLE of } g(\theta) \text{ is } g(\hat{\theta}_{MLE})$$

invariance property of MLE:

$$\hat{\theta}_{MLE} \xrightarrow{e} \theta \quad g(\theta)$$

- The UMVUE (by Lehmann-Scheffe Thm and $\hat{\Delta} = R_1/n - R_2/m$) and MLE of Δ is $\hat{\Delta} = \bar{X} - \bar{Y}$. point estimator Recall. duality between C.I. and testing

 $\hat{\Delta} \sim \text{normal}$ in Thm 1 (LNp.8)

- But, $\hat{\Delta} \neq 0$ is not a strong enough evidence to reject $\Delta = 0$ (Note. $P(\hat{\Delta} \neq 0) = 1$). A better way is to examine if a C.I. of Δ contains 0.

- Q: how to construct an interval estimator for Δ ? interval estimator

Review 2 (pivotal quantity of θ)

Recall LN, CH7, P.33

A pivotal quantity for θ is a function of data X_1, \dots, X_n and the parameter θ , denoted by

$$Q(\mathbf{X}, \theta) = Q(X_1, \dots, X_n, \theta), \leftarrow \text{a r.v., but not a statistic}$$

if the distribution of $Q(\mathbf{X}, \theta)$ is irrelevant to all parameters.

Theorem 4 (confidence interval of Δ , 2-sample normal model)

Under the two-sample normal model (*) in LNp.6,

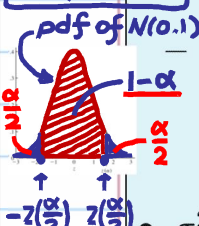
- σ^2 known (σ^2 is not a parameter)

– a pivotal quantity of Δ is

Recall. distribution of $\bar{X} - \bar{Y}$ in Thm 1 (LNp.8)irrelevant to μ_X, μ_Y function of data & Δ only

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

a known constant $\rightarrow \sigma \sqrt{\frac{1}{n} + \frac{1}{m}}$ std. error of $\bar{X} - \bar{Y}$



- a 100(1 - α)% C.I. for Δ is $(\bar{X} - \bar{Y}) \pm z(\alpha/2) \times (\sigma \sqrt{\frac{1}{n} + \frac{1}{m}})$ since $1 - \alpha = P(|Q_{Z,\Delta}| < z(\alpha/2))$ data: fixed, Δ : changed

$$= P\left((\bar{X} - \bar{Y}) - z(\alpha/2)\sigma \sqrt{\frac{1}{n} + \frac{1}{m}} < \Delta < (\bar{X} - \bar{Y}) + z(\alpha/2)\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}\right)$$

- σ^2 unknown (σ^2 is a parameter)

– a pivotal quantity of Δ is

parameter \rightarrow function of data onlyfunction of data & Δ only

$$Q_{T,\Delta} = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{s_p \sqrt{\frac{1}{n} + \frac{1}{m}}} = \frac{[(\bar{X} - \bar{Y}) - \Delta] / (\sigma \sqrt{\frac{1}{n} + \frac{1}{m}})}{\frac{s_p}{\sigma} \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{m+n-2}$$

estimated std. error of $\bar{X} - \bar{Y}$ $\sim \chi^2_{m+n-2}$ (Thm 1, LNp.8) $\sqrt{\left[\frac{(m+n-2)s_p^2}{\sigma^2} \right] \frac{1}{m+n-2}}$ independent (Thm 1, LNp.8)

irrelevant to μ_X, μ_Y, σ^2

TBp.193

C.I. when σ^2 knowna $100(1 - \alpha)\%$ C.I. for Δ is $(\bar{X} - \bar{Y}) \pm t_{m+n-2}(\alpha/2) \times \left(s_p \sqrt{\frac{1}{n} + \frac{1}{m}}\right)$ **Note 2** (A note about the confidence intervals of Δ) t is a distribution symmetric about 0estimated std. error of $\bar{X} - \bar{Y}$

These confidence intervals are of the form $(\text{estimate}) \pm (\text{critical value}) \times [(\text{estimated}) \text{ standard error}]$, evaluates the accuracy of estimator

where the (estimated) standard error of $\bar{X} - \bar{Y}$ is $\sigma_{\bar{X} - \bar{Y}} = \sigma \sqrt{\frac{1}{n} + \frac{1}{m}}$ when σ^2 is known, and is $s_{\bar{X} - \bar{Y}} = s_p \sqrt{\frac{1}{n} + \frac{1}{m}}$ when σ^2 is unknown.

Example 2 (confidence interval of Δ , heat of fusion of ice, cont. Ex.1 in LNp.3)

Statistical modeling: assume 2-sample normal model (*) in LNp.6

- $n = 13, \bar{X}_A = 80.02, s_A = 0.024; m = 8, \bar{X}_B = 79.98, s_B = 0.031$
- $s_p = \sqrt{\frac{12}{19} s_A^2 + \frac{7}{19} s_B^2} = 0.027, s_{\bar{X}_A - \bar{X}_B} = s_p \sqrt{\frac{1}{13} + \frac{1}{8}} = 0.012$
- A 95% confidence interval for $\Delta = \mu_A - \mu_B$ is $(\bar{X}_A - \bar{X}_B) \pm t_{19}(0.025) \times s_{\bar{X}_A - \bar{X}_B} = (0.04) \pm (2.093) \times (0.012) = (0.015, 0.065)$.

Question 3 (how to perform testing of $\Delta=0$?)Note: $0 \notin (0.015, 0.065) \Rightarrow \text{reject } \Delta=0$

- Recall.** duality between confidence interval and hypothesis testing
- Q:** What are the hypothesis testings corresponding to these confidence intervals of Δ ?

Theorem 5 (z-test and t-test for $\Delta=\Delta_0$, 2-sample normal model)

Under the two-sample normal model (*) in LNp.6, consider the null and alternative hypotheses:

$$H_0 : \mu_X - \mu_Y = \Delta = \Delta_0 \quad \text{vs.} \quad H_A : \mu_X - \mu_Y = \Delta \neq \Delta_0$$

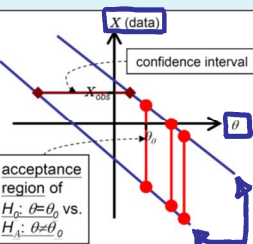
where Δ_0 is a known constant (Note. if $\Delta_0 = 0$, $H_0 : \mu_X = \mu_Y$ vs. $H_A : \mu_X \neq \mu_Y$), and H_A is a two-sided alternative. From the duality between C.I. and testing,

C.I. fixed $|Q_{Z, \Delta}| < z(\alpha/2)$ changed $|Q_{T, \Delta}| < t_{m+n-2}(\alpha/2)$ testing changed $|Q_{Z, \Delta_0}| < z(\alpha/2)$ fixed $|Q_{T, \Delta_0}| < t_{m+n-2}(\alpha/2)$

the corresponding test of these confidence intervals are:

- test statistic
 - σ^2 known: $Z = \frac{(\bar{X} - \bar{Y}) - \Delta_0}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}}$ (cf. $Q_{Z, \Delta}$ in LNp.11) \leftarrow pivotal quantity
 - σ^2 unknown: $T = \frac{(\bar{X} - \bar{Y}) - \Delta_0}{s_p \sqrt{\frac{1}{n} + \frac{1}{m}}}$ (cf. $Q_{T, \Delta}$ in LNp.11) \leftarrow estimated std. error of $\bar{X} - \bar{Y}$

Q: Are they statistics?



$$|Q_{Z, \Delta}| < z(\alpha/2)$$

$$|Q_{T, \Delta}| < t_{m+n-2}(\alpha/2)$$

check Thm 4 (LNp.11)

- null distribution
 - σ^2 known: under H_0 , $Z \sim N(0, 1)$
 - σ^2 unknown: under H_0 , $T \sim t_{m+n-2}$
- level- α rejection region $\Delta = \Delta_0$, Thm 4 (LNp.11)
 - σ^2 known: $|Z| > z(\alpha/2)$, called z-test (reasonable?)
 - σ^2 unknown: $|T| > t_{m+n-2}(\alpha/2)$, called t-test (reasonable?)

Why $\alpha/2$?
 \therefore 2-sided alternative

Note. The t-test (or z-test) rejects H_0 if and only if its corresponding C.I. does not include Δ_0 .