

(1/8)

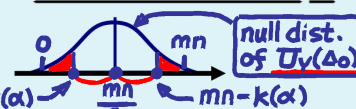
Assume it is an integer

irrelevant to F

- * the test statistic: $U_Y(\Delta_0) = \# \{X_i < Y_j - \Delta_0\} = \# \{Y_j - X_i > \Delta_0\}$,
 (Note: X_i is fixed, Y_j is changed)
- * the acceptance region: $k(\alpha) \leq U_Y(\Delta_0) \leq mn - k(\alpha)$,
 where $k(\alpha)$ is the critical value determined by the significance level α
 (Note. The null distribution of $U_Y(\Delta_0)$ is symmetric about $mn/2$.)

What if $\Delta_0 = 0$?

check Thm 10 (LNp.40)



- By the duality of test and C.I., a $100(1 - \alpha)\%$ confidence interval for Δ is

$$C = \{ \Delta \mid k(\alpha) \leq U_Y(\Delta) \leq mn - k(\alpha) \}.$$

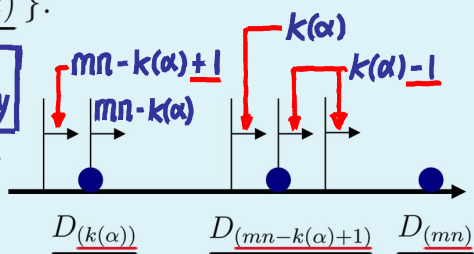
$k(\alpha)$ can be obtained from the critical value of W_Y (check Thm 11, LNp.39) in Table B, TBp. A21

- Let $D_{(1)}, D_{(2)}, \dots, D_{(mn)}$ denote the ordered mn differences $(Y_j - X_i)$'s.

$$C = [D_{(k(\alpha))}, D_{(mn-k(\alpha)+1)}].$$

pivotal quantity

What is its dist.?



- To see this,

- * if $\Delta_0 = D_{(k(\alpha))}$, then $U_Y(\Delta_0) = \# \{Y_j - X_i > \Delta_0\} = mn - k(\alpha)$,
 if $\Delta_0 < D_{(k(\alpha))}$, then $U_Y(\Delta_0) = \# \{Y_j - X_i > \Delta_0\} \geq mn - k(\alpha) + 1$,
 thus, $D_{(k(\alpha))}$ is the leftmost point of the confidence interval C ,
 * if $\Delta_0 < D_{(mn-k(\alpha)+1)}$, then $U_Y(\Delta_0) = \# \{Y_j - X_i > \Delta_0\} \geq k(\alpha)$,
 if $\Delta_0 \geq D_{(mn-k(\alpha)+1)}$, then $U_Y(\Delta_0) = \# \{Y_j - X_i > \Delta_0\} \leq k(\alpha) - 1$,
 thus, $D_{(mn-k(\alpha)+1)}$ is the rightmost point of the confidence interval C .

Example 6 (C.I. for Δ , heat of fusion of ice, cont. Ex.4 in LNp.34 & Ex.5 in LNp.42)

- $n = 13$ (method A), $m = 8$ (method B), $W_B = 51$. Under null, $E(W_B) = 88$.
- Under the significant level $\alpha = 0.05$, the critical value for W_B^* is 60 (Ex.4, LNp.34) \Rightarrow acceptance region: $61 \leq W_B \leq 88 + (88 - 61) = 115$

$$k(\alpha) \Leftrightarrow 25 \leq U_B = W_B - [8(8+1)]/2 = W_B - 36 \leq 79.$$

- After sorting the $mn = 8 \times 13 = 104$ differences $(Y_j - X_i)$'s, we get

$$D_{(k(\alpha)=25)} = -0.07 \quad \text{and} \quad D_{(mn-k(\alpha)+1=80)} = -0.01.$$

$mn - k(\alpha)$
consistent with the test result in Ex.4

normality seems valid in this case (check box plot in LNp.4)

A 95% confidence interval for Δ is $(-0.07, -0.01)$, which does not contain 0.

[$\xleftrightarrow{\text{cf.}}$ the C.I. $(0.015, 0.065)$ given in Ex.2 (LNp.12) $\times (-1) = (0.01, 0.07)$]

– Note that the Δ here is the $-\Delta$ in Ex.2.

under normality

Δ can be defined as "difference of medians" under (\square)

check Note 8 (LNp.39)

– In this case, the C.I. based on the nonparametric model is slightly wider than the one based on the normal model.

– But, the latter C.I. relies on the validity of normality assumption.

Theorem 15 (Bootstrap confidence interval for $\pi_\Delta (\Leftrightarrow \Delta)$)

Consider the nonparametric model (\diamond) in LNp.35 or the nonparametric model (\square) in LNp.27. (Note. (1) (\square) has more models than (\diamond) (2) $\pi_\Delta = P(X < Y)$ is well-defined in (\diamond) and (\square) (3) Δ is well-defined only in (\diamond))

- Bootstrapping is a numerical method that can be used to gain information about the sampling distribution of $\hat{\pi}_\Delta = \frac{1}{mn} (\# \{X_i < Y_j\}) \xrightarrow{e} \pi_\Delta$, and the estimated standard error of $\hat{\pi}_\Delta$.
 (a r.v. only has one obs. of this r.v.)

TBp. 284-285

Under H_0
Under H_A

In bootstrap, we

$$\hat{\pi}_{\Delta} \leftarrow \left\{ \begin{array}{l} X_1, \dots, X_n \sim \text{i.i.d. from } \underline{F} \\ Y_1, \dots, Y_m \sim \text{i.i.d. from } \underline{G} \end{array} \right\} \leftarrow \text{independent} \quad \leftarrow \text{cf.}$$

- replace the true cdf \underline{F} (unknown) by the empirical cdf \hat{F}_n (known) of $(X_1, \dots, X_n) = (x_1, \dots, x_n)$ [\hat{F}_n : assigns x_i 's equal probabilities $1/n$]
- replace the true cdf \underline{G} (unknown) by the empirical cdf \hat{G}_m (known) of $(Y_1, \dots, Y_m) = (y_1, \dots, y_m)$ [\hat{G}_m : assigns y_j 's equal probabilities $1/m$]
- Re-sample (generate data $X'_1, \dots, X'_n, Y'_1, \dots, Y'_m$ using simulation) from this model:

$$\left\{ \begin{array}{l} X'_1, \dots, X'_n \sim \text{i.i.d. from } \hat{F}_n \\ Y'_1, \dots, Y'_m \sim \text{i.i.d. from } \hat{G}_m \end{array} \right\} \leftarrow \text{independent} \quad \leftarrow \text{Are they similar?}$$

joint distribution of these r.v.'s: completely known

 - X'_1, \dots, X'_n is a with-replacement sample from the population $\{x_1, \dots, x_n\}$,
 - Y'_1, \dots, Y'_m is a with-replacement sample from the population $\{y_1, \dots, y_m\}$.
- Repeat the re-sampling procedure many times, say B times, and known data
 - at each time, compute $\hat{\pi}'_{\Delta} = \frac{1}{mn} \#\{X'_i < Y'_j\}$ from $(X'_1, \dots, X'_n, Y'_1, \dots, Y'_m)$
 - this produces a bootstrap sample: $(\hat{\pi}'_{\Delta,1}, \dots, \hat{\pi}'_{\Delta,B})$ \leftarrow can be regarded as a sample of $\hat{\pi}_{\Delta}$
- A histogram of $(\hat{\pi}'_{\Delta,1}, \dots, \hat{\pi}'_{\Delta,B})$ offers an indication of the sampling distribution of $\hat{\pi}_{\Delta}$ (\Rightarrow a $100(1-\alpha)\%$ C.I. of π_{Δ} is $[\hat{\pi}'_{\Delta,(B(\alpha/2))}, \hat{\pi}'_{\Delta,(B(1-\alpha/2))}]$),
- the standard deviation of $(\hat{\pi}'_{\Delta,1}, \dots, \hat{\pi}'_{\Delta,B}) \xrightarrow{e}$ the standard error of $\hat{\pi}_{\Delta}$.

a r.v.

❖ Reading: textbook, 11.2.3

a r.v.

• Comparing paired samples \leftrightarrow Independent samples

• Problem formulation and statistical modeling

observed data (random variables) $\{X_1, \dots, X_n\}$ and $\{Y_1, \dots, Y_n\}$

Such connection does not exist in 2-independent samples.

• X_i 's, Y_i 's are continuous quantities
 • $X - Y$ is meaningful

2-indep samples \leftarrow cf. Data

the comparison of their means is meaningful.

population

Recall. Strata in survey sampling. (LN, CH7, P.59-63)

a stratum

homogeneous within stratum

For example, in human population,

- X_i 's: left eye vision
- Y_i 's: right eye vision of the i th person

For example, in medical study,

- X_i 's: treatment
- Y_i 's: control applied on the i th twins

2-indep samples \leftarrow cf.

s.r.s., $N \rightarrow \infty$:
 without replacement \approx with replacement (\Rightarrow i.i.d.)

• $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) \sim$ i.i.d. with a common continuous joint distribution $\underline{F}(x, y) \leftarrow$ population distribution

• (X_i, Y_i) , i.e., \underline{F} , might not be independent \leftarrow Estimation of a ratio (LN, CH7, P.38~39)

U	block	V
1	1	X_1
1	2	X_2
\vdots	\vdots	\vdots
1	n	X_n
2	1	Y_1
2	2	Y_2
\vdots	\vdots	\vdots
2	n	Y_n

Let random variables Z_1, \dots, Z_n represent the variability of the n members sampled from the population. Independent samples

Assume Z_1, \dots, Z_n are i.i.d. from a population distribution $H(z)$.

Let $\underline{X} = \underline{\phi}(\underline{Z})$ and $\underline{Y} = \underline{\psi}(\underline{Z})$, where $\underline{\phi}, \underline{\psi}$ contain random components, and denote

$(\underline{\phi}, \underline{\psi}): \underline{Z} \rightarrow (\underline{X}, \underline{Y})$
 $(H(z)) \quad (F(x, y))$

– $F(x, y)$: the joint distributions of (X, Y) , not necessarily independent

– $\underline{\mu}_X$ and $\underline{\mu}_Y$: the means of \underline{X} and \underline{Y} , respectively,

– $\underline{\Delta} = \underline{\mu}_X - \underline{\mu}_Y$.

Then, for $1 \leq i \leq n$,

$\begin{cases} X_i = \underline{\phi}(Z_i) \\ Y_i = \underline{\psi}(Z_i) \end{cases}$

Because $Z_1, \dots, Z_n \sim \text{i.i.d. } H(z)$,
 $(X_1, Y_1), \dots, (X_n, Y_n) \sim \text{i.i.d. } F(x, y)$

Further assume that

$\begin{cases} X_i = \underline{\phi}(Z_i) = \underline{\mu}_X + \epsilon_{1i} \\ Y_j = \underline{\psi}(Z_{n+j}) = \underline{\mu}_Y + \epsilon_{2j} \end{cases}$

mean zero in two independent samples case. check LNp 2

(1) $\underline{\phi}(Z) = \underline{\phi}^*(Z) + \underline{\delta}_1$ and $\underline{\psi}(Z) = \underline{\psi}^*(Z) + \underline{\delta}_2$, where $\underline{\phi}^*, \underline{\psi}^*$ are fixed functions and $\underline{\delta}_1, \underline{\delta}_2$ are independent random variables with mean 0

(2) $\underline{Z}, \underline{\delta}_1, \underline{\delta}_2$ are independent

(3) $\underline{\psi}^*(\underline{Z}) = \underline{\phi}^*(\underline{Z}) - \underline{\Delta} \Rightarrow \underline{\Delta} = \underline{\phi}^*(\underline{Z}) - \underline{\psi}^*(\underline{Z})$

$\underline{\mu}_X = E[\underline{\phi}^*(Z)]$
 $\underline{\mu}_Y = E[\underline{\psi}^*(Z)]$

Then,

$\begin{cases} X_i = \underline{\phi}^*(Z_i) + \underline{\delta}_{1i} = \underline{\phi}^*(Z_i) + \underline{\delta}_{1i} = \underline{\mu}_X + [\underline{\phi}^*(Z_i) - \underline{\mu}_X] + \underline{\delta}_{1i} \\ Y_i = \underline{\psi}^*(Z_i) + \underline{\delta}_{2i} = \underline{\phi}^*(Z_i) - \underline{\Delta} + \underline{\delta}_{2i} = \underline{\mu}_Y + [\underline{\phi}^*(Z_i) - \underline{\mu}_X] + \underline{\delta}_{2i} \end{cases}$

Q: What are the sources of variation in ϵ 's and δ 's? If we apply the above formulation to the case of two independent samples, then add (3)

$\begin{cases} X_i = \underline{\mu}_X + \epsilon_{1i} = \underline{\phi}^*(Z_i) + \underline{\delta}_{1i} = \underline{\mu}_X + (\underline{\phi}^*(Z_i) - \underline{\mu}_X) + \underline{\delta}_{1i} \\ Y_j = \underline{\mu}_Y + \epsilon_{2j} = \underline{\psi}^*(Z_{n+j}) + \underline{\delta}_{2j} = \underline{\mu}_Y + (\underline{\phi}^*(Z_{n+j}) - \underline{\mu}_X) + \underline{\delta}_{2j} \end{cases}$

A comparison

– Increase sample sizes: increase information about $\underline{\mu}_X$ and $\underline{\mu}_Y$ (signal)

– 2 independent \rightarrow paired: suppress the variation of error (noise)

Theorem 16 (A brief variance comparison of paired and independent samples) Under (3)

Consider the models in the dashed frames. Under the two models,

– $\epsilon = [\underline{\phi}^*(Z) - \underline{\mu}_X] + \underline{\delta} \Rightarrow \text{Var}(\epsilon) = \text{Var}[\underline{\phi}^*(Z)] + \text{Var}(\underline{\delta}) \geq \text{Var}(\underline{\delta})$

– 2 independent samples ($n = m$) due to the variation of members in the population

– $X_i - Y_j = (\underline{\mu}_X - \underline{\mu}_Y) + (\epsilon_{1i} - \epsilon_{2j}) \equiv \sigma_\epsilon^2$

$= (\underline{\mu}_X - \underline{\mu}_Y) + [\underline{\phi}^*(Z_i) - \underline{\phi}^*(Z_{n+j})] + (\underline{\delta}_{1i} - \underline{\delta}_{2j})$

– $\bar{X} - \bar{Y} = (\underline{\mu}_X - \underline{\mu}_Y) + (\bar{\epsilon}_1 - \bar{\epsilon}_2)$ paired samples

$\Rightarrow \bar{X} - \bar{Y} \xrightarrow{e} \underline{\Delta}$ and $\text{Var}(\bar{X} - \bar{Y}) = (\sigma_{\epsilon_1}^2/n) + (\sigma_{\epsilon_2}^2/n)$

$\text{Var}(D_i) = \sigma_{\delta_1}^2 + \sigma_{\delta_2}^2$
 $D_i \equiv X_i - Y_i = (\mu_X - \mu_Y) + (\delta_{1i} - \delta_{2i})$
 $\bar{D} = \bar{X} - \bar{Y} = (\mu_X - \mu_Y) + (\bar{\delta}_1 - \bar{\delta}_2)$
 $\Rightarrow \bar{X} - \bar{Y} \xrightarrow{e} \Delta = \mu_X - \mu_Y$ and
 $\text{Var}(\bar{X} - \bar{Y}) = (\sigma_{\delta_1}^2/n) + (\sigma_{\delta_2}^2/n)$

paired samples $\xrightarrow{\text{cf.}}$ 2-indep. samples
 larger $\xrightarrow{\text{cf.}}$ smaller $\xrightarrow{\text{cf.}}$

$\text{Var}(\bar{X} - \bar{Y})$ under the 2-independent-sample model with the sample size
 $n' = \frac{\sigma_{\epsilon}^2}{\sigma_{\delta}^2} n \quad (\geq n)$

same estimator of Δ as in 2 indep. samples

- Paired sample is more effective than independent samples in this case.

Theorem 17 (Conditions under which paired sample is more effective)

remove assumption (3)

Consider the models in the dotted frames of LNp.48. Under the two models,

- $E(X) = E[\phi^*(Z) + \delta_1] = E[\phi^*(Z)] = \mu_X$
 $E(Y) = E[\psi^*(Z) + \delta_2] = E[\psi^*(Z)] = \mu_Y$
 - $\text{Var}(X) = \text{Var}[\phi^*(Z) + \delta_1] = \text{Var}[\phi^*(Z)] + \text{Var}(\delta_1) \equiv \sigma_X^2 (= \sigma_{\epsilon_1}^2)$
 $\text{Var}(Y) = \text{Var}[\psi^*(Z) + \delta_2] = \text{Var}[\psi^*(Z)] + \text{Var}(\delta_2) \equiv \sigma_Y^2 (= \sigma_{\epsilon_2}^2)$
 - 2 independent samples ($n = m$)
 - $\text{Cov}(X_i, Y_j) = \text{Cov}[\phi^*(Z_i) + \delta_{1i}, \psi^*(Z_{n+j}) + \delta_{2j}] = 0 \xrightarrow{\text{cf.}}$ paired samples
 - $E(\bar{X} - \bar{Y}) = \mu_X - \mu_Y = \Delta$
 - $\text{Var}(\bar{X} - \bar{Y}) = (\sigma_X^2 + \sigma_Y^2)/n \xrightarrow{\text{cf.}}$ paired samples
- Note. X (or Y) have same marginal dist. in 2-indep. & paired cases.

- paired samples
 - $\text{Cov}(X_i, Y_i) = \text{Cov}[\phi^*(Z_i) + \delta_{1i}, \psi^*(Z_i) + \delta_{2i}] = \text{Cov}[\phi^*(Z_i), \psi^*(Z_i)] \equiv \sigma_{XY}$
 - $\text{Cov} \left\{ \begin{matrix} Z_i \leftrightarrow Z_i \\ Z_i \leftrightarrow \delta_{1i} \\ Z_i \leftrightarrow \delta_{2i} \\ \delta_{1i} \leftrightarrow \delta_{2i} \end{matrix} \right\} \xrightarrow{\text{cf.}}$ 2-indep. samples
 - * Note. We do not observe $(\phi^*(Z_i), \psi^*(Z_i))$'s. But, σ_{XY} can be estimated using (X_i, Y_i) 's data.
 - * Denote the correlation of (X_i, Y_i) by $\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$
 - * Notice that $\rho_{XY} \neq (\leq \text{in absolute value})$
 - $\text{Cor}[\phi^*(Z), \psi^*(Z)] = \frac{\sigma_{XY}}{\sigma_{\phi^*(Z)} \sigma_{\psi^*(Z)}} \xrightarrow{\text{cf.}}$ 2-indep. samples
 - Let $D_i = X_i - Y_i, i = 1, \dots, n$. Then,
 - D_1, \dots, D_n are i.i.d. $\leftarrow \because (X_1, Y_1), \dots, (X_n, Y_n)$ are independent
 - $E(D_i) = \mu_X - \mu_Y$
 - $\text{Var}(D_i) = \text{Var}(X_i) + \text{Var}(Y_i) - 2 \text{Cov}(X_i, Y_i) = \sigma_X^2 + \sigma_Y^2 - 2 \sigma_{XY}$
 - Since $\bar{D} = \bar{X} - \bar{Y} \xrightarrow{e} \Delta$
 - $E(\bar{D}) = \mu_X - \mu_Y = \Delta$
 - $\text{Var}(\bar{D}) = \text{Var}(\bar{X} - \bar{Y}) = (\sigma_X^2 + \sigma_Y^2 - 2 \sigma_{XY})/n$
 $= (\sigma_X^2 + \sigma_Y^2 - 2 \rho_{XY} \sigma_X \sigma_Y)/n \xrightarrow{\text{cf.}}$ 2-indep. samples

Why?

Under (3)

- If $\rho_{XY} > 0$ ($\Leftrightarrow \sigma_{XY} > 0 \Leftrightarrow \text{Cov}[\phi^*(Z), \psi^*(Z)] > 0$), then paired sample is more effective than independent samples.

- When $\psi^*(Z) = \phi^*(Z) - \Delta$,

$$\begin{aligned}\text{Cov}[\phi^*(Z), \psi^*(Z)] &= \text{Cov}[\phi^*(Z), \phi^*(Z) - \Delta] \\ &= \text{Cov}[\phi^*(Z), \phi^*(Z)] > 0.\end{aligned}$$

1/13

- Q: Why are independent samples more effective than paired samples when $\sigma_{XY} < 0$?

- If $\sigma_X^2 = \sigma_Y^2 = \sigma^2$, then in the paired case

$$\sigma_D^2 = \text{Var}(\bar{D}) = [2\sigma^2(1 - \rho_{XY})]/n \quad \text{and} \quad \sigma_{\bar{X}-\bar{Y}}^2 = \text{Var}(\bar{X} - \bar{Y}) = 2\sigma^2/n$$

in the unpaired case. The relative efficiency is $\sigma_D^2 / \sigma_{\bar{X}-\bar{Y}}^2 = 1 - \rho_{XY}$.

- If $\rho_{XY} = 0.5$, a paired design with n pairs of subjects yields the same precision as an unpaired design with 2n subjects per treatment.

- From now on, the analyses of paired data are based on

$$D_i = X_i - Y_i, \quad i = 1, \dots, n.$$

- Statistical modeling for D_i's: $D_1, \dots, D_n \sim \text{i.i.d. } F \Leftarrow \text{one-sample model}$

