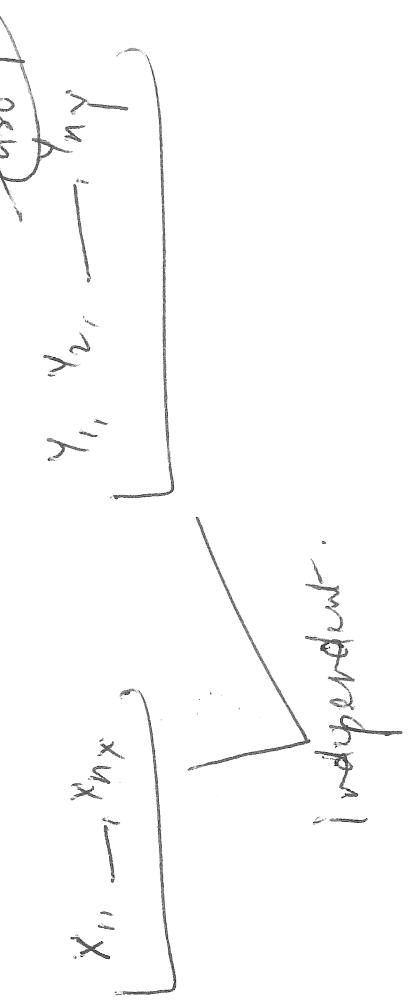


# Independent or paired samples?

## Design 1: (Two independent samples)



independent.

## Design 2: (Paired samples)

D	$(X_1, Y_1)$	$D = X - Y$
1	$(X_1, Y_1)$	$D_1$
2	$(X_2, Y_2)$	$D_2$
⋮	⋮	⋮
n	$(X_n, Y_n)$	$D_n$

Generally: Assume different subjects

in different groups  
and paired.

Given a choice, if feasible, a paired design is better than the two-sample (indep.) design. This is because in case of the former, any diff. in  $X$  and  $Y$  is just due to the effect of the intervention in study. In case of the latter, the effect of the intervention gets confounded with the difference in the groups.

# CI for $\mu_X - \mu_Y$ with paired design

$\rightarrow$  paired index

- Data:  $(X_i, Y_i), i=1, 2, \dots, n$
- Parameter of interest:  $\mu_X - \mu_Y = E[D] = \mu_D$
- Define:  $D = X - Y, D_i = X_i - Y_i, i = 1, \dots, n$   
 $\Rightarrow E[D] = E[X] - E[Y] = \mu_X - \mu_Y$
- Apply one-sample procedure to the differences.

100(1 -  $\alpha$ )% CI for  $\mu_D$  assuming  $D_1, \dots, D_n \sim N(\mu_D, \sigma_D^2)$ :

$$\text{Pivot: } \frac{\bar{D} - \mu_D}{\sigma_D / \sqrt{n}} \sim N(0, 1)$$
$$\left\{ \begin{array}{l} \bar{D} \pm z_{\alpha/2} \frac{\sigma_D}{\sqrt{n}} \\ \bar{D} \pm t_{n-1, \alpha/2} \frac{\sigma_D}{\sqrt{n}} \end{array} \right\}$$

$\sigma_D^2$  is unknown,  
 $\sigma_D^2$  is sample variance of the differences

Approximate  $100(1 - \alpha)\%$  CI for  $\mu_D$  if  $n$  is large:

$$\bar{D} \pm z_{\alpha/2} \frac{s_D}{\sqrt{n}} \sim \bar{D} \pm z_{\alpha/2} \frac{s_D}{\sqrt{n}}$$

Pivot:  $\frac{\bar{D} - \mu_D}{s_D / \sqrt{n}} \sim N(0, 1)$ , provided  $n$  is large.

Q: What is the pivot here?

# CI for $\mu_X - \mu_Y$ with two independent samples (First: normality)

Setup:

$$\text{Population: } X \sim N[\mu_X, \sigma_X^2]$$

$$\text{Population: } Y \sim N[\mu_Y, \sigma_Y^2]$$

$$X_1, \dots, X_{n_X}$$

→ first

→ first

→ first

**Scenario 1:** No assumption regarding  $\sigma_X^2$  and  $\sigma_Y^2$  — they may be equal or unequal.

Estimators of parameters:

- $\mu_X: \bar{X}$
- $\mu_Y: \bar{Y}$
- $\sigma_X^2: S_X^2$  — sample var.
- $\sigma_Y^2: S_Y^2$  — sample var.
- $\mu_X - \mu_Y: \bar{X} - \bar{Y}$

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

Pivot:

$$\frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\hat{SE}(\bar{X} - \bar{Y})} \sim \mathcal{N}(0, 1)$$

$\left| \begin{array}{l} \text{Var}(\bar{X} - \bar{Y}) = \text{Var}(\bar{X}) + \text{Var}(\bar{Y}) \\ \bar{X} \text{ and } \bar{Y} \text{ are independent} \end{array} \right.$

$$\Rightarrow \hat{SE}[\bar{X} - \bar{Y}] = \sqrt{\frac{s_x^2}{n_X} + \frac{s_y^2}{n_Y}}$$

- Satterthwaite's approximation: The distribution of the pivot can be approximated by  $t_{\nu}$  distribution where

$$\nu = \frac{\left( \frac{s_x^2}{n_X} + \frac{s_y^2}{n_Y} \right)^2}{\frac{s_x^4}{n_X(n_X-1)} + \frac{s_y^4}{n_Y(n_Y-1)}}$$

- Approximate  $100(1 - \alpha)\%$  CI for  $\mu_X - \mu_Y$ :

$$(\bar{X} - \bar{Y}) \pm t_{\alpha/2} \hat{SE}(\bar{X} - \bar{Y}),$$

- When  $n_X, n_Y$  are large, the CI is:  $(\bar{X} - \bar{Y}) \pm 2s_{\bar{X}} \hat{SE}(\bar{X} - \bar{Y})$

- Is the normality assumption needed when  $n_X, n_Y$  are large?  $\mathcal{N}(\bar{v})$ .

Normality + equal variance:

boxplot can be used to check this assumption.

**Scenario 2:** Assume common variance, i.e.,  $\sigma_X^2 = \sigma_Y^2 = \sigma^2$ .

- Estimators of  $(\mu_X, \mu_Y, \mu_X - \mu_Y)$ :  $\bar{X}, \bar{Y}, \bar{X} - \bar{Y}$
- Estimate  $\sigma^2$  using the pooled sample variance  $S_p^2$ :

$$S_p^2 = \frac{(n_x-1) S_x^2 + (n_y-1) S_y^2}{n_x + n_y - 2}$$

$$\begin{aligned} E[S_p^2] &= \frac{E[(n_x-1) S_x^2 + (n_y-1) S_y^2]}{n_x + n_y - 2} \\ &= \frac{E[(n_x-1) E[S_x^2] + (n_y-1) E[S_y^2]]}{n_x + n_y - 2} \\ &= \frac{\sigma^2(n_x-1) + \sigma^2(n_y-1)}{n_x + n_y - 2} \\ &= \sigma^2 \end{aligned}$$

$$\text{Var}(S_p^2) = \sigma^2 \left( \frac{1}{n_x} + \frac{1}{n_y} \right)$$

$$\text{SE}(\bar{X} - \bar{Y}) = \sqrt{S_p^2 \left( \frac{1}{n_x} + \frac{1}{n_y} \right)} = \sqrt{S_p^2} \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}$$

• Estimated SE:

Pivot:

$$\frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\text{SE}(\bar{X} - \bar{Y})} = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\frac{\sigma^2}{n_X} + \frac{\sigma^2}{n_Y}}} \sim t_{n_X + n_Y - 2} \quad (\text{Result}).$$

- $100(1 - \alpha)\%$  CI for  $\mu_X - \mu_Y$ :  $(\bar{X} - \bar{Y}) \pm t_{n_X + n_Y - 2, \alpha/2} \text{SE}(\bar{X} - \bar{Y})$ .
- When  $n_X, n_Y$  are large, the CI is:  $(\bar{X} - \bar{Y}) \pm 2\bar{s}_p \text{SE}(\bar{X} - \bar{Y})$ .
- Is the normality assumption needed when  $n_X, n_Y$  are large?  $\mathcal{N}_0$ .

## Recap

Comparing two population means:

① Paired samples : Apply one-sample procedures to the differences.

② Independent samples :

Confidence interval for  $\mu_1 - \mu_2$

Population normal

Yes

No assumption  
of  $\sigma_1^2$  or  $\sigma_2^2$   
regarding

No

Are  $n_1$  and  $n_2$  large?

NO

YES

Take advanced  
course

$(\bar{X} - \bar{Y}) \pm t_{n_1+n_2-2} \cdot s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$

$(\bar{X} - \bar{Y}) \pm z_{0.95} \cdot \sqrt{\frac{s_x^2}{n_1} + \frac{s_y^2}{n_2}}$

$(\bar{X} - \bar{Y}) \pm t_{n_1+n_2-2} \cdot s_p \sqrt{\frac{s_x^2}{n_1} + \frac{s_y^2}{n_2}}$

where  $s_p^2$  is obtained by Satterthwaite approx

# Large-sample CI for $p_X - p_Y$ with two independent samples

**Setup:**

$$\text{Pop: } X \sim \text{Bernoulli}(p_X)$$

$$\text{Pop: } Y \sim \text{Bernoulli}(p_Y)$$

$X_1, \dots, X_{n_X}$   
 $Y_1, \dots, Y_{n_Y}$

$\hat{p}_X = \frac{1}{n_X} \sum_{i=1}^{n_X} X_i$   
 $\hat{p}_Y = \frac{1}{n_Y} \sum_{i=1}^{n_Y} Y_i$

• Estimators of  $(p_X, p_Y, p_X - p_Y)$ :

$\hat{p}_X, \hat{p}_Y, \hat{p}_X - \hat{p}_Y$

$$\text{SE}(\hat{p}_X - \hat{p}_Y) = \sqrt{\frac{p_X(1-p_X)}{n_X} + \frac{p_Y(1-p_Y)}{n_Y}}$$

• Estimated SE:

$$\widehat{\text{SE}}(\hat{p}_X - \hat{p}_Y) = \sqrt{\frac{\hat{p}_X(1-\hat{p}_X)}{n_X} + \frac{\hat{p}_Y(1-\hat{p}_Y)}{n_Y}}$$

**Pivot:**

$$(\hat{p}_X - \hat{p}_Y) \pm \widehat{\text{SE}}(\hat{p}_X - \hat{p}_Y)$$

• Approximate  $100(1 - \alpha)\%$  CI:

$$\left( \hat{p}_X - \hat{p}_Y \right) - \left( \hat{p}_X - \hat{p}_Y \right) \sim N(0, 1)$$

provided both  $n_X$  and  $n_Y$  are large.

Confidence interval for  $p_1 - p_2$

Are  $n_1$  and  $n_2$  large?

NO

Take advanced course

YES

$$(\hat{p}_1 - \hat{p}_2) \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

## CI for a function of $\theta$

Issue: Suppose  $(L, U)$  is a  $100(1 - \alpha)\%$  CI for  $\theta$ . How to get a  $100(1 - \alpha)\%$  CI for  $g(\theta)$ , where  $g$  is a monotonically increasing function of  $\theta$ ?

Have:

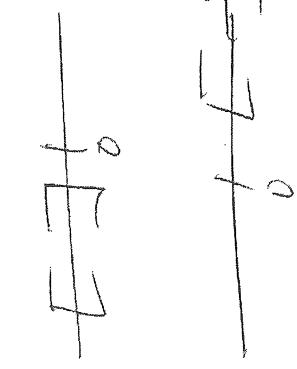
$$\Pr [L \leq \theta \leq U] = 1 - \alpha \quad (\text{for all } \theta)$$
$$\Pr [g(L) \leq g(\theta) \leq g(U)] = 1 - \alpha$$

$\Rightarrow [g(L), g(U)]$  is a  $100(1 - \alpha)\%$  CI for  $g(\theta)$ .



Consider the following scenarios:  $\alpha$  for  $p_1, p_2$

Case 1:



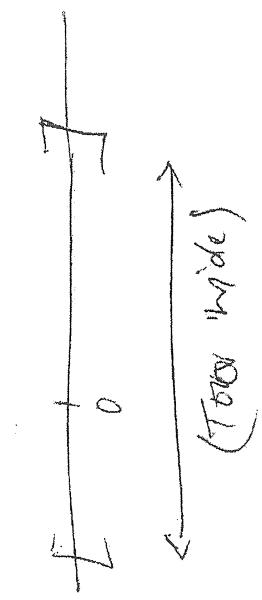
Case 2:



Case 3:



Case 4:



Case 5:

Case 5: Need more data  
(so that  $\alpha$  is not  
that wide)

Case 6:

Case 6: may still conclude that  
 $p_1 = p_2$ , but this is  
a borderline (need  
more data for a more  
exact conclusion)

From R: 95% CI for  $\hat{p}_X - \hat{p}_Y$ :  $[ -0.05, 0.05 ]$

(b) What assumptions, if any, did you make to compute the interval in (a)? Do the assumptions seem reasonable?

Since  $\sigma$  is right in the middle of the CI,  
we can conclude that  $\hat{p}_X = \hat{p}_Y$ .

both  $\hat{p}_X$  and  $\hat{p}_Y$   
are large.

## Example 1

The data below show the sugar content (as a % of weight) of several national brands of children's and adults' cereals.

$X = \text{Sugar content in typical children's cereal}$  —  $\mu_X = E[X]$

$Y = \text{Sugar content in typical adult's cereal}$  —  $\mu_Y = E[Y]$

**Children's cereals:** 40.3, 55, 45.7, 43.3, 50.3, 45.9, 53.5, 43, 44.2, 44, 47.4, 44, 33.6, 55.1, 48.8, 50.4, 37.8, 60.3, 46.5

**Adults' cereals:** 20, 30.2, 2.2, 7.5, 4.4, 22.2, 16.6, 14.5, 21.4, 3.3, 6.6, 7.8, 10.6, 16.2, 14.5, 4.1, 15.8, 4.1, 2.4, 3.5, 8.5, 10, 1, 4.4, 1.3, 8.1, 4.7, 18.4

- (a) Is it reasonable to assume that each sample comes from a normal distribution? No, may be ok for children's cereals but not for adults' cereals
- (b) Can the variances of the two distributions be assumed to be equal? Justify your answer.