

# AMS 572 Data Analysis I

## Power and sample size for two population means

Pei-Fen Kuan

Applied Math and Stats, Stony Brook University

Based on  $Z$  test

# Power calculation

$$H_0 : \mu_1 - \mu_2 = \Delta_1$$

$$H_a : \mu_1 - \mu_2 = \Delta_2 \neq \Delta_1$$

$$\text{power} = P(\text{reject } H_0 | H_a)$$

$$\begin{aligned} &= P(Z_0 \geq z_{\alpha/2} | \mu_1 - \mu_2 = \Delta_2) + P(Z_0 \leq -z_{\alpha/2} | \mu_1 - \mu_2 = \Delta_2) \\ &= P\left(\frac{\bar{X} - \bar{Y} - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \geq z_{\alpha/2} | \mu_1 - \mu_2 = \Delta_2\right) \\ &\quad + P\left(\frac{\bar{X} - \bar{Y} - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \leq -z_{\alpha/2} | \mu_1 - \mu_2 = \Delta_2\right) \\ &= P\left(\frac{\bar{X} - \bar{Y} - \Delta_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \geq z_{\alpha/2} - \frac{\Delta_2 - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} | \mu_1 - \mu_2 = \Delta_2\right) \\ &\quad + P\left(\frac{\bar{X} - \bar{Y} - \Delta_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \leq -z_{\alpha/2} - \frac{\Delta_2 - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} | \mu_1 - \mu_2 = \Delta_2\right) \\ &= P\left(Z \geq z_{\alpha/2} - \frac{\Delta_2 - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right) + P\left(Z \leq -z_{\alpha/2} - \frac{\Delta_2 - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right) \end{aligned}$$

# Sample size determination for a given margin of error $E$

Based on exact or the large sample approximate z-test ;

$$n_1 = n_2 = n$$

Suppose the margin of error is  $E$  with probability  $(1 - \alpha)$

$$P(|(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)| \leq E) = 1 - \alpha$$

$$\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

$$P(-E \leq \bar{X} - \bar{Y} - (\mu_1 - \mu_2) \leq E) = 1 - \alpha$$

$$P\left(\frac{-E}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \leq \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \leq \frac{E}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right) = 1 - \alpha$$

# Sample size determination for a given CI length $L$

# Sample size determination for a given power

$$H_0 : \mu_1 - \mu_2 = \Delta_1$$

$$H_a : \mu_1 - \mu_2 = \Delta_2 > (\text{or } <) \Delta_1$$

$$1 - \beta = \text{power} = P(\text{reject } H_0 | H_a)$$

$$= P(Z_0 \geq z_\alpha | \mu_1 - \mu_2 = \Delta_2)$$

$$= P\left(\frac{\bar{X} - \bar{Y} - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}}} \geq z_\alpha | \mu_1 - \mu_2 = \Delta_2\right)$$

$$= P\left(\frac{\bar{X} - \bar{Y} - \Delta_2}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}}} \geq z_\alpha - \frac{\Delta_2 - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}}} | \mu_1 - \mu_2 = \Delta_2\right)$$

$$z_\alpha - \frac{\Delta_2 - \Delta_1}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}}} = -z_\beta$$

# Sample size determination for a given power

One-tailed test, exact Z

$$n = \frac{(z_{\alpha} + z_{\beta})^2 (\sigma_1^2 + \sigma_2^2)}{(\Delta_2 - \Delta_1)^2}$$

One-tailed test, approximate Z

$$n \doteq \frac{(z_{\alpha} + z_{\beta})^2 (\sigma_1^2 + \sigma_2^2)}{(\Delta_2 - \Delta_1)^2}$$

Two-tailed test, exact or approximate Z

$$n \doteq \frac{(z_{\alpha/2} + z_{\beta})^2 (\sigma_1^2 + \sigma_2^2)}{(\Delta_2 - \Delta_1)^2}$$

Example: A new method of making concrete blocks has been proposed. To test whether or not the new method increases the compressive strength, 5 sample blocks are made by each method.

New Method	14	15	13	15	16
Old Method	13	15	13	12	14

- (a) Get a 95% CI for the mean difference of the 2 methods.
- (b) At  $\alpha = 0.05$ , can you conclude the new method is better?

Write a SAS and R program for part (b).



# SAS Code

```
data block ;  
input method $ strength ;  
datalines ;  
new 14  
new 15  
new 13  
new 15  
new 16  
old 13  
old 15  
old 13  
old 12  
old 14  
;  
run ;
```

# SAS Code

```
proc univariate data=block normal plot ;  
class method ;  
var strength ;  
run ;
```

```
proc ttest data=block sides=U ;  
class method ;  
var strength ;  
run ;
```

The SAS System				
Method	Variances	DF	t Value	Pr > t
Pooled	Equal	8	1.66	0.0673
Satterthwaite	Unequal	8	1.66	0.0673

Equality of Variances				
Method	Num DF	Den DF	F Value	Pr > F
Folded F	4	4	1.00	1.0000

## R Code and Output

```
> new <- c(14,15,13,15,16)
> old <- c(13,15,13,12,14)
> var.test(new,old)
```

F test to compare two variances

data: new and old

F = 1, num df = 4, denom df = 4, p-value = 1

alternative hypothesis: true ratio of variances is not equal to

95 percent confidence interval:

0.1041175 9.6045299

sample estimates:

ratio of variances

1

# R Code and Output

```
> shapiro.test(new)
```

Shapiro-Wilk normality test

data: new

W = 0.96086, p-value = 0.814

```
> shapiro.test(old)
```

Shapiro-Wilk normality test

data: old

W = 0.96086, p-value = 0.814

# R Code and Output

```
> t.test(new,old,var.equal=TRUE,alternative='greater')
```

Two Sample t-test

```
data:  new and old
```

```
t = 1.6641, df = 8, p-value = 0.06733
```

```
alternative hypothesis: true difference in means is greater than
```

```
95 percent confidence interval:
```

```
-0.1409392          Inf
```

```
sample estimates:
```

```
mean of x mean of y
```

```
14.6      13.4
```

Example: An experiment was done to determine the effect on dairy cattle of a diet supplement with liquid whey. While no differences were noted in milk production between the group with a standard diet (hay + grain + water) and the experimental group with whey supplement (hay + grain + whey), a considerable difference was noted in the amount of hay ingested. For a 2-tailed test with  $\alpha=0.05$ , determine the approximate number of cattle that should be included in each group if we want  $\beta \leq 0.1$  for  $|\mu_1 - \mu_2| \geq 0.5$ . Previous study has shown  $\sigma \approx 0.8$ .

Example: Do fraternities help or hurt your academic progress at college? To investigate this question, 5 students who joined fraternities in 1998 were randomly selected. It was shown that their GPA before and after they joined the fraternities are as follows.

Student	1	2	3	4	5
Before	3	4	3	3	2
After	2	3	3	2	1
Diff	1	1	0	1	1

Test the hypothesis at  $\alpha=0.05$



$$H_0 : \mu_d = 0$$

$$H_a : \mu_d \neq 0$$

Assumption : the difference follows a normal distribution.

$$\bar{X}_d = 0.8, s_d = 0.447, n = 5, \alpha = 0.05$$

$$\text{Test statistic : } T_0 = \frac{\bar{X}_d - 0}{s_d / \sqrt{n}} \stackrel{H_0}{\sim} t_{n-1}$$

$$|T_0| = 4.02 > t_{4,0.025} = 2.776$$

We reject  $H_0$  at  $\alpha=0.05$  and conclude fraternities does hurt

# SAS Code

```
data frat ;  
input before after ;  
diff = before - after ;  
datalines ;  
3 2  
4 3  
3 3  
3 2  
2 1  
;  
run ;  
  
proc univariate data=frat normal ;  
var diff ;  
run ;
```

# R Code and Output

```
> before <- c(3,4,3,3,2)
> after <- c(2,3,3,2,1)
> shapiro.test(before-after)
```

Shapiro-Wilk normality test

```
data:  before - after
W = 0.55218, p-value = 0.000131
```

```
### normality assumption is violated,
### thus the appropriate approach is to
### use non-parametric approach from Chapter 14.
```

# R Code and Output

```
> t.test(before,after,paired=TRUE)
```

Paired t-test

data: before and after

t = 4, df = 4, p-value = 0.01613

alternative hypothesis: true difference in means is not equal to

95 percent confidence interval:

0.244711 1.355289

sample estimates:

mean of the differences

0.8

# R Code and Output

```
> d <- c(1,1,0,1,1)
> t.test(d)
```

One Sample t-test

```
data:  d
t = 4, df = 4, p-value = 0.01613
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 0.244711 1.355289
sample estimates:
mean of x
      0.8
```