

두 평균의 비교

두 평균의 비교: 독립표본과 대응표본

- ◆ 귀무가설은 두 변수가 독립이라는 것이고, 대립가설은 두 변수가 상호관련이 있다는 것

- $H_0: X \perp Y$ $H_1: X \sim Y$

- ◆ 이변량 분석 방법의 선택

- 두 변수가 모두 범주형인 경우: 교차분석
- 두 변수가 모두 크기가 의미를 갖는 수치형 혹은 박쥐형인 경우: 상관관계 분석
- 한 변수는 범주형이고 다른 한 변수는 수치형 혹은 박쥐형인 경우: ?
 - ❖ 예) 성별로 주량의 차이가 있는지 알아보려고 한다. 변수 X 는 범주형으로 남자는 1 그리고 여자는 2로 구분하고, 변수 Y 는 수치형으로 주량을 나타냄

두 평균의 비교: 독립표본과 대응표본

- ◆ 남자 100명, 여자 100명을 각각 표본으로 무작위 선출한다.
- ◆ 이들에게 각각 주량이 얼마인지 물어보고 남성집단과 여성집단 각각의 평균 주량을 구한다.
- ◆ 남성표본의 평균 주량을 \bar{y}_1 이라고 하고, 여성표본의 평균 주량을 \bar{y}_2 라고 하자. 변수 X 와 변수 Y 가 독립이라면, 두 모평균 과 의 관계는 어떠하겠는가?
- ◆ 성별 구분과 주량이 관련 없다면, 남성집단 주량의 모평균 과 여성집단의 모평균 는 같을 것이다.
- ◆ 반대로 독립이 아니라면, 두 모평균 값은 다를 것이다. 따라서 귀무가설 과 대립가설은 다음과 같다.
 - $H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$
- ◆ 남성집단과 여성집단처럼 독립적인 두 모집단의 평균을 비교하기 위해서, 각 모집단에서 추출된 표본을 독립표본(independent sample)이라고 한다.

두 평균의 비교: 독립표본

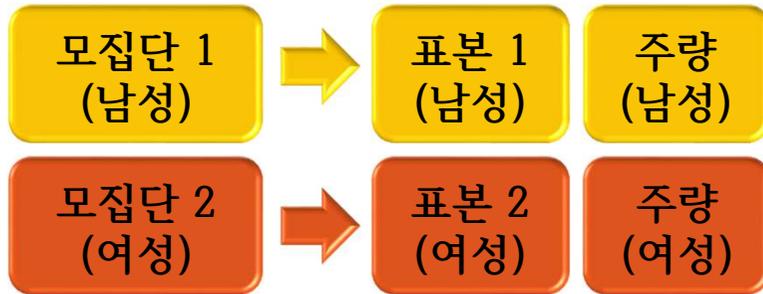
- ◆ 남자 100명, 여자 100명을 각각 표본으로 무작위 선출한다.
- ◆ 이들에게 각각 주량이 얼마인지 물어보고 남성집단과 여성집단 각각의 평균 주량을 구한다.
- ◆ 남성표본의 평균 주량을 \bar{y}_1 이라고 하고, 여성표본의 평균 주량을 \bar{y}_2 라고 하자. 변수 X 와 변수 Y 가 독립이라면, 두 모평균 과 의 관계는 어떠하겠는가?
- ◆ 성별 구분과 주량이 관련 없다면, 남성집단 주량의 모평균 과 여성집단의 모평균 는 같을 것이다.
- ◆ 반대로 독립이 아니라면, 두 모평균 값은 다를 것이다. 따라서 귀무가설 과 대립가설은 다음과 같다.

$$\bullet H_0: \mu_1 = \mu_2 \qquad H_1: \mu_1 \neq \mu_2$$

- ◆ 남성집단과 여성집단처럼 독립적인 두 모집단의 평균을 비교하기 위해서, 각 모집단에서 추출된 표본을 독립표본(independent sample)이라고 한다.

두 평균의 비교: 독립표본과 대응표본

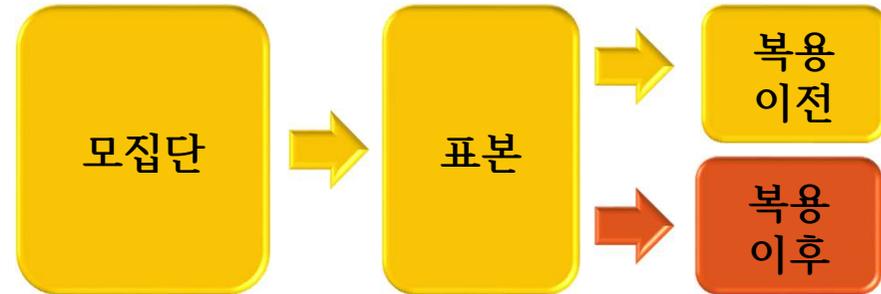
◆ 독립표본(배타적인 두 집단)



◆ 예) 남성과 여성의 주량

◆ 예) 한국인과 미국인의 수명

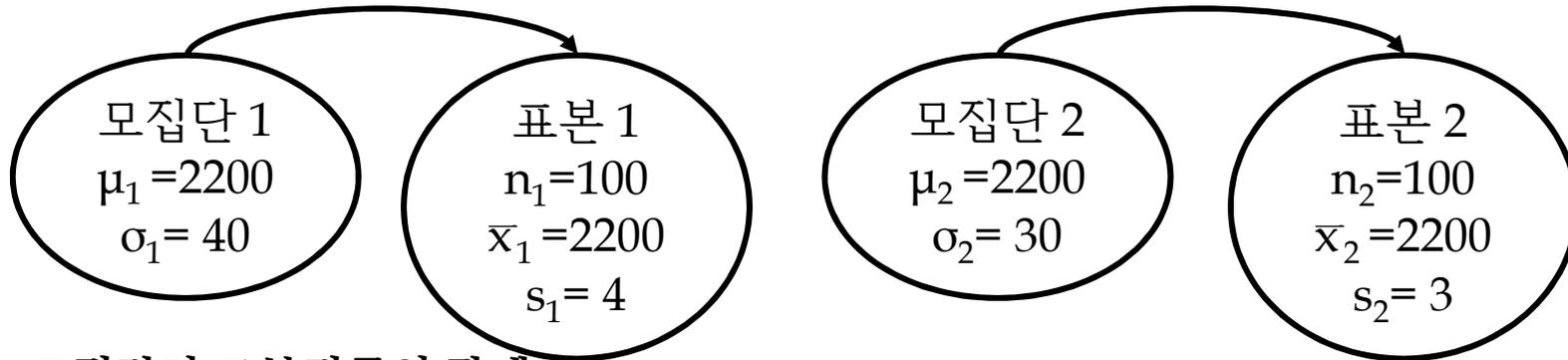
◆ 대응표본 (동일한 집단)



◆ 예) 다이어트 복용 전후의 변화 여부

◆ 예) 광고시청 전후 브랜드에 대한 태도 변화

두 표본평균 차이의 표본분포



◆ 모집단과 표본평균의 관계

● $\bar{X} = \mu, s = \frac{\sigma}{\sqrt{n}}$

◆ 두 표본평균 차이의 분포 (모집단의 분산을 아는 경우)

● $\bar{X}_1 - \bar{X}_2$ 의 평균 : $\mu_1 - \mu_2$, $\bar{X}_1 - \bar{X}_2$ 의 표준편차 : $\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

● $(\bar{X}_1 - \bar{X}_2) \sim N\left(\mu_1 - \mu_2, \left(\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right)^2\right), (\bar{X}_1 - \bar{X}_2) \sim N(0, 5^2)$

두 표본평균 차이 검정의 종류

- ◆ Case 1 : 모분산을 알고 있는 경우 (비현실적)
 - ◆ Case 2 : 모분산을 모르고 등분산인 경우
 - ◆ Case 3 : 모분산을 모르고 등분산이 아닌 경우
 - ◆ Case 4: 대응표본(쌍표본)의 검정
-
- ◆ Case 2, Case 3은 등분산 여부인지를 1차적으로 검정(Levene's test for equality of variance)하고 그 후 적합한 식에 따라 '두 평균의 차이 유무에 대한 검정 통계량'을 계산함

CASE 1 : 모분산을 알고 있는 경우

◆ 하루 섭취하는 음식물 칼로리에 있어서, 남성과 여성의 차이가 있는지를 알아보려고 한다. 남성집단의 분산은 40^2 이고, 여성집단의 분산은 30^2 이라고 알려져 있다. 남성집단의 평균을 그리고 여성집단의 평균을 라고 할 때, 귀무가설과 대립가설은 다음과 같다.

● $H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$ 혹은 $H_0: \mu_1 - \mu_2 = 0$ $H_1: \mu_1 - \mu_2 \neq 0$

◆ 두 집단의 모평균을 추정하기 위해, 각 집단으로부터 100개씩의 표본을 추출하였다. 그 결과 남성집단의 평균은 2,205 이고, 여성집단의 평균은 2,195 이었다. 그러면 이를 근거로 두 집단의 모평균에 차이가 있다고 말할 수 있는가?

◆ 두 집단의 모평균에 차이 $\sim N(2205-2195, 5^2)$

◆ 95% 신뢰구간

◆ $(\bar{x}_1 - \bar{x}_2) \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = 10 \pm (1.960)5 = [0.2, 19.8]$

CASE 2: 모분산을 모르고 등분산인 경우

- ◆ 모분산을 모르고 등분산이라는 가정이 가능한 경우 (표본의 분산을 통해 모집단이 서로 등분산이라고 판단되는 경우)
- ◆ 두 모집단 평균의 차이의 분포는 t-분포를 따름
- ◆ 분산의 경우 통합분산(pooled variance)를 통해 정해짐

- $$s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}$$

- ◆ 차이의 분포는

- $$(\bar{x}_1 - \bar{x}_2) \sim t\left(\mu_1 - \mu_2, s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right); n_1 + n_2 - 2\right)$$

CASE 3 : 모분산을 모르고 등분산이 아닌 경우

◆ 대표본($n \geq 30$)인 경우 \rightarrow 대수의 법칙(정규분포처럼 취급)

$$\bullet (\bar{x}_1 - \bar{x}_2) \sim N\left(\mu_1 - \mu_2, \left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)\right)$$

◆ 소표본인 경우 \rightarrow t분포로

$$\bullet (\bar{x}_1 - \bar{x}_2) \sim t\left(\mu_1 - \mu_2, \left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right); v'\right), \quad v' = \frac{\left[\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right]^2}{\frac{(s_1^2/n_1)}{n_1-1} + \frac{(s_2^2/n_2)}{n_2-1}}$$

◆ 실제 분석에서는 n 의 크기에 상관없이 t 분포로 처리 ($\because n$ 이 커지면 t 분포는 정규분포에 근사)

CASE 4: 대응표본(쌍표본)의 검정

◆ 대응표본(Paired sample): 특정 처리 전과 후의 데이터

● $H_0: \mu_D = 0$ $H_1: \mu_D \neq 0$

◆ 이를 검정하기 위해, 다음과 같은 n 개의 대응표본을 추출하였다.

$$(x_1, y_1) (x_2, y_2) (x_3, y_3) \dots (x_n, y_n)$$

◆ $d_i = x_i - y_i$ 로 정의하고 계산

● 대응표본간 차이 $\mu_D \sim t\left(\bar{d}, \sum_{i=1}^n \frac{(d_i - \bar{d})^2}{n-1}; n-1\right)$ 단, $\bar{d} = \sum_{i=1}^n d_i / n$

SPSS의 활용: 독립표본 T 검정

ID	Gender	BFM	FFM
1	1	15.40	36.70
2	2	19.10	42.80
3	3	8.80	35.60
4	4	14.90	31.20
5	5	13.60	40.80
6	6	17.60	39.40

SPSS의 활용: 독립표본 T 검정

집단통계량

Gender	N	평균	표준편차	평균의 표준오차
LBMD 0	191	-.142	1.3447	.0973
1	630	-.415	1.3058	.0520
FBMD 0	191	.640	.9523	.0689
1	630	-.144	.9655	.0385

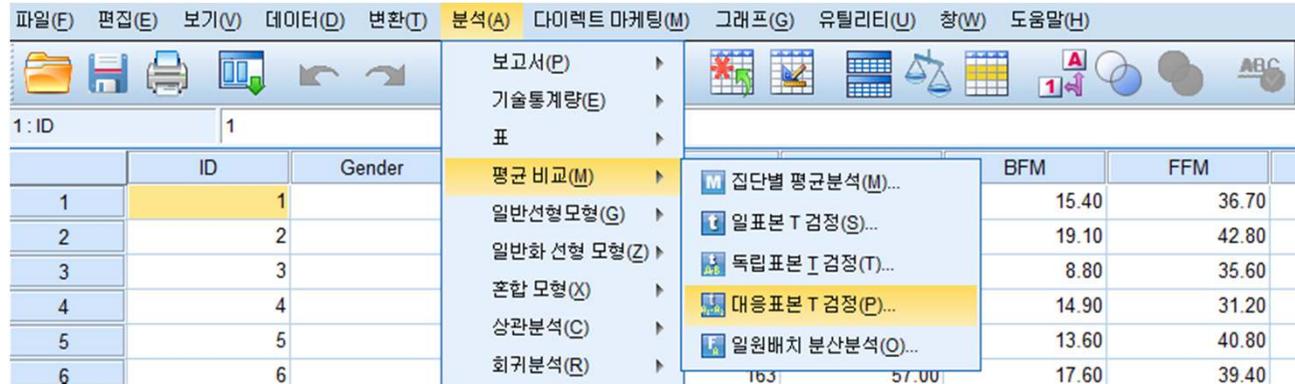
독립표본 검정

		Levene의 등분산 검정		평균의 동일성에 대한 t-검정						
		F	유의확률	t	자유도	유의확률 (양쪽)	평균차	차이의 표준오차	차이의 95% 신뢰구간	
									하한	상한
LBMD	등분산이 가정됨	.478	.490	2.508	819	.012	.2724	.1086	.0592	.4856
	등분산이 가정되지 않음			2.469	306.593	.014	.2724	.1103	.0553	.4895
FBMD	등분산이 가정됨	.270	.603	9.863	819	.000	.7841	.0795	.6281	.9402
	등분산이 가정되지 않음			9.936	317.554	.000	.7841	.0789	.6289	.9394

◆ Mean Comparisons between Male and Female Group

Variables	Gender		p-value
	Male (n=191)	Female (n=630)	
LBMD T-score	-0.142 ± 1.345	-0.415 ± 1.306	0.012
FBMD T-score	0.640 ± 0.952	-0.144 ± 0.966	0.000

SPSS의 활용: 대응표본 T 검정과 단일표본 T 검정



대응표본 통계량

	평균	N	표준편차	평균의 표준오차
대응 1 LBMD	-.351	821	1.3192	.0460
대응 1 FBMD	.039	821	1.0174	.0355

대응표본 상관계수

	N	상관계수	유의확률
대응 1 LBMD & FBMD	821	.656	.000

대응표본 검정

	대응차					t	자유도	유의확률 (양측)
	평균	표준편차	평균의 표준오차	차이의 95% 신뢰구간				
				하한	상한			
대응 1 LBMD - FBMD	-.3901	1.0077	.0352	-.4591	-.3210	-11.091	820	.000

비모수 검정: 독립표본과 대응표본

파일(F) 편집(E) 보기(V) 데이터(D) 변환(T) 분석(A) 다 Direkt 마케팅(M) 그래프(G) 유틸리티(U) 창(W) 도움말(H)

1: ID 1

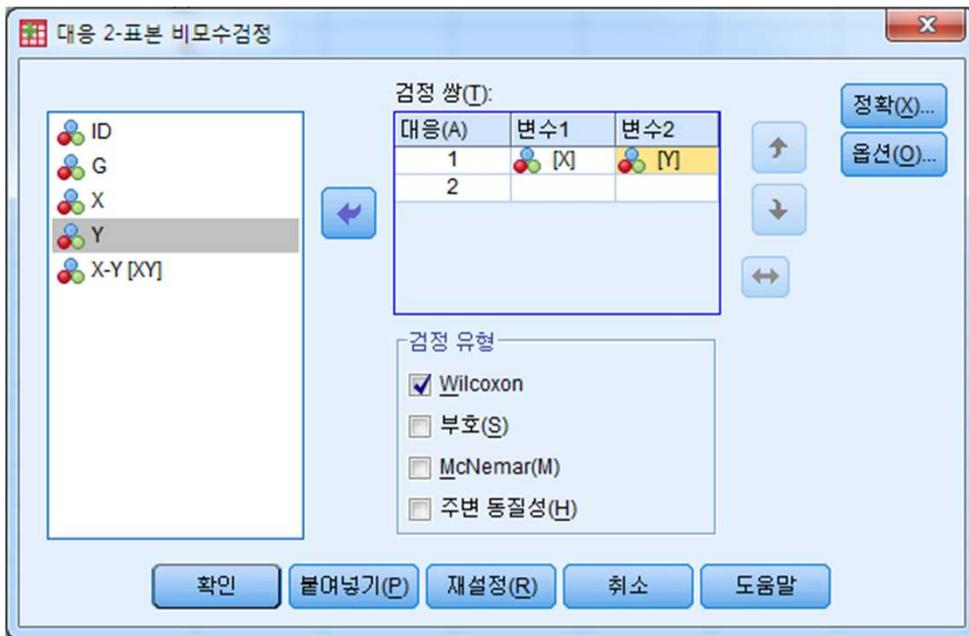
ID	G
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	
22	
23	

보고서(P) >
 기술통계량(E) >
 표 >
 평균 비교(M) >
 일반선형모형(G) >
 일반화 선형 모형(Z) >
 혼합 모형(X) >
 상관분석(C) >
 회귀분석(R) >
 로그선형분석(O) >
 신경망(W) >
 분류분석(Y) >
 차원 감소(D) >
 척도(A) >
비모수 검정(N) >
 예측(T) >
 생존확률(S) >
 다중응답(U) >
 결측값 분석(V)... >
 다중 대입(T) >
 복합 표본(L) >
 품질 관리(Q) >
 ROC 곡선(V)... >

Y	XY	변수	변수	변수
160	10			
130	15			
130	0			
165	5			
130	15			
130	5			
150	10			
135	5			
155	0			
170	5			
150	0			
0	0			
0	0			
0	0			
0	0			
160				
175				
140				
135				

▲ 일표본(O)...
 ▲ 독립 표본(I)...
 ▲ 대응 표본(R)...
레거시 대화 상자(L) >
 카이제곱검정(C)...
 이항(B)...
 런(R)...
 일표본 K-S(1)...
 독립 2-표본(2)...
 독립 K-표본(K)...
 대응 2-표본(L)...

대응표본과 비모수 검정: 일측슨 부호 순위합 검정



순위

		N	평균순위	순위합
Y - X	음의 순위	15 ^a	8.00	120.00
	양의 순위	0 ^b	.00	.00
	동률	5 ^c		
	합계	20		

- a. $Y < X$
- b. $Y > X$
- c. $Y = X$

검정 통계량^b

	Y - X
Z	-3.458 ^a
근사 유의확률(양측)	.001

- a. 양의 순위를 기준으로.
- b. Wilcoxon 부호순위 검정

독립표본과 비모수 검정: 맨-휘트니 U 검정



순위

	G	N	평균순위	순위합
X-Y	0	10	10.60	106.00
	1	10	10.40	104.00
합계		20		

검정 통계량^a

	X-Y
Mann-Whitney의 U	49.000
Wilcoxon의 W	104.000
Z	-.078
근사 유의확률(양측)	.938
정확한 유의확률 [2*(단측 유의확률)]	.971 ^a

a. 동렬에 대해 수정된 사항이 없습니다.

b. 집단변수: G