

# 통계이론 및 분석방법 소개

## 이론 중심으로

hmkang98@naver.com

와이즈인컴퍼니

# 통계자료분석

# 통계자료분석에 필요한 통계 개념

## 확률의 공리(probability axioms)

- $P(\Omega) = 1$
- $P(A) = P(A_1) + \cdots + P(A_n) \quad (A_i \cap A_j = \phi, \quad i \neq j)$
- $0 \leq P(A) \leq 1$

## 통계적 추론(statistical inference)

- 추정(estimation)
  - 신뢰구간(confidence interval)
- 검정(testing)
  - 가설 설정
  - 유의수준( $\alpha$ )
  - 기각역(critical region)
  - 검정통계량(test statistic)
  - 유의확률(p-value)

# 통계자료분석에 필요한 통계 개념

## 자료의 형태

- 수치형 자료(numerical data), 양적자료(quantitative data)
  - 연속형자료(continuous data) : 관측가능한 값이 연속적인 경우로 몸무게, 키 등이 있음
  - 이산형자료(discrete data) : 교통사고건수와 같이 관측가능한 값이 셀 수 있는 경우
- 범주형 자료(categorical data), 질적자료(qualitative data)
  - 순위형 자료(ordinal data) : 각 범주의 순서가 의미가 있는 경우로, 매우좋다, 좋다, 보통이다, 나쁘다, 매우 나쁘다 등과 같은 경우
  - 명목형 자료(nominal data) : 각 범주간 순서가 의미없는 경우로 남자, 여자 등과 같은 경우

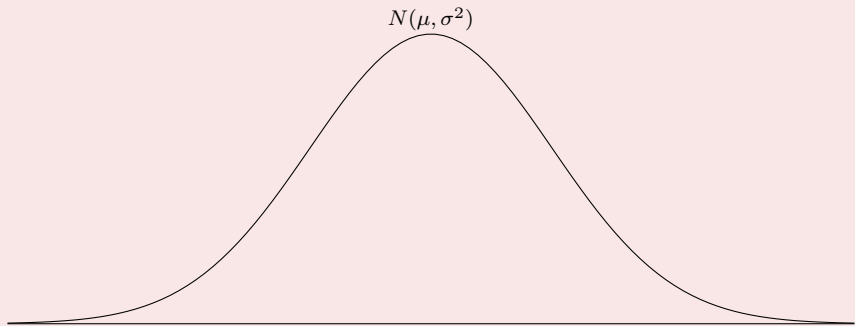
# 통계분포

# 정규분포

## 정규분포(normal distribution)

분포의 위치를 나타내는 평균  $\mu$ 와 분포의 크기를 나타내는 표준편차  $\sigma$ 로 표현하는 분포로 그래프의 형태는 종 모양(bell curve)이고 표현은  $N(\mu, \sigma^2)$ 이다.

## 정규분포 그림

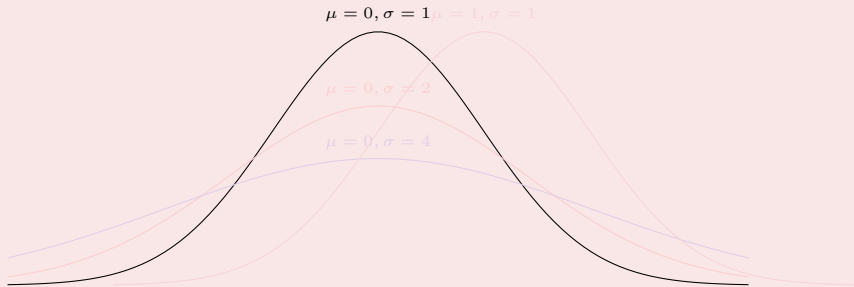


# 정규분포

## 정규분포의 성질

- ♣ 평균에 대하여 좌우로 대칭, 최빈값=중앙값=평균
- ♣ 평균  $\mu$ 는 중심위치를 나타내고, 표준편차  $\sigma$ 는 평균  $\mu$ 로부터 퍼져있는 정도를 나타냄

## 평균 $\mu$ 와 표준편차 $\sigma$ 가 다른 정규분포 그림

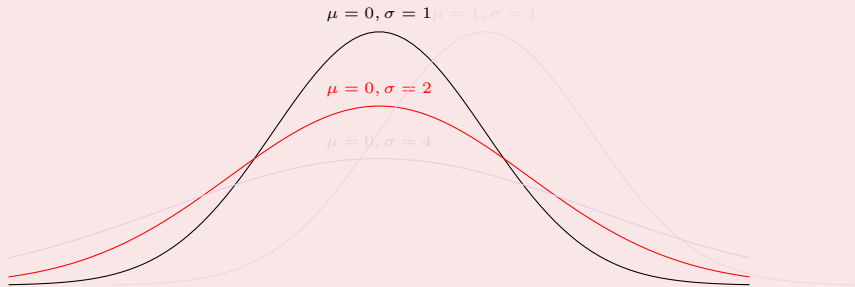


# 정규분포

## 정규분포의 성질

- ♣ 평균에 대하여 좌우로 대칭, 최빈값=중앙값=평균
- ♣ 평균  $\mu$ 는 중심위치를 나타내고, 표준편차  $\sigma$ 는 평균  $\mu$ 로부터 퍼져있는 정도를 나타냄

## 평균 $\mu$ 와 표준편차 $\sigma$ 가 다른 정규분포 그림



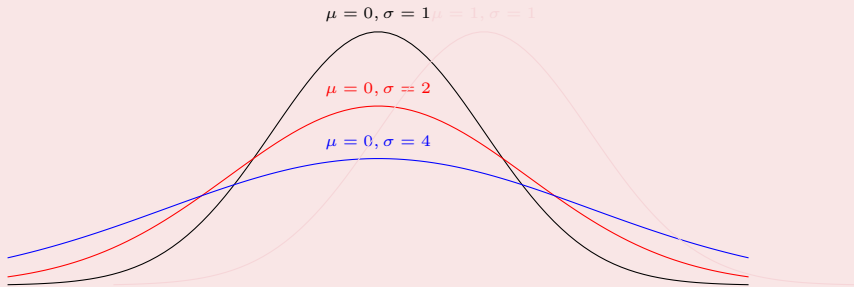


# 정규분포

## 정규분포의 성질

- ♣ 평균에 대하여 좌우로 대칭, 최빈값=중앙값=평균
- ♣ 평균  $\mu$ 는 중심위치를 나타내고, 표준편차  $\sigma$ 는 평균  $\mu$ 로부터 퍼져있는 정도를 나타냄

## 평균 $\mu$ 와 표준편차 $\sigma$ 가 다른 정규분포 그림

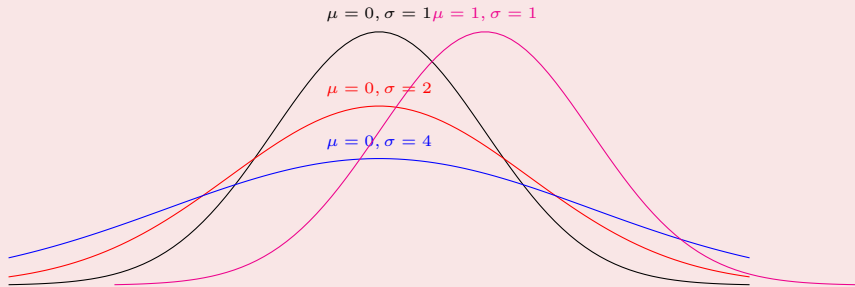


# 정규분포

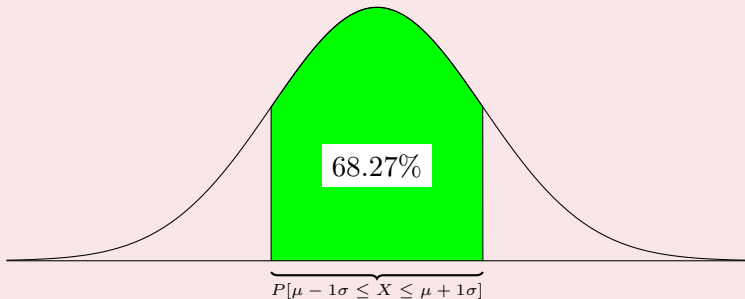
## 정규분포의 성질

- ♣ 평균에 대하여 좌우로 대칭, 최빈값=중앙값=평균
- ♣ 평균  $\mu$ 는 중심위치를 나타내고, 표준편차  $\sigma$ 는 평균  $\mu$ 로부터 퍼져있는 정도를 나타냄

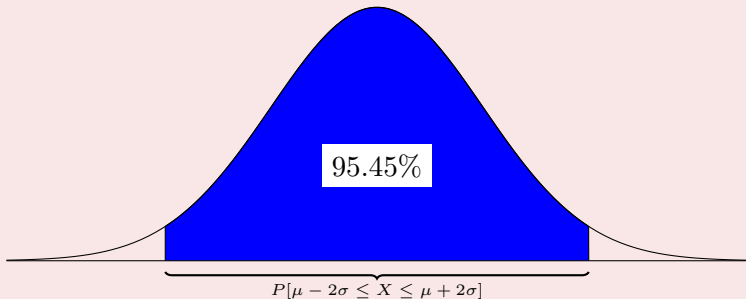
## 평균 $\mu$ 와 표준편차 $\sigma$ 가 다른 정규분포 그림



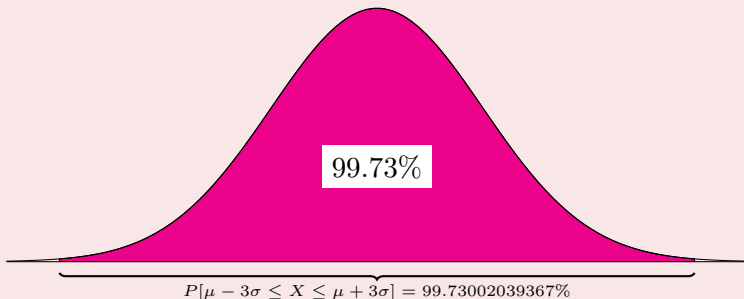
확률변수  $X$ 가 정규분포를 따를 때 많이 알려진 확률



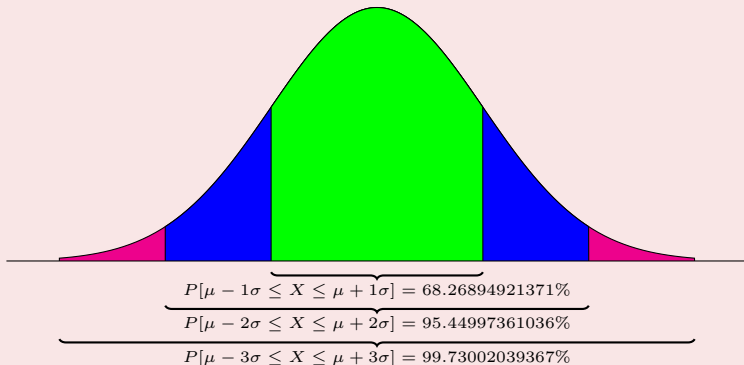
확률변수  $X$ 가 정규분포를 따를 때 많이 알려진 확률



확률변수  $X$ 가 정규분포를 따를 때 많이 알려진 확률



확률변수  $X$ 가 정규분포를 따를 때 많이 알려진 확률



## 확률변수 $X$ 가 정규분포를 따를 때 많이 알려진 확률

♣  $P[\mu - 1\sigma \leq X \leq \mu + 1\sigma] = 68.26894921371\%$

♣  $P[\mu - 2\sigma \leq X \leq \mu + 2\sigma] = 95.44997361036\%$

♣  $P[\mu - 3\sigma \leq X \leq \mu + 3\sigma] = 99.73002039367\%$

♣  $P[\mu - 4\sigma \leq X \leq \mu + 4\sigma] = 99.99366575163\%$

♣  $P[\mu - 5\sigma \leq X \leq \mu + 5\sigma] = 99.99994266969\%$

♣  $P[\mu - 6\sigma \leq X \leq \mu + 6\sigma] = 99.99999980268\%$

♣  $P[\mu - 7\sigma \leq X \leq \mu + 7\sigma] = 99.99999999974\%$

## t-분포 예

&lt;표-2&gt; 신생아의 주수별 변인의 기술 통계 및 t-검정 결과

		N	mean	SD	SE	t
체중	full-term	35	3.15	.34	0.057	.000*
	preterm	33	2.50	.41	0.071	
		68	2.83	.49	0.060	
두위	full-term	35	33.37	1.48	0.250	.004*
	preterm	33	32.09	2.03	0.354	
		68	32.75	1.87	0.227	
흉위	full-term	35	31.97	1.63	0.275	.000*
	preterm	33	29.62	2.38	0.415	
		68	30.83	2.34	0.283	
신장	full-term	35	47.94	2.21	0.374	.001*
	preterm	33	45.50	3.25	0.567	
		68	46.76	3.01	0.365	
심박수	full-term	35	129.26	8.40	1.420	.040*
	preterm	33	134.30	11.10	1.932	
		68	131.71	10.06	1.220	
호흡수	full-term	35	47.20	4.48	0.757	.946
	preterm	33	47.27	4.32	0.753	

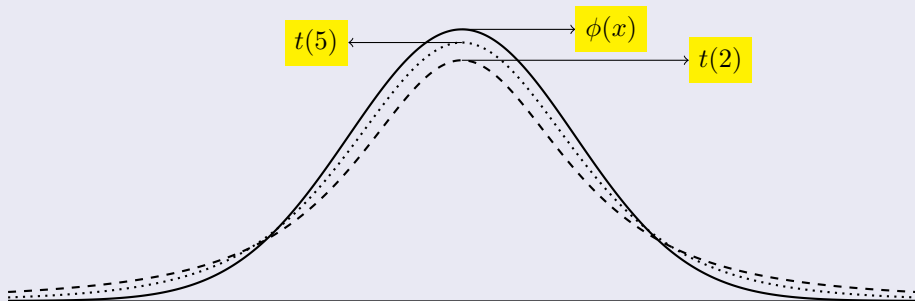


## t - 분포 소개

- 1908년 영국의 과학자 고셋(W. C. Gosset)이 Biometrika에 소개
- 이 논문은 필명인 Student라는 이름으로 출간
- 모집단의 분포가 정규분포이고 표본의 크기가 작은 경우 표본평균  $\bar{X}$ 의 분포
- 모수는 자유도(degree of freedom)이며 모양은 종 형태이고, 0에서 대칭인 분포
- 자유도가 커지면 정규분포에 가깝게 됨
- $$f(t) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi\nu}\Gamma(\frac{\nu}{2})} \cdot \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}} \quad \text{for } -\infty < t < \infty$$

# 표준정규분포와 t - 분포 비교

그림으로 표준정규분포와 t - 분포 비교



# 표준정규분포와 $t$ - 분포 비교

## 애니메이션으로 표준정규분포와 $t$ - 분포 비교

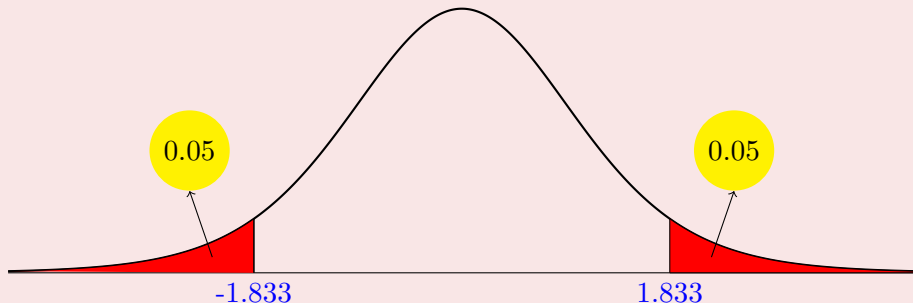
그림: 애니메이션으로 정규분포와  $t$ -분포 비교

## 보기 : t - 분포표(그래프와 함께)

**문제** : 자유도가 9인 t 분포를 갖는 통계량 t에 대하여

$P(-b \leq t \leq b) = 0.9$ 를 만족시키는 b를 찾아라.

**풀이** : t 분포는 0을 중심으로 대칭이므로  $(-\infty, -b)$ 와  $(b, \infty)$ 에 각각 5%의 확률이 존재한다. 자유도가 9인 t 분포에서 상위 5%의 확률인 제 95% 백분위수를 구하면 1.833이 된다. 즉,  $t_{0.05}(9) = 1.833$ 이 된다.



# $\chi^2$ 분포 ( $\chi^2$ distribution)

## $\chi^2$ 분포 ( $\chi^2$ distribution) 소개

- $\chi^2$  분포는 1900년 Karl Pearson이 제안한 분포
- 감마분포의 특별한 분포
- 감마분포

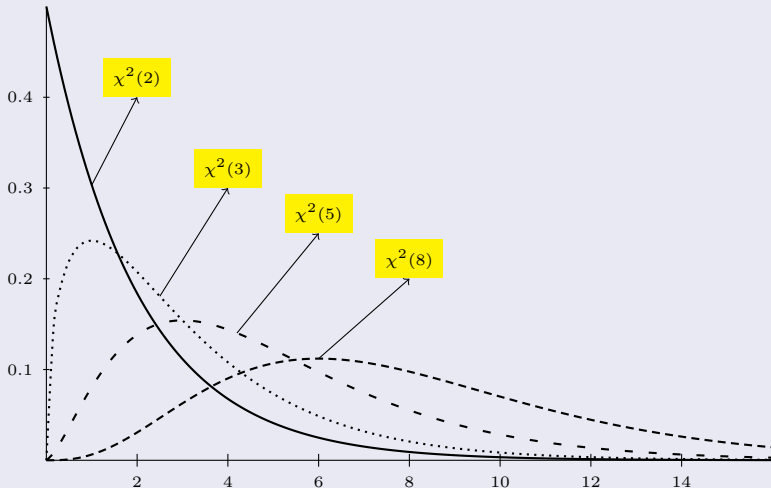
$$Gamma(\alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}$$

- $\chi^2$  분포

$$Gamma(x, \alpha = v/2, \beta = 2) = \frac{1}{\Gamma(v/2)2^{v/2}} x^{\frac{v}{2}-1} e^{-x/2}$$

- $\chi^2$  분포는  $\chi^2$  검정과 분산 추론에 사용
- 모수(parameter)는 자유도가 있고, 모양(shape)은 비대칭형

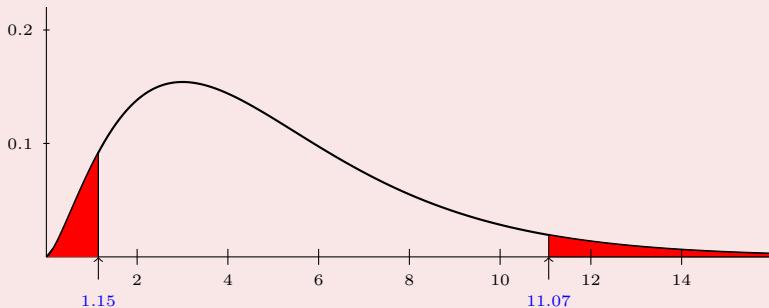
## 그림으로 본 자유도에 따른 $\chi^2$ 분포들



## 보기 : $\chi^2$ 분포표

**문제** :  $\chi^2$  분포표에서 자유도가 5인 상, 하위 5%의 확률을 주는 값을 찾아라.

**풀이** :  $\chi^2$  분포표를 보면  $\chi_{0.05}^2(5) = 1.15$ 이고  $\chi_{0.95}^2(5) = 11.07$ 이다. 이 값들을 그래프에 표현하면 다음과 같다



## F-분포 예

표 9. 과제, 성별, 교육수준에 따른 발화당 단어수의 반복측정 분산분석

	제곱합	자유도	F
개체 내			
과제	1727.449	1,817	216.598***
과제 * 성별	46.040	1,817	5.773**
과제 * 교육수준	1.875	3,634	.118
과제 * 교육수준 * 성별	10.205	3,634	.640
개체 간			
성별	37.949	1	3.535
교육수준	101.887	2	4.745**
성별 * 교육수준	.707	2	.033

\*\*  $p < .01$ , \*\*\*  $p < .001$ .

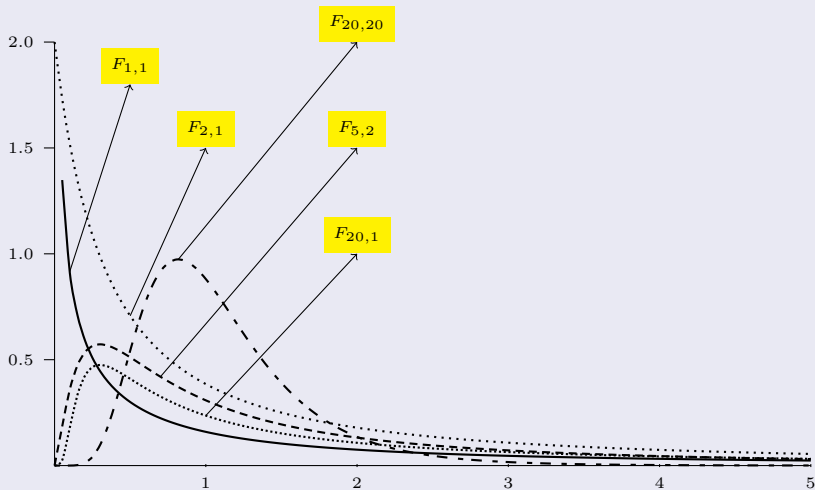


## F - 분포의 특징

- 통계학자 피셔(R. A. Fisher)가 제안한 확률분포
- 세 집단의 평균 비교에 사용
- 두 집단의 분산비율 추론에 사용(한 집단 분산 추론은  $\chi^2$ 분포)
- $$f(x; \nu_1, \nu_2) = \frac{\Gamma\left(\frac{\nu_1 + \nu_2}{2}\right)}{\Gamma\left(\frac{\nu_1}{2}\right)\Gamma\left(\frac{\nu_2}{2}\right)} \left(\frac{\nu_1}{\nu_2}\right)^{\frac{\nu_1}{2}} \frac{x^{\frac{\nu_1 - 2}{2}}}{\left[1 + \left(\frac{\nu_1}{\nu_2}\right)x\right]^{\frac{\nu_1 + \nu_2}{2}}}$$
- F - 분포는 서로 독립인  $\chi_{\nu_1}^2, \chi_{\nu_2}^2$  일때  $F = \frac{\chi_{\nu_1}^2/\nu_1}{\chi_{\nu_2}^2/\nu_2}$
- 카이제곱과 마찬가지로 양의 구간에서만 확률값을 갖는다.
- 분포의 모양은 비대칭형이다.
- $F_{\alpha, \nu_1, \nu_2}$ 는  $1/F_{1-\alpha, \nu_2, \nu_1}$ 과 같다.
- t 분포와 F 분포의 관계 :  $t_{\alpha/2}^2(n) = F_{\alpha, 1, n}$ ,  $t^2 = \frac{p(n-1)}{n-p} F_{p, n}$

# F - 분포

## 그림으로 본 자유도에 따른 F분포들



## 애니메이션으로 F-분포 변화 보기

그림: 애니메이션으로 F-분포 변화 보기( $\alpha = 0.05$ )

# 통계적 추론

## 어느 개인병원에 방문하는 환자수가

- \* 10일 동안 살펴본 결과 하루 평균 200명이었다.
- \* 내원하는 환자수의 오차한계가 30명이었다면 이 병원에 내원하는 환자수는 170명에서 230명이라고 예측할 수 있겠다.

## 다음 달 이 병원에 내원한 하루 평균 환자수는

- \* 220명이었다면 환자수가 증가했다고 할 수 있는가?
- \* 250명이었다면 환자수가 증가했다고 할 수 있는가?

## 위의 결과로

- \* 내원 환자수는 하루 평균 200명인데 오차가 있으니 220명은 있을 수 있는 경우이다.
- \* 내원 환자수는 하루 평균 200명인데 250명은 환자수가 증가했다고 할 수 있겠다.

# 표본에서 표본평균을 구하는 과정

## 표본 추출 예

연구자가 병원에 내원한 전체 인원 7000 명에서 100 명을 임의로 뽑아 전체 환자의 몸무게 평균( $\mu$ )을 예측한다고 하자.

- ① 대부분의 경우 100 명의 조사는 1 회만 한다.
- ② 실제로 7000 명에서 100 명을 뽑는 회수는  $\binom{7000}{100} \approx 7.98 \times 10^{82}$  이다.
- ③ 위의 개수 만큼 표본 평균을 구할 수 있고 표본평균의 평균과 분산 등을 계산할 수 있다.
- ④ 연구자가 1회 조사하여 얻은 표본평균( $\bar{X}$ )은 수많은 표본평균 중 하나를 구한 것이다.

# 추정

# 모평균에 $\mu$ 대한 추정

## 하나의 값으로 예측하는 점추정(point estimation)

- 관심이 있는 전체집단의 평균을 모집단의 평균  $\mu$ 라고 부르고 하나의 값으로 예측
- 모집단의 평균  $\mu$ 은 알 수 없으니
- 평균이  $\mu$ , 표준편차가  $\sigma$ 인 모집단에서 임의로 표본  $X_1, \dots, X_n$ 을 뽑아
- 표본평균  $\bar{X}$ 을 계산하여 전체집단의 평균  $\mu$ 라고 예측한다.
- 모집단의 평균  $\mu$ 와 일부집단인 표본집단의 평균  $\bar{X}$ 는 같지 않을 것이다.
- 모집단의 특성을 나타내는 값을 모수(parameter), 표본집단의 특성을 나타내난 값을 통계량(statistic)이라고 한다.



# 모평균에 $\mu$ 대한 추정

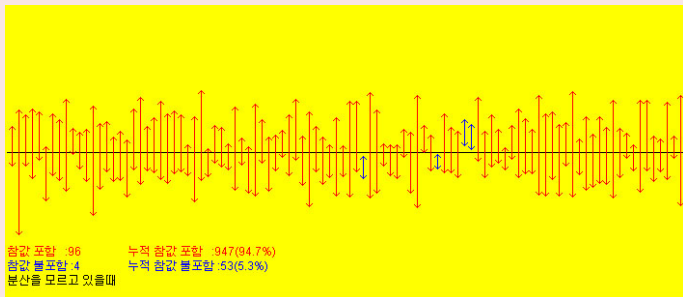
## 구간추정(interval estimation)

- **구간추정(interval estimation)** : 표본평균의 분포를 사용하여 표본으로부터 모집단의 평균이 포함될 것이라 예상되는 구간을 추정
- **신뢰구간(confidence interval)** : 구간추정에서 제시하는 구간
- **신뢰구간**은 (하한값, 상한값)의 형태로 구성
- 모집단에서 추출한 표본마다 계산되는 **신뢰구간**은 서로 다를 수 있음
- 가장 확실한 신뢰구간 :  $(-\infty, \infty) \rightarrow$  어떤 정보도 제공하지 못함
- **신뢰수준(confidence level)** : 모수가 **신뢰구간**에 포함될 확률로 보통 90%, 95%, 99%를 사용
- **신뢰수준** 또는 **신뢰도**는  $100(1 - \alpha)\%$  또는  $1 - \alpha$ 로 표시
- 모집단의 평균  $\mu$ 에 대한  $100(1 - \alpha)\%$  **신뢰구간**은 표본평균  $\pm$  표본오차(sampling error)로 나타낸다.

# 모평균에 $\mu$ 대한 추정

## 신뢰구간의 의미

- 평균  $\mu$ 가 0인 모집단에서 표본을 뽑아 신뢰구간을 계산한 결과
- 전체 신뢰구간 중 모평균을 포함하는 신뢰구간은  $947/1000 = 0.947$
- 신뢰구간은 모집단의 평균이 하나의 신뢰구간에 포함될 확률이  $100(1 - \alpha)\%$ 가 아니라 신뢰구간을 계속하여 계산할 경우 모집단의 평균이 포함되는 신뢰구간이  $100(1 - \alpha)\%$ 에 가깝게 된다.



신뢰도 : 95    표본의 크기 : 4    다시 그리기

# 검정

# 모평균 $\mu$ 에 대한 검정

## 가설검정이 필요한 한 가지 예

어느 도시의 보건당국에서 여러 성인병을 유발하는 높은 콜레스테롤 수치를 낮춤으로 그 도시의 의료비용을 절감하려고 한다. 그래서 이 도시는 콜레스테롤 수치를 줄이는 캠페인을 1년간 대대적으로 벌였다. 이 캠페인이 성인의 콜레스테롤 양을 줄이는데 효과가 있었는지 검증하고자 그 도시의 성인 40명을 대상으로 콜레스테롤 수치를 측정하였다. 캠페인을 시작할 때, 이 도시 성인의 콜레스테롤 수치는 평균이 200(mg/dl)이고 표준편차는 24(mg/dl)인 분포였다고 알려져 있다.

캠페인을 진행한 후 성인 40명을 뽑아 콜레스테롤 수치의 평균( $\bar{X}$ )을 계산하였다. 이 경우 성인의 콜레스테롤 양이 줄었는지 판단하려면 어떻게 해야 하는가?

- 1 콜레스테롤 수치가 낮아졌는지 판단하기 위하여 모든 성인을 조사하는 것은 불가능
- 2 캠페인 후 그 도시의 성인 중  $n$ 명 표본추출
- 3 캠페인 후 성인의 콜레스테롤 수치 모평균  $\mu$ 가 200mg/dl보다 작다고 할 수 있는지 통계적으로 판단
- 4  $\bar{X}$ 가 얼마나 작아야  $\mu$ 가 200mg/dl보다 작다고 주장할 수 있을까?

# 모평균 $\mu$ 에 대한 검정

## 가설검정(testing statistical hypothesis)

모평균에 대한 가설이 적합한지 추출한 표본으로 판단

### 가설

- **귀무가설(null hypothesis,  $H_0$ )** : 대립가설의 반대 가설로 가설검정에서 기준이 되는 가설
- **대립가설(alternate hypothesis,  $H_1$ )** : 입증하여 주장하고자 하는 가설

### 가설설정 적용사례

- \*  $H_0$  : 도시의 성인 콜레스테롤 수치 평균  $\mu$ 는 200mg/dl이다.
- \*  $H_1$  : 도시의 성인 콜레스테롤 수치 평균  $\mu$ 는 200mg/dl보다 작다.

# 모평균 $\mu$ 에 대한 검정

## 오류의 종류

실제의 상태 \ 검정의 결론	$H_0$ 기각 못함 ( $H_1$ 채택 안됨)	$H_0$ 기각 ( $H_1$ 채택)
	옳은 결론 제 2종 오류 확률	제 1종 오류 확률 옳은 결론
$H_0$ 참 $H_0$ 거짓		

- **제 1종의 오류 확률(type I error probability,  $\alpha$ )** : 실제 상태는 귀무가설( $H_0$ )이 옳지만 잘못 판단하여 귀무가설( $H_0$ )을 기각하는 오류 확률
- **제 2종의 오류 확률(type II error probability,  $\beta$ )** : 실제 상태는 귀무가설( $H_0$ )이 옳지 않지만 잘못 판단하여 귀무가설( $H_0$ )을 기각하지 않는 오류 확률

# 모평균 $\mu$ 에 대한 검정

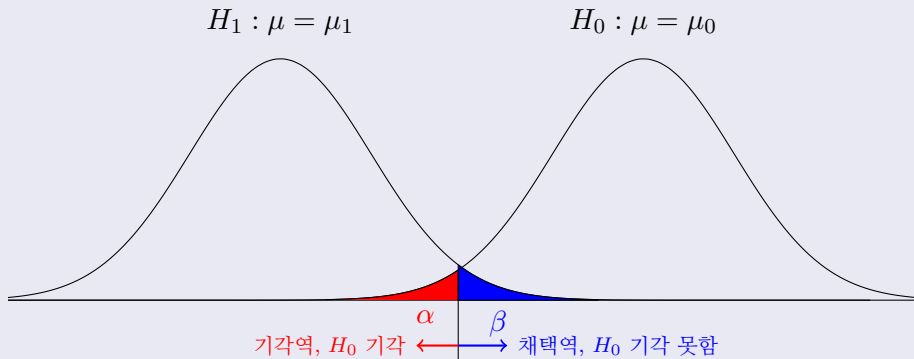
## 오류의 적용사례

실제의 상태 \ 검정의 결론	$\mu = 200\text{mg/dl}$	$\mu < 200\text{mg/dl}$
	옳은 결론 제 2종 오류	제 1종 오류 옳은 결론
$\mu = 200\text{mg/dl}$		
$\mu < 200\text{mg/dl}$		

- ✱ 제 1종의 오류( $\alpha$ ) : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200임에도 불구하고 잘못 예측하여 200보다 작다고 판단하는 오류
- ✱ 제 2종의 오류( $\beta$ ) : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200보다 작음에도 불구하고 잘못하여 200으로 판단하는 오류

# 모평균 $\mu$ 에 대한 검정

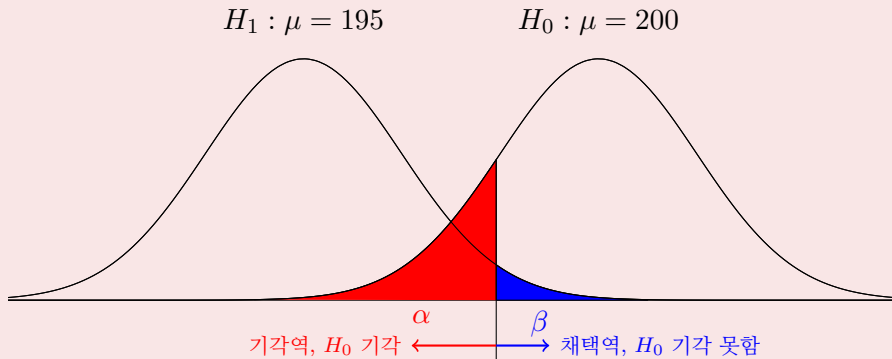
## 모평균 $\mu$ 에 대한 검정 - $\alpha, \beta$





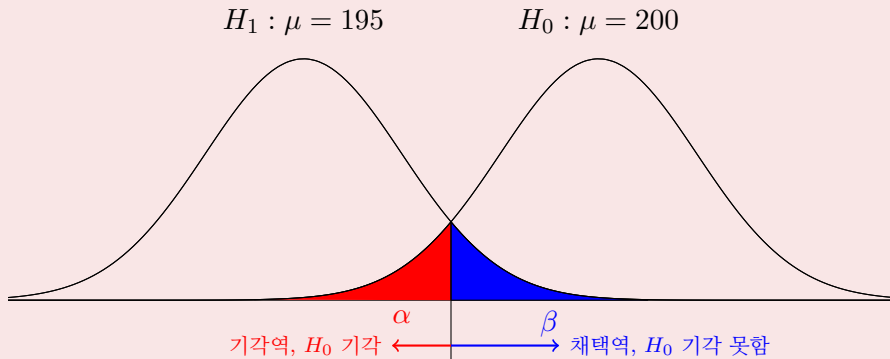
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - $\alpha, \beta$



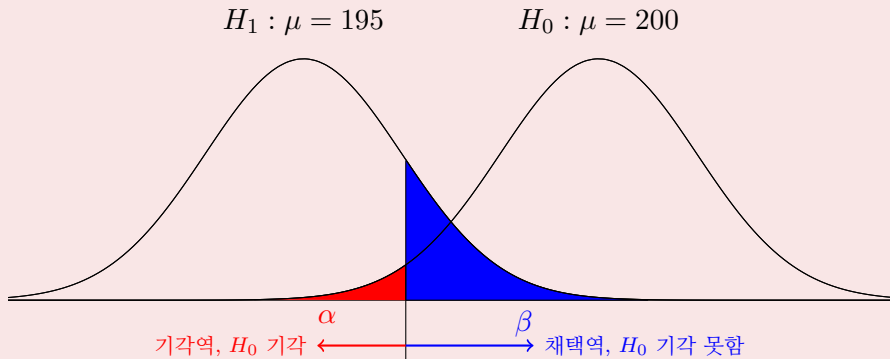
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - $\alpha, \beta$



# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - $\alpha, \beta$



# 모평균 $\mu$ 에 대한 검정

## 검정에 사용하는 용어들

- **유의수준(significance level,  $\alpha$ )** :  $H_0$ 가 옳다고 검정의 결론을 내렸으나 잘못 판단하여  $H_0$ 를 기각하게 되는 오류의 최대허용한계로 많은 경우 (1%, 5%, 10%) 등을 유의수준으로 사용하며 연구자가 스스로 판단하여 허용오류를 결정함
- **기각역(critical region)** :  $\bar{X}$ 가 취하는 구간 중에서  $H_0$ 를 기각하게 되는 구간이며 유의수준  $\alpha$ 로 계산(표현은  $R: \bar{X} \leq c$ 의 형태)
- **검정통계량(test statistic)** : 모집단의 일부분인 표본으로부터 검정의 결론( $H_0$ 를 기각하거나,  $H_0$ 를 기각하지 못하거나)을 내리는데 사용하는 통계량(예 콜레스테롤 양의 평균  $\bar{X}$ ,  $\frac{\bar{X}-\mu_0}{\sigma/\sqrt{n}}$ )
- **유의확률(p-값, p-value)** : 유의수준은 연구자가 직접 결정하나 유의확률은 표본에서 계산된 검정통계량의 관측치에서  $H_0$ 를 기각하는 최소의 유의수준으로 관측값의 유의수준으로 부르기도 함

# 모평균 $\mu$ 에 대한 검정

## 검정에 사용하는 용어들

- **유의수준(significance level,  $\alpha$ )** : 연구자가 일어날 수 있는 경우인지, 드물게 일어나는 경우인지 판단하기 위한 그 기준을 설정한다. 묵시적으로 5%가 표준이다.
- **기각역(critical region)** : 유의수준이 설정되면 드물게 일어나는 경우의 영역을 구한다.
- **검정통계량(test statistic)** : 표본평균
- **유의확률(p-값, p-value)** : 표본평균이 귀무가설  $H_0$  집단에서 발생할 누적확률이다.

# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 통계적 가설검정 단계

- ❶ 귀무가설( $H_0$ )과 대립가설( $H_1$ )을 설정한다. 두 가설은 서로 상반된 가설로 대부분 연구자 또는 조사자는 주장하고자 하는 대립가설( $H_1$ )이 의미있도록 하기 위하여 귀무가설( $H_0$ )을 기각하기 원한다.
- ❷ 유의수준  $\alpha$ 에 대하여 귀무가설( $H_0$ )을 기각하는 영역인 기각역을 설정한다.
- ❸ 검정통계량과 유의확률을 계산한다.
- ❹ 검정통계량이 기각역에 속하면(유의확률이 유의수준보다 작으면) 귀무가설( $H_0$ )을 기각한다.

# 모평균 $\mu$ 에 대한 검정

## 가설 설정

**귀무가설( $H_0$ )**은 기존에 연구 결과 알려진 사실이므로 특정한 분포이다. 따라서 모수(parameter)가 알려져 있으며 이 값에 대한 분포가 연구 결과 주장하려는 **대립가설( $H_1$ )**이 통계적으로 유의한지 결정한다.

모평균  $\mu$ 에 대한 가설검정에서

- **귀무가설( $H_0$ )**은  $H_0 : \mu = \mu_0$

일 때 **대립가설( $H_1$ )**은 세 가지로 설정할 수 있다. 또한 검정하는 방법에 따라 단측검정(one sided test)과 양측검정(two sided test)로 나눈다.

- $H_1 : \mu > \mu_0 \rightarrow$  one sided test
- $H_1 : \mu < \mu_0 \rightarrow$  one sided test
- $H_1 : \mu \neq \mu_0$  즉  $H_1 : \mu > \mu_0$  또는  $H_1 : \mu < \mu_0 \rightarrow$  two sided test

가설검정은 유의확률로 판단한다. 그러나 실제로 가설이 무엇인지 모르고 판단하는 경우가 있다. 자료에서 분석결과를 얻기 전에 가설이 무엇인지 반드시 알고 있어야 한다.

# 모평균 $\mu$ 에 대한 검정

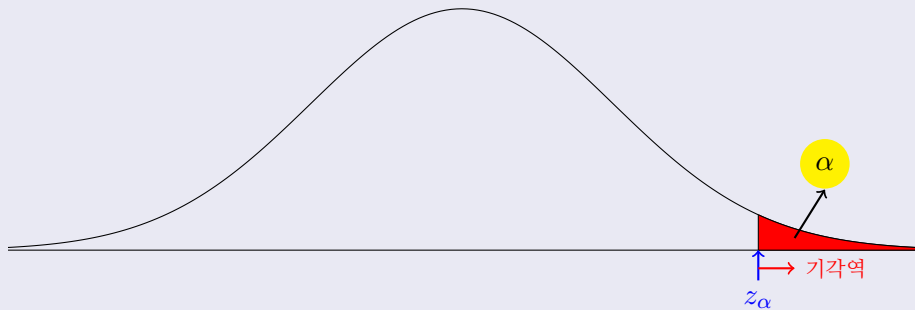
## 기각역(reject area)과 유의수준(significance level, $\alpha$ )

- 1 연구자는 스스로 귀무가설을 기각하는 유의수준을 설정
- 2 연구자가 결정한 유의수준에 대한 기각역을 계산
- 3 기각역은 대립가설의 종류에 따라 세 가지로 설정할 수 있음
  - $H_1 : \mu > \mu_0$  일 때  $R : \bar{X} > c$  이거나  $R : \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = Z > z_{\alpha/2}$
  - $H_1 : \mu < \mu_0$  일 때  $R : \bar{X} < c$  이거나  $R : \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = Z < -z_{\alpha/2}$
  - $H_1 : \mu \neq \mu_0$  일 때  $R : |\bar{X}| > c$ 로  $R : \bar{X} > c$  or  $\bar{X} < -c$  이거나  $R : \left| \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right| = |Z| > z_{\alpha/2}$ 로  $R : Z > z_{\alpha/2}$  or  $Z < -z_{\alpha/2}$
- 4 실제로 분석에서 기각역은 관심의 대상이 아니고, 유의수준은 묵시적으로 0.05이며, “분석결과 유의확률(p-value)이 0.05보다 작으면 의미있는 결과가 나온 것이다.” 라고 표현하는 것이 이 부분이다.



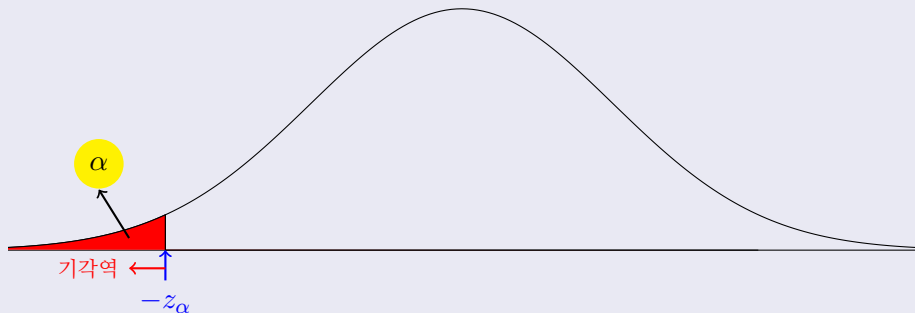
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 단측검정( $H_1 : \mu > \mu_0$ )



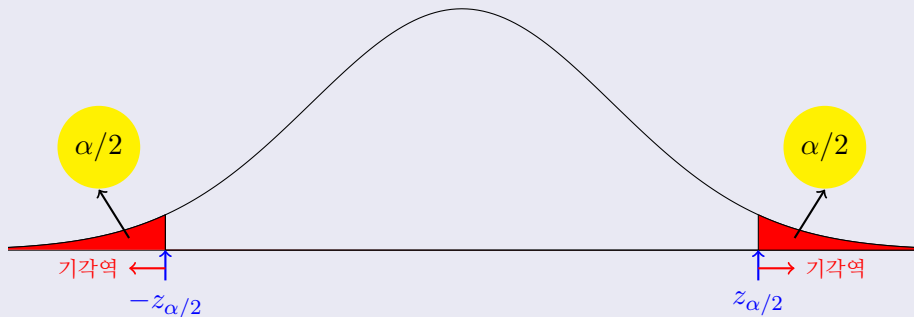
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 – 단측검정( $H_1 : \mu < \mu_0$ )



# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 – 양측검정( $H_1 : \mu \neq \mu_0$ )



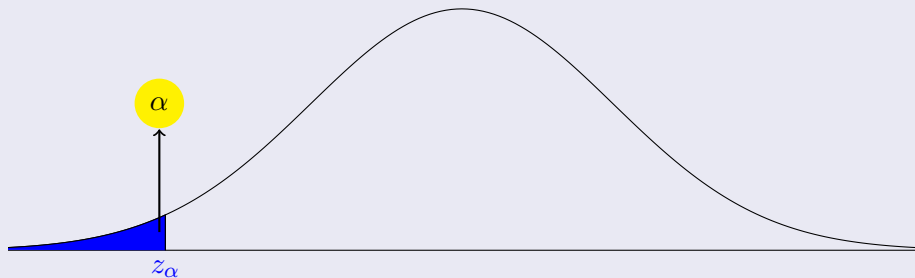
# 모평균 $\mu$ 에 대한 검정

## 유의확률 계산

- 유의확률은 세 개의 대립가설에 대하여 기각역의 부등호 방향과 동일하게 검정통계량의 관측값으로 계산
- 양측검정인 경우, 검정통계량의 관측값에서 계산한 결과에 2를 곱함
- 유의확률은 손으로 계산하기가 쉽지 않으므로 컴퓨터를 이용
- 대부분 통계프로그램이 양측검정 결과를 제시하기 때문에, 단측검정인 경우에는 2로 나눈 값이 유의확률 값임
- $H_1 : \mu > \mu_0$  일 때  $P[\bar{X} > \bar{x}]$  이거나  $P[\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} = Z > z_0]$
- $H_1 : \mu < \mu_0$  일 때  $P[\bar{X} < \bar{x}]$  이거나  $P[\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} = Z < -z_0]$
- $H_1 : \mu \neq \mu_0$  일 때  $P[\bar{X} > \bar{x}] \times 2$  or  $P[\bar{X} < \bar{x}] \times 2$  이거나  $P[Z > z_0] \times 2$  or  $P[Z < -z_0] \times 2$

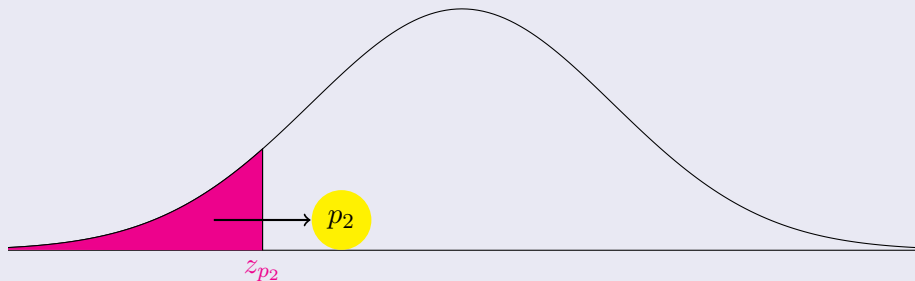
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 유의수준( $\alpha$ )과 유의확률(p-value)



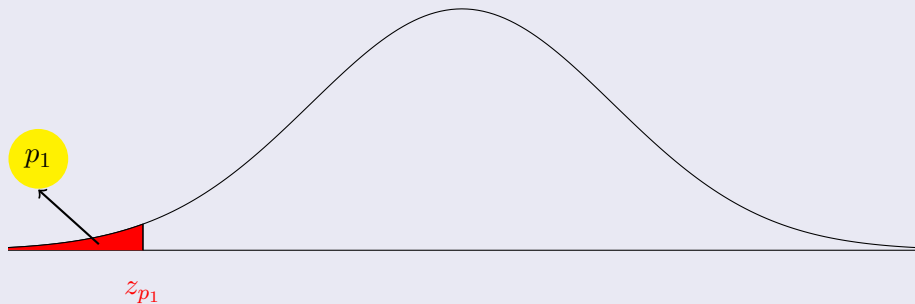
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 유의수준( $\alpha$ )과 유의확률(p-value)



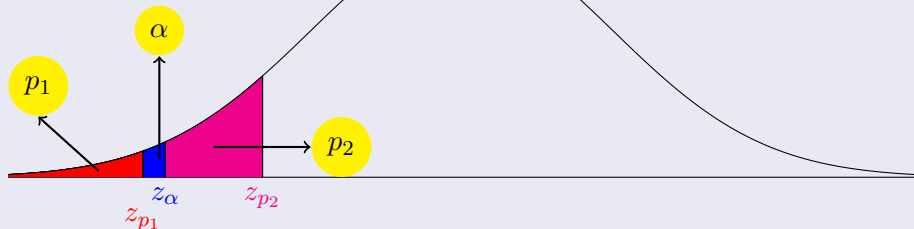
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 유의수준( $\alpha$ )과 유의확률(p-value)



# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 유의수준( $\alpha$ )과 유의확률(p-value)





# 모평균 $\mu$ 에 대한 검정

유의수준( $\alpha$ ) > 유의확률( $p$ -값)이면 귀무가설  $H_0$  기각

# 모평균 $\mu$ 에 대한 검정

## 통계적 가설검정 예제

### ① 가설 설정

- \*  $H_0$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl이다.
- \*  $H_1$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl보다 작다.

### ② 유의확률, 유의수준, 검정통계량, 기각역

- \* 유의수준  $\alpha = 0.05$  (사용자가 설정)
- \* 기각역  $R$  :  $Z \leq -1.645$  or  $\bar{X} \leq 193.76$  ( $\alpha = 0.05, \sigma = 24, n = 40$ )
- \* 검정통계량 : ①  $\bar{X} = 195$ , ②  $\bar{X} = 192$  (표본에서 계산)
- \* 유의확률 : ①  $P[\bar{X} \leq 195] = 0.0938$ , ②  $P[\bar{X} \leq 192] = 0.0175$

### ③ 결론 : 유의수준 $\alpha = 0.05$ 에서

- ① 귀무가설( $H_0$ )을 기각 못함
- ② 귀무가설( $H_0$ )을 기각

# 모평균 $\mu$ 에 대한 검정

## 통계적 가설검정 예제

### ❶ 가설 설정

- ※  $H_0$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl이다.
- ※  $H_1$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl보다 작다.

### ❷ 유의확률, 유의수준, 검정통계량, 기각역

- ※ 유의수준  $\alpha = 0.05$  (사용자가 설정)
- ※ 기각역  $R : Z \leq -1.645$  or  $\bar{X} \leq 193.76 (\alpha = 0.05, \sigma = 24, n = 40)$

$$\begin{aligned} 0.05 &= P[Z \leq -1.645] = P\left(\frac{\bar{X} - 200}{24/\sqrt{40}} \leq -1.645\right) \\ &= P[\bar{X} \leq 200 - 1.645 \times 24/\sqrt{40}] = P[\bar{X} \leq 193.76] \end{aligned}$$

- ※ 검정통계량 : ❶  $\bar{X} = 195$ , ❷  $\bar{X} = 192$ (표본에서 계산)
- ※ 유의확률 : ❶  $P[\bar{X} \leq 195] = 0.0938$ , ❷  $P[\bar{X} \leq 192] = 0.0175$

### ❸ 결론 : 유의수준 $\alpha = 0.05$ 에서

- ❶ 귀무가설( $H_0$ )을 기각 못함
- ❷ 귀무가설( $H_0$ )을 기각

# 모평균 $\mu$ 에 대한 검정

## 통계적 가설검정 예제

### ① 가설 설정

- \*  $H_0$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl이다.
- \*  $H_1$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl보다 작다.

### ② 유의확률, 유의수준, 검정통계량, 기각역

- \* 유의수준  $\alpha = 0.05$  (사용자가 설정)
- \* 기각역  $R : Z \leq -1.645$  or  $\bar{X} \leq 193.76 (\alpha = 0.05, \sigma = 24, n = 40)$
- \* 검정통계량 : ①  $\bar{X} = 195$ , ②  $\bar{X} = 192$  (표본에서 계산)
- \* 유의확률 : ①  $P[\bar{X} \leq 195] = 0.0938$ , ②  $P[\bar{X} \leq 192] = 0.0175$

$$\begin{aligned} &= P[\bar{X} \leq 195] = P\left(\frac{\bar{X} - 200}{24/\sqrt{40}} \leq \frac{195 - 200}{24/\sqrt{40}}\right) \\ &= P[Z \leq -1.3176] = 0.0938 \end{aligned}$$

### ③ 결론 : 유의수준 $\alpha = 0.05$ 에서

- ① 귀무가설( $H_0$ )을 기각 못함
- ② 귀무가설( $H_0$ )을 기각

# 모평균 $\mu$ 에 대한 검정

## 통계적 가설검정 예제

### ① 가설 설정

- \*  $H_0$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl이다.
- \*  $H_1$  : 성인 콜레스테롤 수치의 평균  $\mu$ 는 200mg/dl보다 작다.

### ② 유의확률, 유의수준, 검정통계량, 기각역

- \* 유의수준  $\alpha = 0.05$  (사용자가 설정)
- \* 기각역  $R : Z \leq -1.645$  or  $\bar{X} \leq 193.76 (\alpha = 0.05, \sigma = 24, n = 40)$
- \* 검정통계량 : ①  $\bar{X} = 195$ , ②  $\bar{X} = 192$  (표본에서 계산)
- \* 유의확률 : ①  $P[\bar{X} \leq 195] = 0.0938$ , ②  $P[\bar{X} \leq 192] = 0.0175$

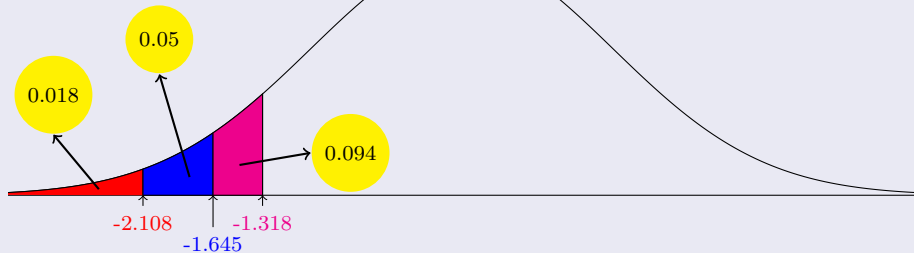
$$\begin{aligned} &= P[\bar{X} \leq 192] = P\left(\frac{\bar{X} - 200}{24/\sqrt{40}} \leq \frac{192 - 200}{24/\sqrt{40}}\right) \\ &= P[Z \leq -2.1082] = 0.0175 \end{aligned}$$

### ③ 결론 : 유의수준 $\alpha = 0.05$ 에서

- ① 귀무가설( $H_0$ )을 기각 못함
- ② 귀무가설( $H_0$ )을 기각

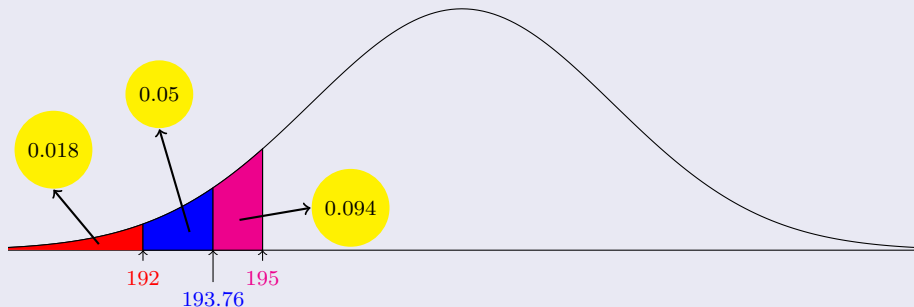
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 유의확률과 유의수준(표준화한 경우)



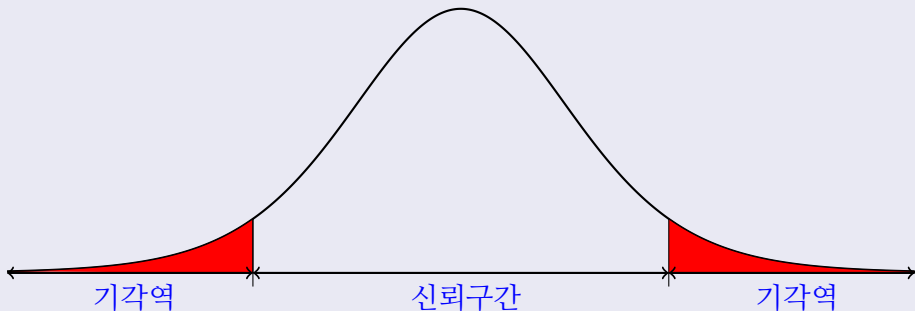
# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 유의확률과 유의수준



# 신뢰구간과 양측검정

## 신뢰구간과 양측검정의 관계





# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정( $n \geq 30$ ) 예제

어느 다이어트 방법을 소개하는 책자에서 주장하기를 그 다이어트 방법을 이용하면 5주 동안 10kg 넘게 체중을 줄일 수 있다고 한다. 그 다이어트 방법을 이용한 56명을 대상으로 5주 동안의 체중 감소량을 조사하였더니 평균이 10.5kg, 표준편차가 4.5kg이었다고 한다. 이 자료에 근거하여 그 책자의 주장이 옳다고 할 수 있는지 유의수준 5%로 검정하라.

① 가설 설정 :  $H_0 : \mu = 10 \quad vs \quad H_1 : \mu > 10$

② 통계량 또는 유의확률 계산 :

$$Z = \frac{\bar{X} - \mu_0}{s/\sqrt{n}} = \frac{\bar{X} - 10}{s/\sqrt{56}}, \quad R : Z \geq z_\alpha = z_{0.05} = 1.645, \quad z = \frac{10.5 - 10}{4.5/\sqrt{56}} = 0.83$$
$$P\text{-value} = P[Z \geq z] = P[Z \geq 0.83] = 0.2033$$

③ 결론 : 검정통계량의 관측값이 기각역에 포함되지 않기 때문에 귀무가설을 기각할 수 없다. 또는 유의확률이 유의수준보다 크기 때문에 귀무가설을 기각할 수 없다.

# 모평균 $\mu$ 에 대한 검정

## 모평균 $\mu$ 에 대한 검정 - 예제를 그래프로

