

문제 1 Central Limit Theorem

중심극한정리는 확률변수의 합 형태(sum of random variables)의 극한분포를 손쉽게 구할 수 있도록 해주기에 통계학에서 가장 자주 사용하는 정리입니다. 이 문제에서는 중심극한정리의 정의와 그 활용에 대해 짚어보겠습니다.

(1-1) 📖 중심극한정리(Central Limit Theorem)의 정의와 그 의미를 서술시오.

📖 통계학입문(3판) 7장 참고

📖 Hogg(8판) 4장 2절, 5장 3절 참고

Def: 평균이 μ 이고 분산이 σ^2 인 독립인 랜덤한 서로 동질인 확률변수 X_1, X_2, \dots, X_n 이 존재할 때, $\frac{\sum (X_i - \mu)}{\sigma}$ 은 n 이 ∞ 에 가까워질수록 $N(0, 1)$ 인 표준정규분포를 따른다.

$$\left(\begin{array}{l} \text{As } n \rightarrow \infty \\ \bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right) \end{array} \right)$$

Meaning: 표본의 크기 n 이 충분히 크다면, 표집량이 어떤 분포를 갖고있든지 $N(\mu, \sigma^2/n)$ 을 따른다.

→ 따라서 서로 다른 표본의 통계량을 통하여 $N(\mu, \sigma^2/n)$ 을 가정하고 추정할 수
있다는 점에 의해서 오는 정리이다.

(1-2) 📖 중심극한정리가 통계적 추론 중 “구간추정”에서 어떻게 유용한지 서술하시오.

📖 Hogg(8판) 4장 2절 참고

(1-3) 📖 중심극한정리를 이용하여 모평균에 대한 근사신뢰구간을 만들 때, 표준오차($\sqrt{\text{Var}(\bar{X})}$) 부분의 모분산을 표본분산으로 대체할 수 있는 이유를 수식적으로 증명하시오.

📖 표본분산 s^2 는 모분산 σ^2 로 확률수렴한다는 사실을 이용할 수 있습니다.

📖 Slutsky's theorem을 이용할 수 있습니다.

📖 Hogg(8판) 5장 1~3절 참고

(1-2) **모평균** | 정규분포를 따르지 않아도 신뢰구간 확정이 가능하다.

CLT에 의해 표본이 크면 표본평균의 분포는 정규분포를 따르므로, $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$ 을 이용해 근사신뢰구간을 구할 수 있다.

예) 신뢰도 95%로 μ (모평균)을 추정하면 $\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ 로 구간확정이 가능하다.
(모산을 모를 경우 S)

$$\begin{aligned} (1-3) \quad S_n^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2 = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}_n^2 \right) \\ &= \frac{n}{n-1} (E(x_i^2) - \mu^2) \end{aligned}$$

$$\lim_{n \rightarrow \infty} S_n^2 = \lim_{n \rightarrow \infty} \frac{n}{n-1} (E(x_i^2) - \mu^2) = E(x_i^2) - \mu^2 = \sigma^2$$

As $n \rightarrow \infty$,

\therefore 표본분산 S_n^2 은 모분산 σ^2 에 수렴한다.

따라서 CLT가 성립하면 표본분산 S^2 이 σ^2 에 확률수렴하므로 대치하여 사용가능하다.

문제 2 Student's Theorem

스튜던트 정리는 통계적 추정에서 필요한 정리 중 하나로, 표본평균과 표본분산이 어떤 분포를 갖는지 알려줍니다. 이 문제에서는 스튜던트 정리의 내용을 어떻게 수식적으로 유도할 수 있는지 알아보겠습니다.

스튜던트 정리는 다음과 같이 총 4개의 내용으로 구성되어 있습니다.

- ① $\bar{X} \approx N(\mu, \frac{\sigma^2}{n}) \rightarrow$ <문제 1>에서 증명함
- ② 표본평균 \bar{X} 와 표본분산 s^2 은 서로 독립이다.
- ③ ???
- ④ ???

(2-1) ③의 내용을 쓰고 증명하시오.

Hogg(8판) 3장 6절 참고

무작위표본 X_1, \dots, X_n 이 독립적으로 동일하게(independently and identically distributed) 평균이 μ 이고 분산이 σ^2 인 정규분포를 따를 때, 자유도가 n 인 카이제곱분포를 따르는 새로운 확률변수 V 를 아래와 같이 두어 증명에 활용할 수 있습니다.

$$V = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2(n)$$

$$\textcircled{3} \quad \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

$$\text{Proof: } Z^2 = \left(\frac{X - \mu}{\sigma} \right)^2, \text{ so, } V = \sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 = \sum_{i=1}^n Z_i^2$$

$$V = \sum_{i=1}^n \left(\frac{X_i - \bar{X} + \bar{X} - \mu}{\sigma} \right)^2 = \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma} \right)^2 + n \left(\frac{\bar{X} - \mu}{\sigma} \right)^2 = \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma} \right)^2 + \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2$$

$$= S^2(n-1) \times \frac{1}{\sigma^2} + \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2$$

$$\text{by ①, } \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right) \sim N(0, 1) \text{ (표준화)이므로 } \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2 \sim \chi^2(1)$$

$$\text{by ②, } \frac{S^2(n-1)}{\sigma^2} \text{과 } \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2 \text{ 독립이며 } \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2 \text{가 자유도와 1인 카이제곱분포를}$$

따라서 V 는 자유도와 n 인 카이제곱 분포를 따른다. $\frac{(n-1)S^2}{\sigma^2}$ 은 자유도가 $(n-1)$ 인 카이제곱분포를 따른다.

(2-2) ④의 내용을 쓰고, (2-1)을 이용하여 증명하시오.

Hogg(8판) 3장 6절 참고

t분포의 정의에 따르면, 표준정규분포를 따르는 확률변수와 카이제곱분포를 따르는 확률변수를 이용하여 t분포를 유도할 수 있습니다.

$$\textcircled{4} \quad T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$$

$$\text{proof: } T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \cdot \frac{S}{S} = \frac{\bar{X} - \mu}{S/\sqrt{n}} \cdot \frac{\sqrt{(n-1)S^2}}{\sqrt{(n-1)S^2}}$$

$$\text{by ① } \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim N(0, 1) \quad , \quad \text{by 2-1 } \frac{S^2(n-1)}{\sigma^2} \sim \chi^2(n-1)$$

$$T \text{ 분포의 정의: } \uparrow \quad W, V \text{ 독립, } W \sim N(0, 1), \quad V \sim \chi^2(r)$$

$$T = \frac{W}{\sqrt{V/r}} \sim t(r)$$

$$\text{따라서 } \frac{\bar{X} - \mu}{S/\sqrt{n}} = W, \quad \frac{(n-1)S^2}{\sigma^2} = V \text{ 라 하면, } r = n-1 \text{ 이 되며}$$

$$T = \frac{\frac{\bar{X} - \mu}{S/\sqrt{n}}}{\sqrt{\frac{(n-1)S^2}{\sigma^2} \cdot \frac{1}{n-1}}} = \frac{W}{\sqrt{V/r}} \sim t(r) \quad (= t(n-1))$$

문제 3 t-test

t검정은 모집단이 정규분포를 따르지만 모표준편차를 모를 때, 모평균에 대한 가설검정 방법입니다. 대개 두 집단의 모평균이 서로 차이가 있는지 파악하고자 할 때 사용하며, 표본평균의 차이와 표준편차의 비율을 확인하여 통계적 결론을 도출합니다.

(3-1) 어떤 학우가 DSL 학회원(동문 포함)의 평균 키가 DSL 학회원이 아닌 사람의 평균 키보다 크다는 주장을 하여, 실제로 그러한지 통계적 검정을 수행하려고 합니다. 며칠간 표본을 수집한 결과 다음의 결과를 얻었다고 합니다.

표본 수 : 각 101명

측정에 응한 DSL 학회원들의 평균 키 : 178.5cm / 표준편차 : 7.05cm

측정에 응한, DSL 학회원이 아닌 사람들의 평균 키 : 179.9cm / 표준편차 : 7.05cm

(a) 귀무가설과 대립가설을 설정하시오. $\text{Let DSL 평균키} = \mu_A, \text{ not DSL 평균키} = \mu_B$
 $\text{DSL 표준편차} = \sigma_A, \text{ not DSL 표준편차} = \sigma_B$

$$H_0: \mu_A - \mu_B = 0 \quad H_1: \mu_A - \mu_B > 0$$

(b) 유의수준 5%에서의 가설검정을 수행하고 결론을 도출하시오.

통계학입문(3판) 7장 참고

어떤 검정통계량이 어떤 분포를 따르는지, 언제 귀무가설을 기각하는지 정해야 합니다.

가설검정: $\bar{X}_A = 178.5, \bar{X}_B = 179.9, n_A = 101, n_B = 101, S_A = 7.05, S_B = 7.05, \alpha = 0.05$
 by CLT, 모분산 unknown. 표본수 미지이므로 z 검정 사용. $\sigma_A^2 = S_A^2, \sigma_B^2 = S_B^2$ 로 모분산 대체
 검정통계량 $z = \frac{(\bar{X}_A - \bar{X}_B) - (\mu_A - \mu_B)}{\sqrt{\hat{\sigma}_A^2/n_A + \hat{\sigma}_B^2/n_B}} = -1.411167$
 $|z| < z_{0.05} = 1.645$ 이므로 reject H_0 . $-1.411167 < 1.645$ 이므로 not reject H_0

결론:
 DSL 학회원들의 평균 키가
 DSL이 아닌 사람들의 평균 키와
 크게 차이가 없다.

(3-2) 신촌 연세로를 지나는 버스 노선의 이용객 수가 '차 없는 거리 해제(2022.10.09.)' 이후 유의미하게 증가했는지 파악하기 위해, 우선 2022년 9월의 평균 이용객 수와 2022년 11월의 평균 이용객 수가 유의미한 차이를 보이는지 통계적 검정을 수행하려고 합니다. <elementary1.ipynb>

(a) 귀무가설과 대립가설을 설정하시오. $\text{Let 해제전 평균 이용객수} = \mu_A, \text{ 해제후 평균 이용객수} = \mu_B$

$$H_0: \mu_A - \mu_B = 0 \quad H_1: \mu_A - \mu_B < 0$$

(b) 파이썬 scipy의 stats 패키지를 활용하여 유의수준 5%에서의 가설검정을 수행하고 결론을 도출하시오.