

QnA.

1. A3Rsu vs. A3Rmi

이 위치에는 행과 ϵ 의 첫번째 행과 ϵ 의 첫번째 원소가 더해진다.

A3Rsu: $E(X'\epsilon) = 0 \rightarrow X = \begin{pmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{pmatrix}$ $\epsilon = \begin{pmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{pmatrix}$ $\rightarrow X'\epsilon = \begin{pmatrix} x_{11}\epsilon_1 + x_{12}\epsilon_2 + \dots + x_{1n}\epsilon_n \\ \vdots \\ x_{n1}\epsilon_1 + x_{n2}\epsilon_2 + \dots + x_{nn}\epsilon_n \end{pmatrix}$

* 파란색은 행렬곱 방향을 의미.

$X' = \begin{pmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{pmatrix}$

= $\begin{pmatrix} \sum_j x_{1j}\epsilon_j \\ \vdots \\ \sum_j x_{nj}\epsilon_j \end{pmatrix}$: Data 상에서 같은 행의 곱의 합

= $\epsilon_1 x_{11} + \epsilon_2 x_{12} + \dots + \epsilon_n x_{1n}$ (본 배며.)
(ϵ 들이 가중치 역할)

- $E(X'\epsilon) = 0$ 은 Data 상에서 같은 행의 x 와 ϵ 의 곱의 합이 0이 나와야 하며 이는 x 와 ϵ 간에 covariance 혹은 correlation이 없음을 의미하기도 함.
대신에 Data 수가 커지고 $\frac{1}{n}$ 로 나눠줘야 함.
- A3Rmi는 어떤 Data(X)가 와도 ϵ 은 평균적으로 $E(\epsilon|X) = 0$ 이 되는 것을 의미하기에 소표본에서도 x 와 ϵ 간 관계가 없음을 생각하면 된다.
- A3Rsu는 대표본에서 삼관이 있다고 받아들일 수 있다.

2. Non indep 부분. $\rightarrow P.11 \quad Var(\bar{x}) = \frac{\sigma^2}{n} \left(\frac{N-n}{N-1} \right)$

* 초기 분포를 따르는 변수의 분산 $V(X) = n \cdot \frac{k}{M} \cdot \frac{M-k}{M} \cdot \frac{M-n}{M-1}$
 $n \cdot p \cdot (1-p)$

- 변수들의 sample 간 iid가 깨져서 $Var(X_i) + Var(X_j)$ 뿐만 아니라 $2Cov(X_i, X_j)$ 가 생기고 Variance가 커진다. (참고로 Cov는 음수가 나올 수 있으나 이전 샘플의 성격을 거의 유지할 것이므로 양수로 볼 수 있음)

- $n \ll N$ 일 때 Correction factor가 1에 가까워진다는 것은 수학에서도 알 수 있지만 (검정치는 $n \rightarrow 1$)

위 논리를 바탕으로 해석할 때는 Cov가 생길 표본 수가 줄어들어 분산의 크기가 많이 증가하지 않기 때문으로 생각할 수 있다.

- Fancier CLT? $x_i \sim iid$ If $x_i \sim \text{not iid}$, Not fancier.

\rightarrow Lindberg-Levy's CLT

Not fancier CLT: L-F CLT

: 이항 분포의 분산 공산
indep이 깨졌으므로 수정됨.

Question 6.

문제의 의미: Regression의 해석 순서

→ Regression 혹은 Statistics의 경합. ① Predictive ② Analytic

예. Why? 예측은 하변한측 Data의 특성과 불확실성이 커지기 때문.
한국의 하계팀 측면에서 특정 고객층이나 요인에 대한 상관성
혹은 인과성을 고려할 수 있음.

* 복습: 1방의 R²부분 / 저 1층 외곽 2층 요인 ($\alpha \uparrow (1-\beta) \downarrow$ 무슨 의미?)

Multicollinearity: $\frac{1}{1 - \text{corr}(X_i, X_j)}$ 은 두 변수의 관계성에 따른 조판. Data 자체의 Multicollinearity를 확인하고 싶으면 $\text{Def}(X)$ 를 확인하자.

* Multicollinearity가 아닐 때 X의 변동이 낮으면 ($X \rightarrow \text{Constant}$: 변동이 낮음) $(X'X) \downarrow (X'X)^{-1} \uparrow V(\beta) = \sigma^2(X'X)^{-1} \uparrow t\text{-value} = \frac{\hat{\beta} - \beta_0}{\text{se}(\hat{\beta})} \downarrow$

X의 2nd moment를 고려해봐.

단축 검정 하면.

* 단축 검정이 양측 검정보다 좋은 이유? 변동을 알고 있을 때, $H_0: \beta = 0$ 이라는 가설을 확률이 높아지니까.

양측 검정을 한다면 이는 사실상 0의 확률을 $\frac{\alpha}{2}$ 로 두는 것과 동등하다.

Question 7.

문제의 의미: Dummy variable의 취와 Interaction term의 의미.

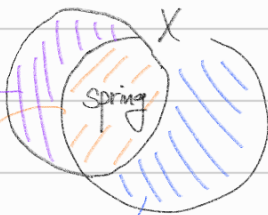
/ test나 CI를 구하기 위해서는 분산(2nd moment)이 필요하다. → 통계학은 다 평균과 분산 예이다? (대부분 결과 분산에 기반한다는 얘기)

* Dummy variable interaction.

ex) X: 2개 방음장. sp: 봄 dummy 변수

Y: 옷 판매량 (dependant)

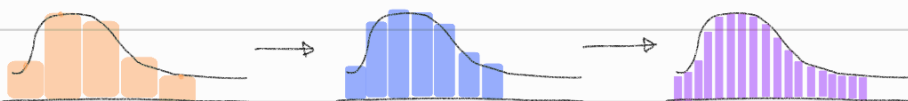
$$Y = \beta_0 + \beta_1 X + \beta_2 sp + \beta_3 X \cdot sp$$



: 변수 X에서 sp가 같이 실행하는 비중을 따진 편미분은 것이라고 생각하자.

* 복습: 분산이 0이 취되는 이유는? 오차 항에 다시 보아 혹은 오차와 강의 자료.

* 이산형이 Data 수가 커질수록 연속형으로 수렴?



시험 내용.

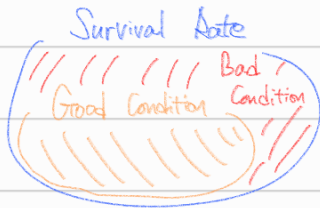
6. Categorical Data.

6.1 Categorical input with categorical input.

• odds ratio: odds 4/2. odd? 전체 분의 비율이 아닌(Risk) 상대 위험의 수에 대한 비율 이구나? 컨트롤 한에 대해 알아보기.

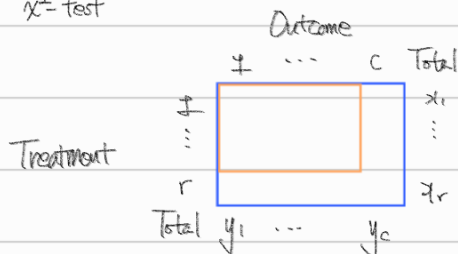
6.1.1 Simpson's paradox. (Question 7)

* Categorical Variable 밖에 고려할 변수가 남아있다는 것 \rightarrow Multivariate Regression.



6.1.2 χ^2 -test

• 자유도:



$\chi^2((r-1)(c-1))$ 인 이유는

각 행과 열의 합이 정해져 있기 때문에
경우의 수가 1씩 줄어들기 때문이다.

6.2 ANOVA

6.2.3 Interpretations and assumptions.

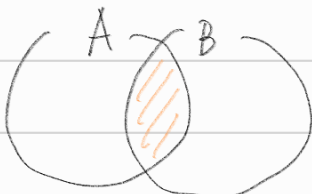
• 모든 Data의 2nd moment인 분산이 동일하다는 것에 하미 (Homogeneity) 2nd moment가 동일하다는 것.
 \rightarrow 분산을 보겠다는 것.

교차 살펴보기.

6.2.4 Some extensions of ANOVA (Question 7번)

• $SSA + SSB + SSAB$

\hookrightarrow Interaction term이 추가 혹은 설명하는 2nd moment를 의미.



6.2.6 Kruskal-Wallis : the non-parametric version.

• Wilcoxon signed rank test / Mann-Whitney U test : comparing medians

↓ ANOVA.

Kruskal-Wallis