

2023 0412

/ 이진 데이터란?

- 브리스톨 병원 예 (PPDAC)  $\Rightarrow$  Problem: - 브리스톨 병원의 어린이 심장 수술 사망률 다른 병원과 비교해서 현저히 높음?
  - 어떤이런?
  - 심장 수술이란?
  - 사망의 기준?

plan: - 심장 수술 database  $\rightarrow$  빅2

Data: HES, CSR 이라는 자료 사용 but, 재복수 사건과 사망연속자 다중

Analysis: HES, CSR 2위 차주로 사용  $\rightarrow$  사망률 예측

Conclusion: 예측 기준 컷점이지만 실제 사망자는 02명

- binary data (이진 데이터): 즉 가치 값으로 다뤄진 데이터. (ex. 1=사망, 0=생존)
  - 평점은 비율로 표시된다. (ex. 사망률)

- '결과'만 중요한 것으로 전한 (ex. 사망률  $\Rightarrow$  긍정 메시지 프레임  
생존률  $\Rightarrow$  긍정 메시지 프레임)

$\therefore$  긍정 프레임, 절제-생존전 요약 모두 활용하여 올바르게 정보 전달 해야함!

- 막대 그래프의 기준점을 고려하여 해석해야함

## 2. 범주형 자료의 소개

- 범주형 자료란?: - Variable: 주어진 상황에 따라 다른 값을 가지는 특징치

$\hookrightarrow$  왜 다양한 범주들 값을 가지는 변수

- 성 X 범주: 국적, 성별
- 성 O 범주: 종교
- 일련 2중 범주: BMI 기준 비만 정도

\* 파생차분 결과 사용해야 마라!!

$\therefore$  숫자를 외야만 안 수 있는 상황과 자료임.  $\rightarrow$  상황과 능력이 필요  
 $\rightarrow$  이런 막대 그래프를 사용하라.

- 상대 위험도  $\longleftrightarrow$  절대 위험도  
 위험도가 있는 절대 위험도 6%, 7%

Contro group 절대 위험도

$\therefore$  상대 위험도가 높더라도 절대 위험도가 작을 경우에는 큰 문제가 아닐 수 있다.

- 기댓값: 주어진 집합에서 특정 사건에 일어날 개략 예측값

- 예: 30%에서 많이 사용됨 (ex. 한국 월드컵 우승 확률의  $\frac{1}{2}$ 인 경우  $\frac{1/30}{1-1/30} \approx \frac{1}{30} \frac{\text{우승의 경우 수}}{\text{우승 X 경우 수}}$ )

- 예 =  $\frac{\text{위험도가 있는 집합의 예}}{\text{Control group의 예}}$  # 절대 위험도와 작은 경우 예비 = 상대 위험도

### 3. 연속형 자료 분석

- 범형

- SSI 0

- SSI X

- 수형

- 이산형: 값이 정수

- 연속형: 값이 실수

- 연속형 자료의 시각화

- strip-chart

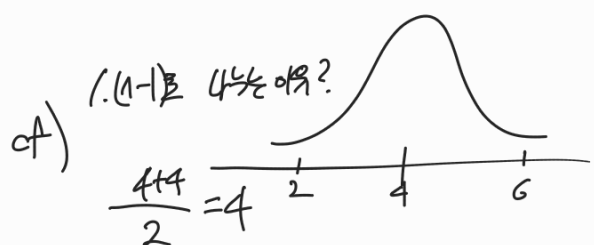
- box plot

- 히스토그램

- 데리비 차분법과 32 변환

- 자료의 수리 분석

- 중앙
  - 평균  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$
  - 중위값
  - 최빈수
- 포괄
  - 범위
  - IQR
  - 표준편차  $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$
  - 분산



- 자료도

2. 제곱근을 절대값으로 사용하면?

- 제곱근을 절대값에 사용되어 있음

- 아웃라이어: 한 개의 데이터 값의 변태 크게 작거나 큰 통계량  
ex. IQR, 중앙값

- strip chart (dot plot): 자료들이 겹칠수록 점의 색이 진해짐

- 히스토그램: 구간 나누어 그림  $\therefore$  구간의 개수, 크기, 시작점이 중요함.

- Modality  $\longleftrightarrow$  Skewness  
 $\rightarrow$  봉우리 몇 개냐  $\rightarrow$  왜칭력이나  
ex. 히스토그램