

# 5강 선형회귀

공분산과 상관관계는 두 데이터 간의 관계를 밝히기 위함이다.

단순 숫자가 아니라 진짜 방정식을 구하고 싶은 것.

두 데이터 간에 어떤 관계가 있는지 정확히 눈으로 나타내고 싶다.

그래서 한 변수가 변할 때,

다른 변수가 정확히 어떤 값이 나오는지 알고 싶은 거임

실제로 선형회귀, 두 데이터 간의 정확한 방정식을 구하는 것 자체로도

큰 연구 결과일 수 있다. 그리고 그라, 그라가 결코 좋은게 아니다.

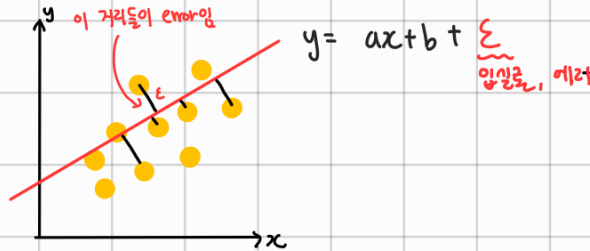
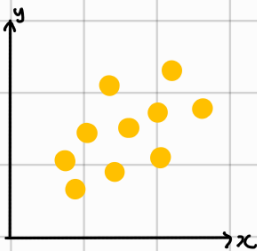
궁극적으로 어떤 표리를 냈을 때, 최대한 단순한게 좋은거임

원래대로 돌아가다, **방정식** 돌아가다 라는 의미

- 회귀: **에러**를 제거한 두 데이터 간의 관계를 도출

• 선형 회귀:  $y = ax + b$  (1차식)

(= 단순회귀)



• 선형 회귀 유도 (기울기)

$$y = \underset{\substack{\text{기울기} \\ \downarrow}}{a}x + \underset{\substack{\text{방정식} \\ \downarrow}}{b} + \varepsilon$$

$$\text{cov}(x, y) = \text{cov}(x, ax + b + \varepsilon)$$

$$= \text{cov}(x, ax) + \cancel{\text{cov}(x, b)} + \text{cov}(x, \varepsilon)$$

b는 상수기 때문에  
어떻게 변하든  
같은 값이 나온다.

⇒ 정지된 상태

⇒ 독립 = 0

변수와 에러는 절대 상관관계가 있어서는 안된다.

있으면 그건 에러가 아님.

어떤 요인이 반드시 숨어있음.

만약 내가 움직일 때마다 에러가 생기면 그걸 야기하는 다른 요인이 있는거임.

그걸 찾아야지 그걸 ε로 리부하면 안됨

$$= a \text{cov}(x, x) = a V(x)$$

자기자신과 cov는 분산과 같은말임

$$\text{cov}(x, y) = a V(x) \text{니까 } a = \frac{\text{cov}(x, y)}{V(x)}$$

## - 선형회귀 유도 (영결편)

· 영의 모형.  $y = ax + b + \varepsilon$

· 기대값 (Expected Value) 의  $E$   $E(y) = E(ax + b + \varepsilon)$

· 영을 계측하여 평균을 내면  
회귀로 따라가겠다. 라는 가정.

$$= E(ax) + E(b) + E(\varepsilon)$$

b가 잔여  
값의  
평균은 b임.

에러들의 평균은 0임.



noise가 0이 안되고 한쪽으로 편향된다면 없애줘야하는 요인임.

$$= aE(x) + b \text{ 이다.}$$

그럼  $b = \underbrace{E(y)}_{\text{영의 평균}} - \underbrace{aE(x)}_{\text{x의 평균} \times a}$

$$b = E(y) - \frac{\text{cov}(x, y)}{V(x)} E(x)$$

## - 예제

$$y = ax + b$$

$$y = \frac{\text{cov}(x, y)}{V(x)} x + E(y) - \frac{\text{cov}(x, y)}{V(x)} E(x)$$

문제)  $\underbrace{[1, 2, 3]}_{=x}$  과  $\underbrace{[1, 3, 5]}_{y \text{ 일때}}$ 의 선형회귀 식은?

$$E(x) = 2$$

$$E(y) = 3$$

$$\text{cov}(x, y) = \frac{4}{3}$$

$$V(x) = \frac{2}{3}$$

$$y = \frac{\frac{4}{3}}{\frac{2}{3}} x + 3 - 2 \cdot 2$$

$$= 2x - 1$$

확인 => 1 넣으면 1, 2 넣으면 3, 3 넣으면 5나옴

- 선형회귀 기울기와 상관계수의 관계

$$y = ax + b$$



선형회귀  
기울기

$$a = \frac{\text{cov}(x, y)}{V(x)} = \frac{\text{cov}(x, y)}{\sigma(x)^2} = \frac{\text{cov}(x, y)}{\sigma(x) \sigma(x)}$$

$\swarrow$   $x, y$ 의 관계  $\searrow$   
 $\swarrow$  분산  $\searrow$   $\sqrt{V} = \sigma$   
 (데이터가 얼마나 퍼져 있는지)  $V = \sigma^2$

상관계수

$$r = \frac{\text{cov}(x, y)}{\sigma(x) \sigma(y)}$$

$\leftarrow \text{cov를}$   
 $\leftarrow \sigma(x) \sigma(y) \text{로 나눈 값.}$

여기까지 비슷하니 때문에, 아래식 유도 가능

$$r \cdot \frac{\sigma(y)}{\sigma(x)} = a$$

$r$ 에  $x$ 의 표준편차를 나누고

$y$ 의 표준편차를 곱하면  $\rightarrow$  선형회귀의 기울기가 나온다.