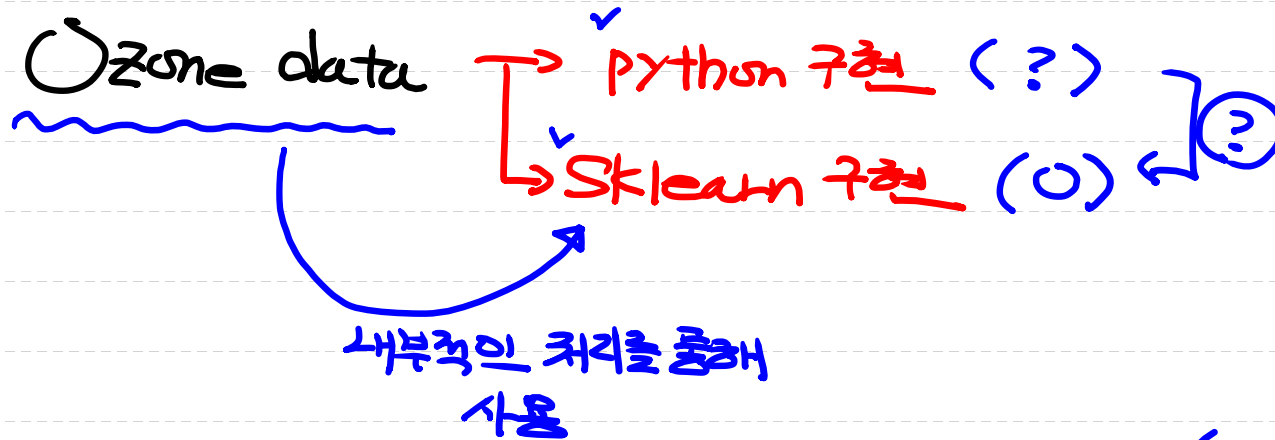


• 03/30



(절차상)
① Missing value 처리

삭제 처리 (전체 데이터가 충분히 많고 Missing value가 5% 이하)

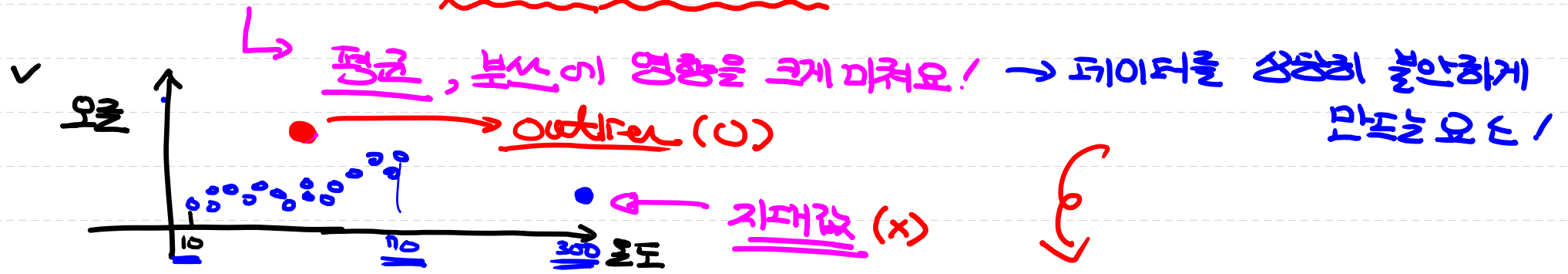
대체 처리

- 평균, 중앙값, 최대, 최대, 최빈
- 머신러닝 기법으로 보충 (더 좋은 방식)

→ Missing value가 조금씩 있을 때

② 이상치 처리

이상치는 값이 일반적인 자료 데이터에 비해 편차가 큰 데이터



이상치 판별
(outlier)

노으로 확인 (그래프로 그려서) →

Boxplot
Scatter

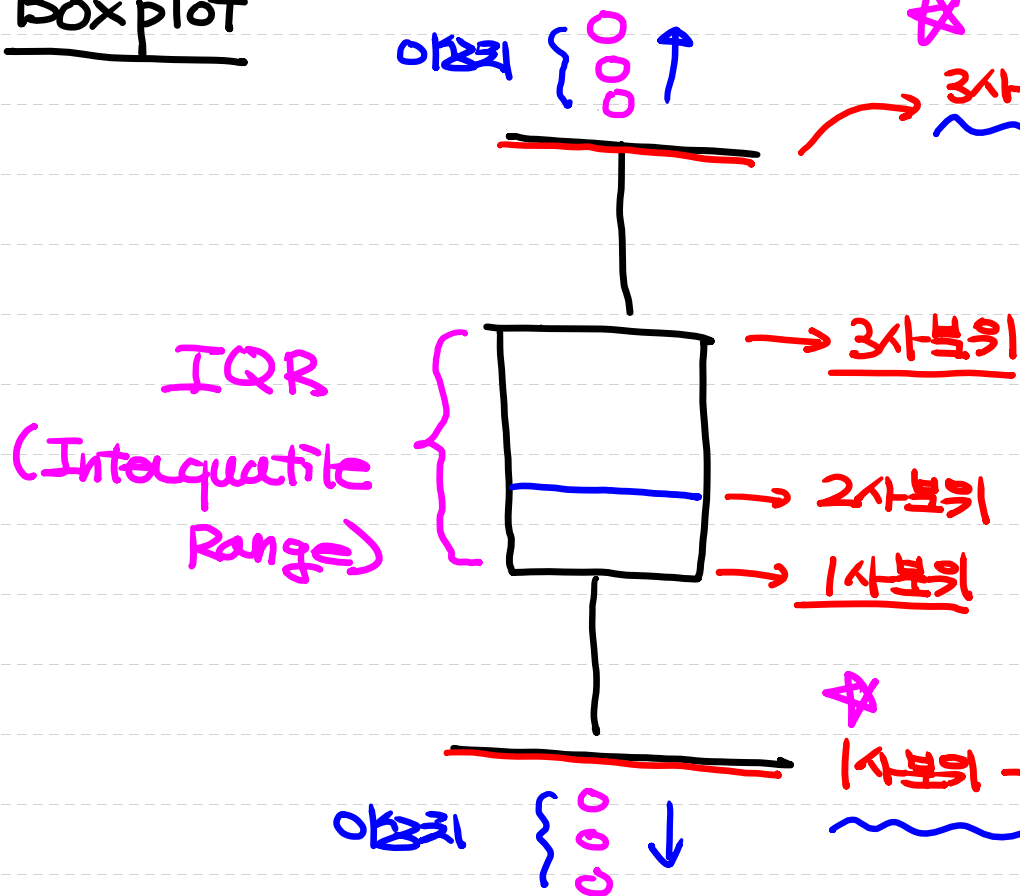
계산을 통해 이상치 파악

많은 기법이 존재하지만

● Turkey Fence ✓

● 표준치 방법 (z-score)

Boxplot



$3\text{사분위} + (\text{IQR value} \times 1.5)$

$\text{IQR value} = 3\text{사분위값} - 1\text{사분위값}$

$\text{IQR value} \times 1.5$



$1\text{사분위} - (\text{IQR value} \times 1.5)$

표준점수

(표준점수, Z-score)

≡ 이분산 이상치 판별 방법

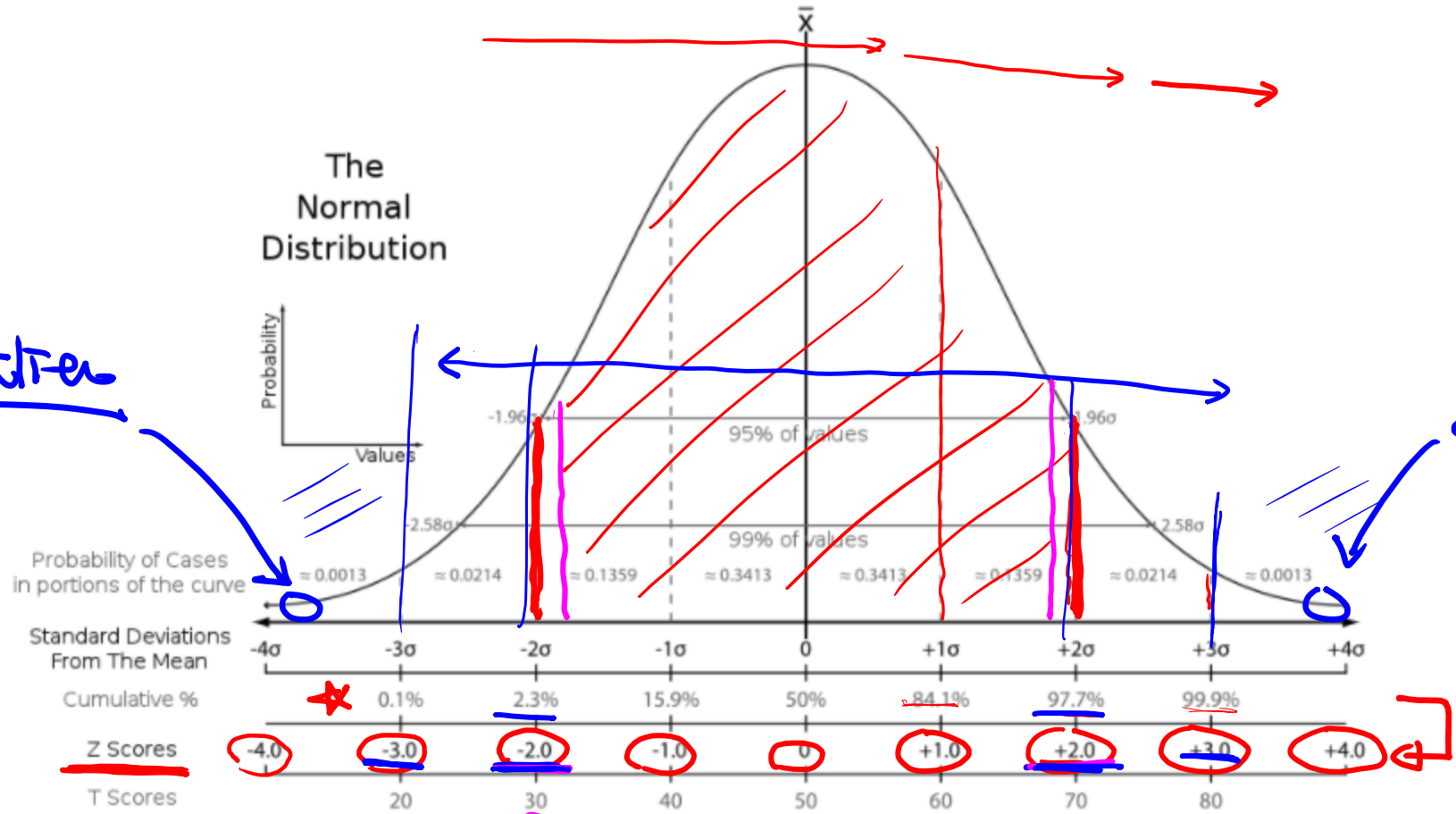


$$\text{Z-score} = \frac{x - \text{평균}}{\text{표준편차}}$$

[1, 2, 3, 4, ..., 22.1]

↓ ↓ ↓ ↓ ↓ ↓ ↓
□ □ □ ○ ○ → Z-score

outlier



이분산의 기준점

● 이상치를 특별히 알고 있어요 !!

→ 하지만 그 값이 실제 "이상치" 인지는 꼭 확인이 필요 !! (?)

그러면 우리의 문제로 돌아와서 !!

✓ Python 구현.
✓ Sklearn 구현.] 구현 결과의 차이
"이상치" 때문인지를 확인해 보세요 !!



확인 결과

→ 아직도 [Sklearn 구현
Python 구현] 차이가 심해요.

↳ 이상치 때문에 문제가
뜬 건 아니었어요 ☹

★ 정규화 (Normalization)

Sklearn은 자체적으로 정규화 진행
(내부적으로)

“데이터가 가진 feature들의 Scale은 맞추려는 작업”

집을 사러가요
 (아파트)

방 개수 2 ~ 10 개

연식(월) 1 ~ 240 개월

(Normalization)

정규화 기법

Min-Max Scaling

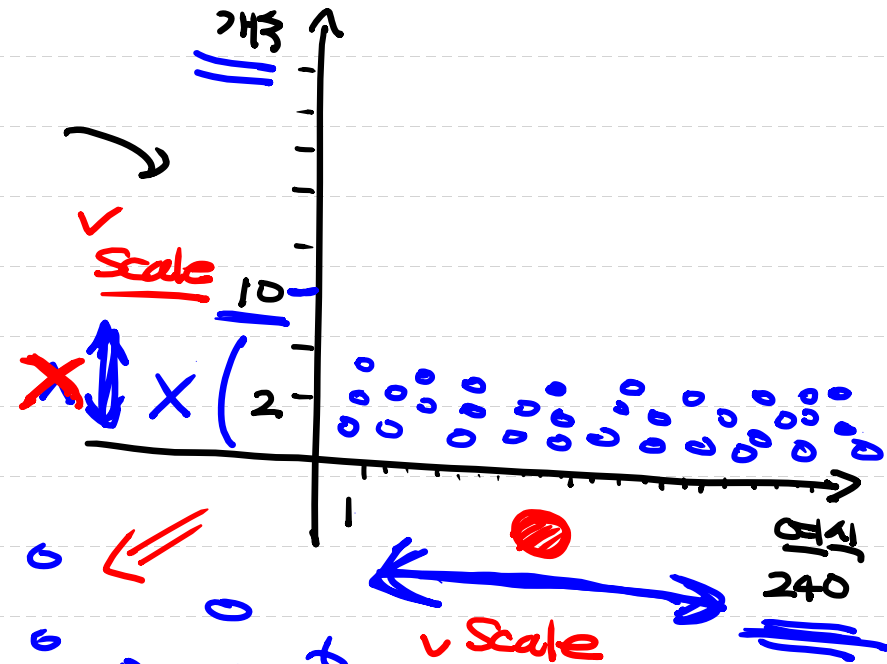
Standardization (표준화)

→ Z-score

$$\frac{x - \mu}{\sigma}$$

최소 → 0
 최대 → 1

최대값
 최소값을 이용하기 때문에 “0/공차에 민감”



“-2 ~ 2” 사이로
 맞추고
 양수, 음수
 있고
 이걸로
 영향을
 잘
 받아요

- ① Simple Linear Regression (단순 선형 회귀) → 독립변수가 1개
- ② Multiple Linear Regression (다중 선형 회귀) → 독립변수가 여러개 (Multi-variable)

우리가 찾고 있는 Linear Regression → Classical Linear Regression

독립변수 3개 → $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

$\hat{y} = \beta_0 + \sum_{i=1}^p \beta_i x_i$

$\hat{y} = wx + b$

$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$

$\hat{y} = X \cdot W + b$

(2차원) (2차원) 1개

$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

$\hat{y} = \beta_0 + \sum_{i=1}^p \beta_i x_i$

$\hat{y} = wx + b$

$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$

$\hat{y} = X \cdot W + b$

(2차원) (2차원) 1개

$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

$\hat{y} = \beta_0 + \sum_{i=1}^p \beta_i x_i$

$\hat{y} = wx + b$

$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$

$\hat{y} = X \cdot W + b$

(2차원) (2차원) 1개

$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$

$\hat{y} = \beta_0 + \sum_{i=1}^p \beta_i x_i$

$\hat{y} = wx + b$

$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$

$\hat{y} = X \cdot W + b$

(2차원) (2차원) 1개