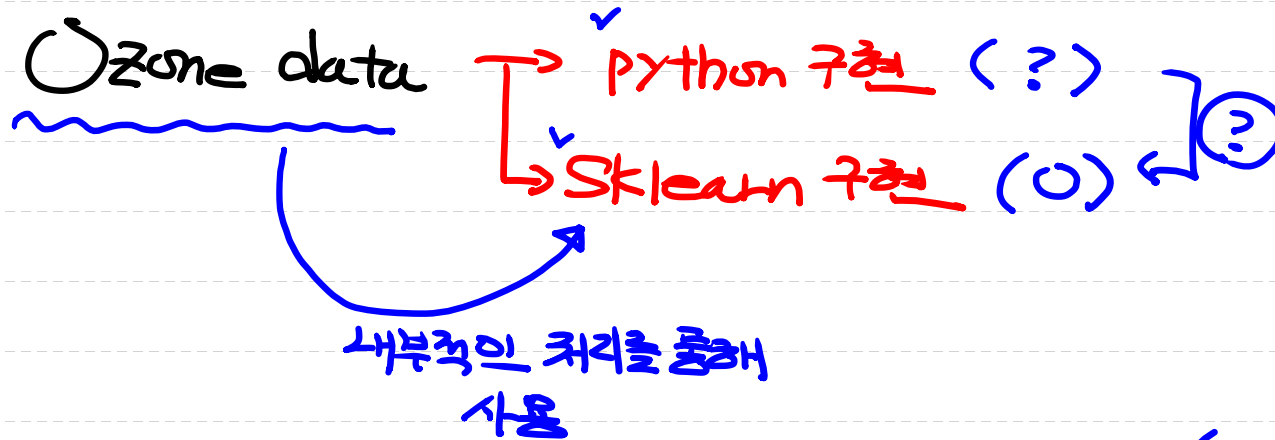


• 03/30



(절대값)  
① Missing value 처리

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

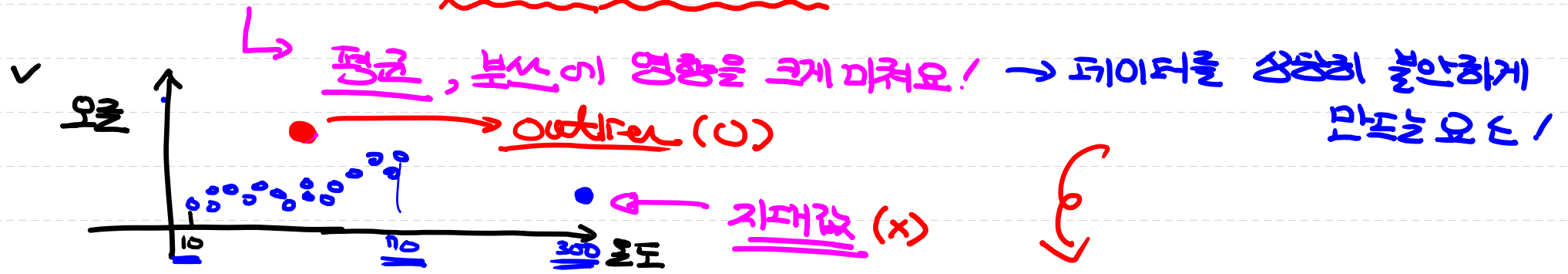
대체 처리

- 평균, 중앙값, 최대, 최대, 최빈
- 머신러닝 기법으로 학습 (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

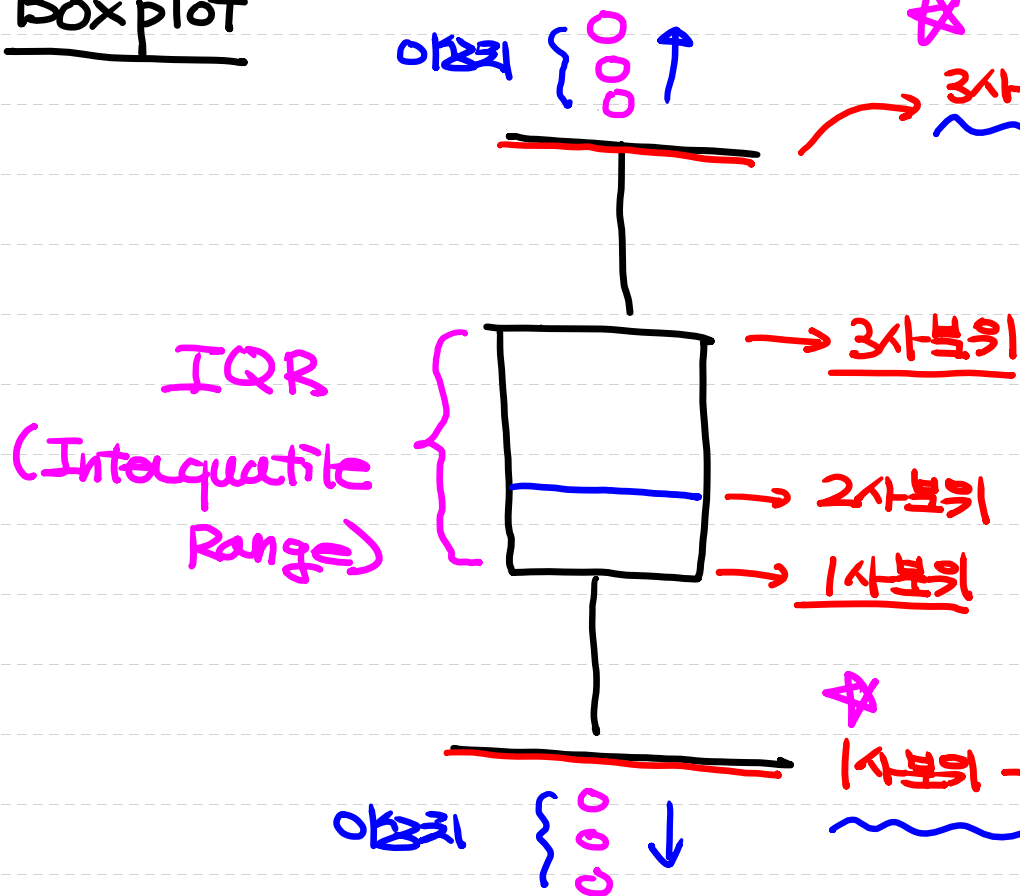
계산을 통해 이상치 파악

많은 기법이 존재하지만

● Turkey Fence ✓

● 표준치 방법 (z-score)

Boxplot



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

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

$\text{IQR value} \times 1.5$



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

표준점

≡ 이분산 이상치 판별 방법

(표준점, Z-score)

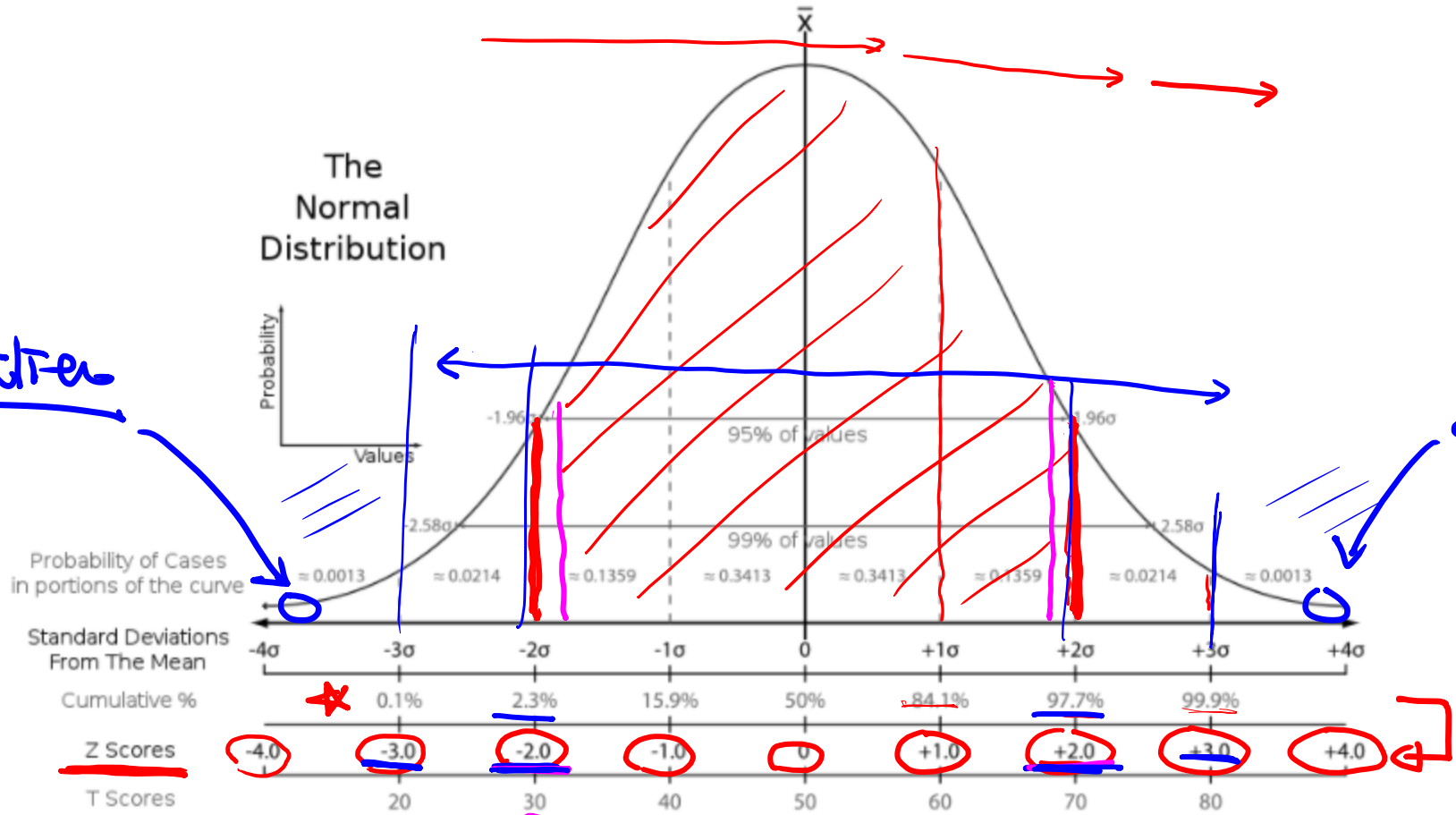


$$\text{Z-score} = \frac{x - \text{평균}}{\text{표준편차}}$$

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

↓ ↓ ↓ ↓ ↓ ↓ ↓  
□ □ □ ○ ○ → Z-score

outlier



outlier

이분산의 기준점

① 이산치를 도출할 수 있습니다!!

→ 하지만 그 깊이 있게 "이해" 인지는 꼭꼭 필요!! (?)

그럼 우리의 문제로 돌아와서 !!

✓ Python 구문.

✓ Sklearn 7월

## 구현결과의 크리치

"아쉬" 부담인지를 확인해 보세요!!



## 학의 결과

→ 아키텍처 [ Sklearn 구현  
Python 구현 ]

차이점 심화됨.

↳ 이상차이로 인 문제  
풀진 아니었어요 ㅠㅠ

# ★ 정규화 (Normalization)

Sklearn은 자체적으로 정규화 진행  
(사부적으로)

“데이터가 가진 feature들의 Scale은 맞추려는 작업”

집을사려가요  
(아파트)

방의개수

2 ~ 10개

연식(월)

1 ~ 240개월

(Normalization)

정규화 기법

Min-Max Scaling

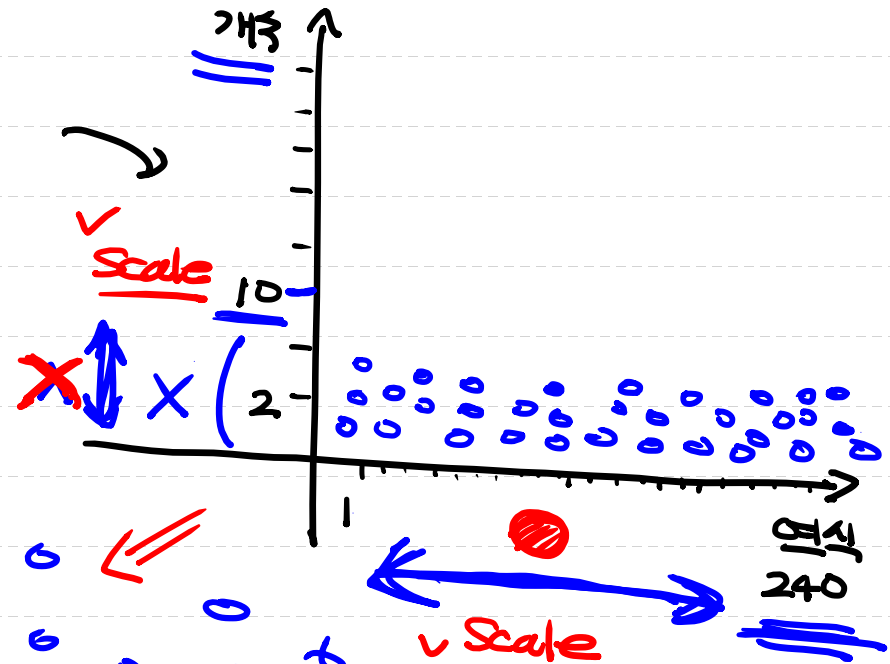
Standardization (표준화)

→ Z-score

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

최소 → 0  
최대 → 1

최대값  
최소값을 이용하기 때문에 “0/1사이”에 민감



“-2 ~ 2” 사이 민감  
양수, 음수  
있고 이상치에  
영향을 잘 받아요