

이상치 (Outlier)

데이터 분석 전처리 (판다스)

이상치가 발생하는 이유?

1. 데이터 입력 오류(Data Entry Errors)
2. 측정오류(Measurement Errors)
3. 샘플링 오류(Sampling Errors)
4. 자연 발생 이상치(Natural Outliers)
5. 데이터 처리 오류(Data Processing Errors)

What is an Outlier?

An outlier is a value that is much larger or smaller in a set of data.

For example, here is a set of data:

3, 35, 37, 38, 40, 42, 44, 76

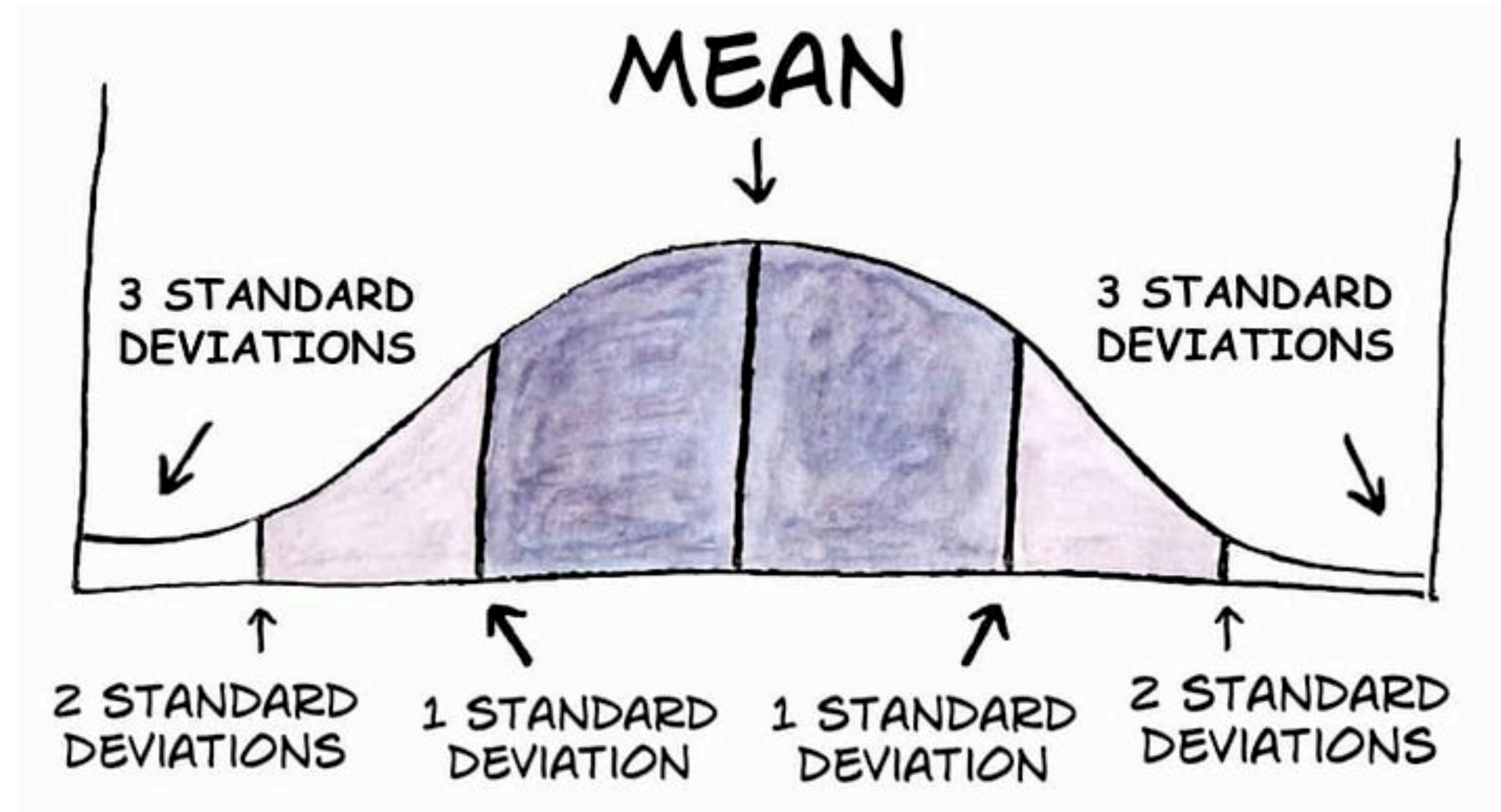
The outliers are 3 and 76 in this set of data as most of the data is in the 30s and 40s.

Outliers can alter data sets and change the mean, median and mode.

It's worth identifying outliers in your data to see whether they help you interpret patterns and information you're collecting.

ink saving Eco

Outlier detection



MAD

$$MAD = \text{median}(|x_i - m|)$$

Median Absolute Deviation
Calculator

$$MAD = \text{median}(|x_i - m|)$$

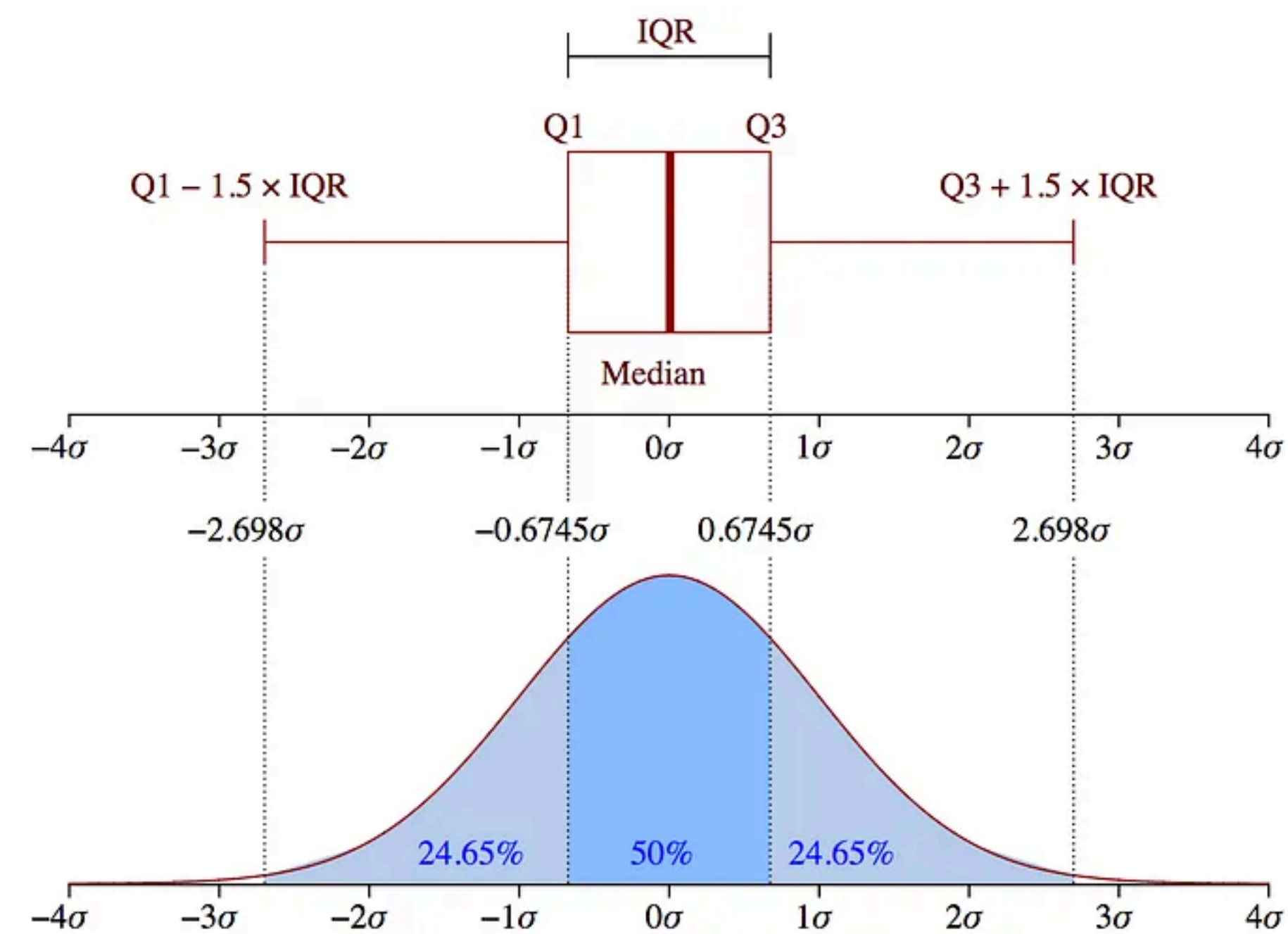
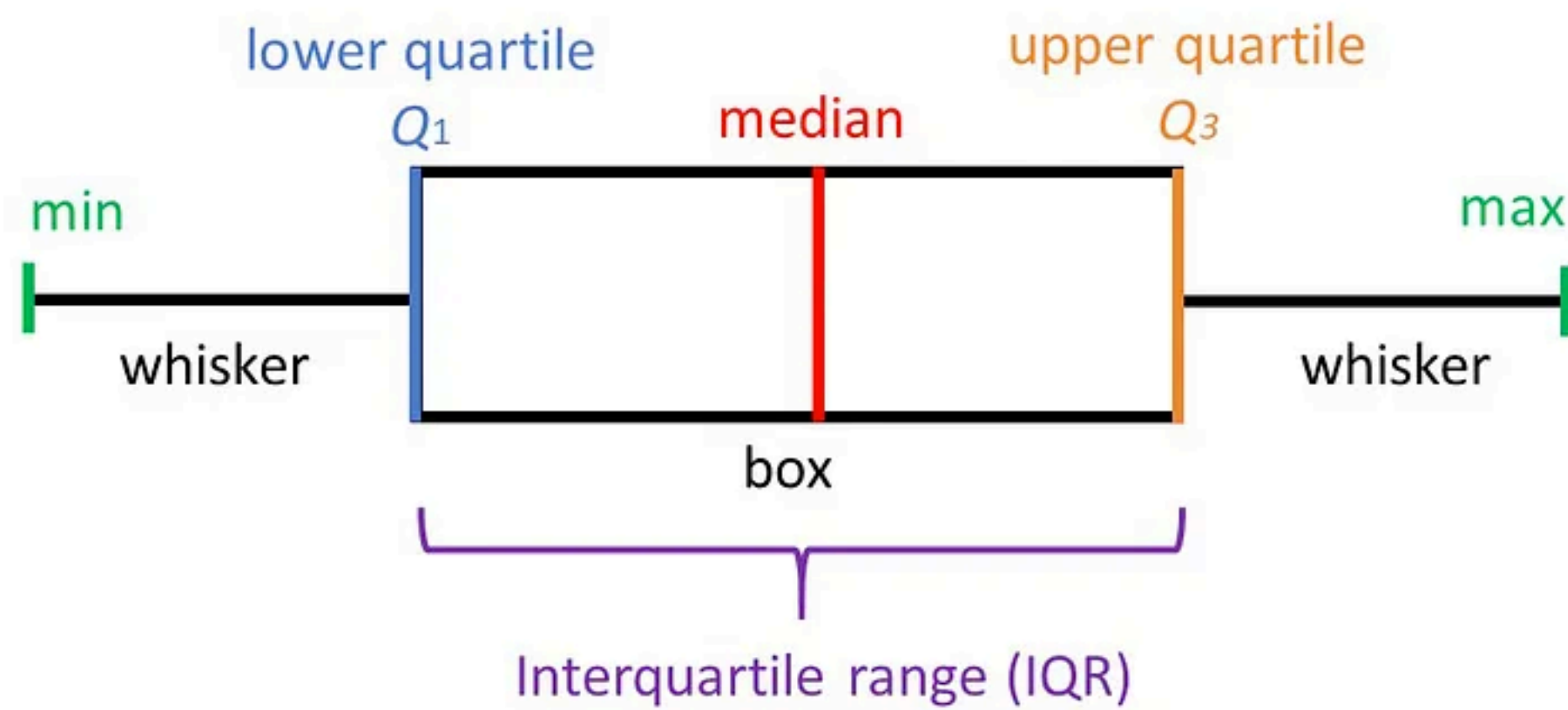
Data (enter up to 50 numbers)

x_1

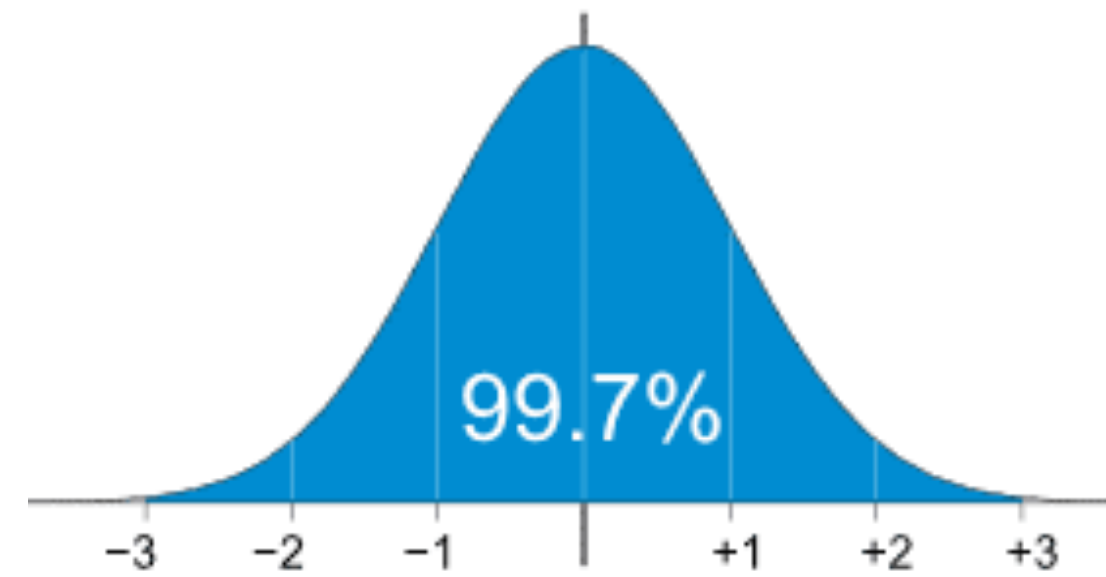
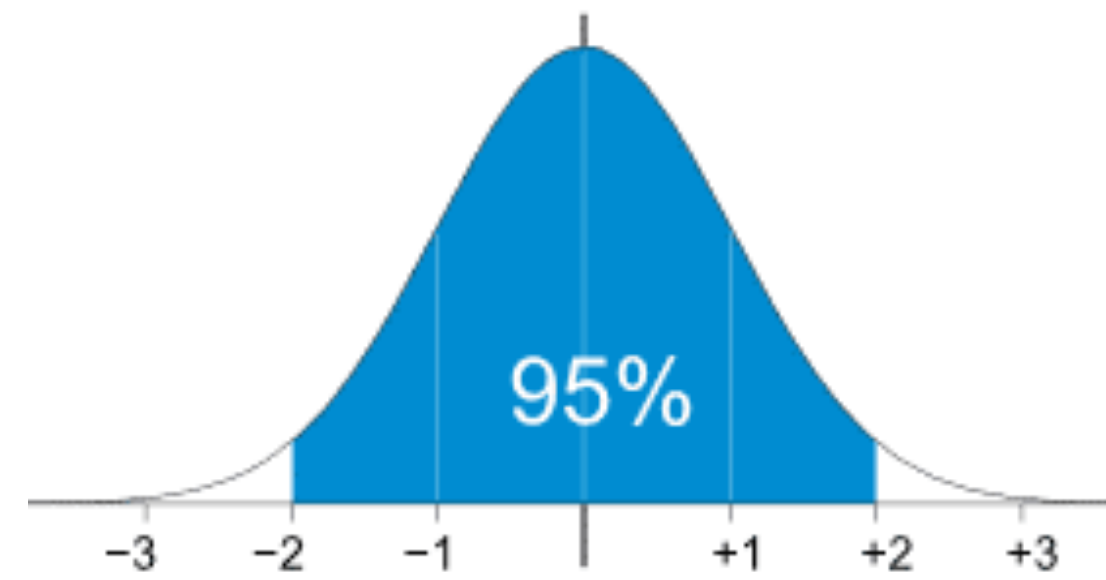
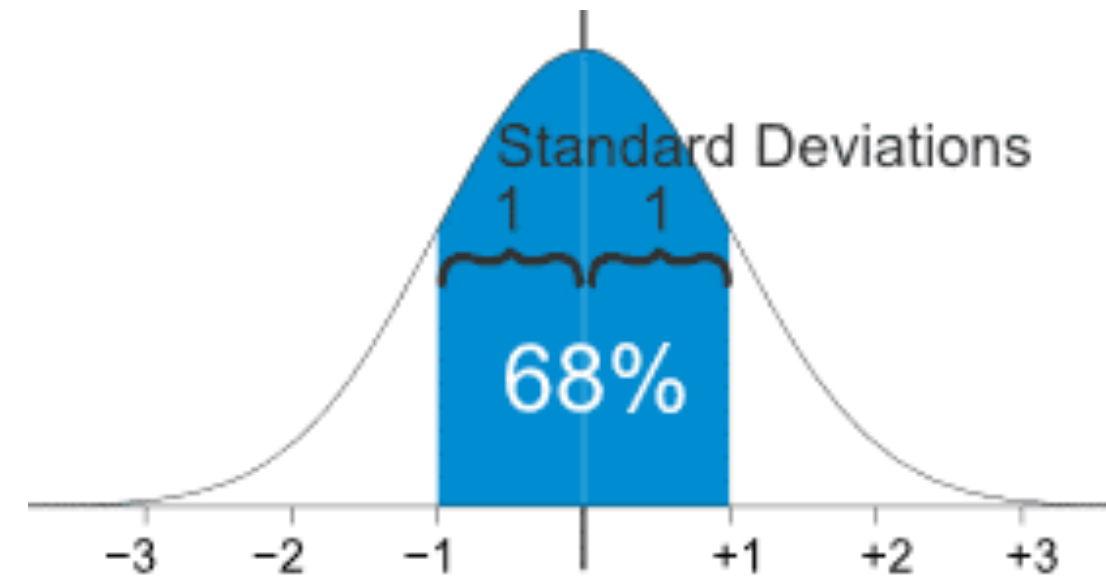
5

—

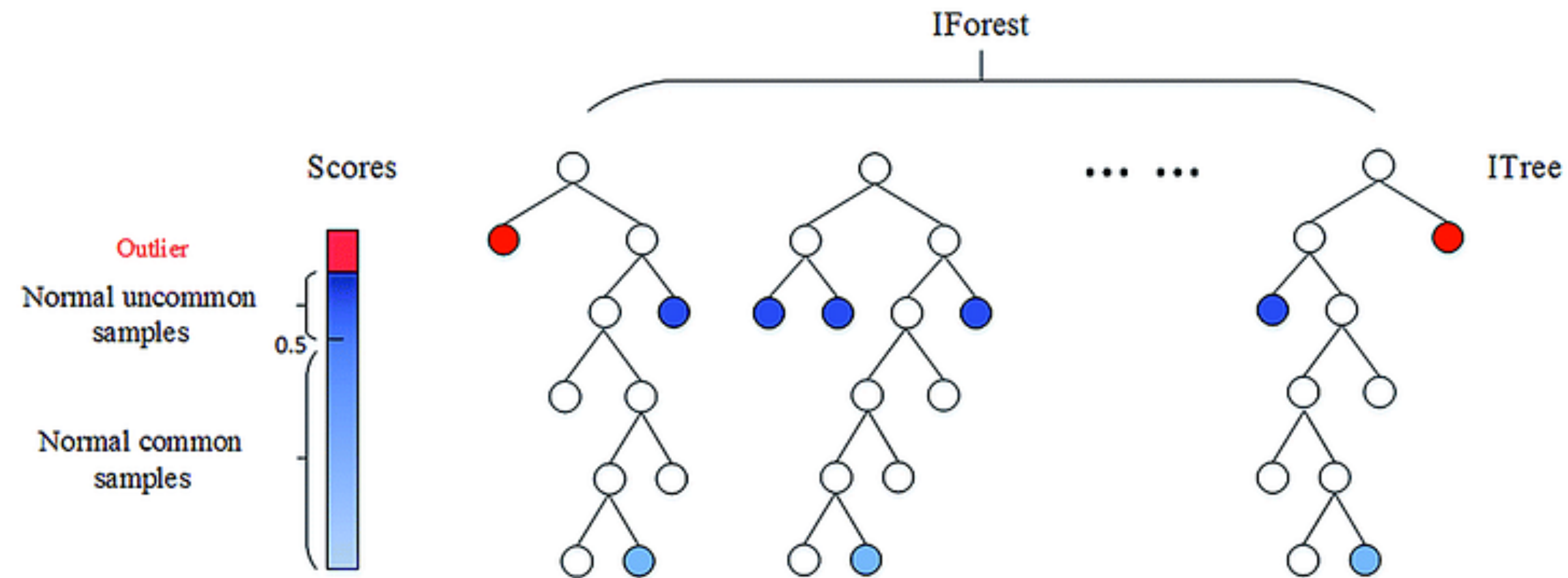
IQR



Standard Deviation

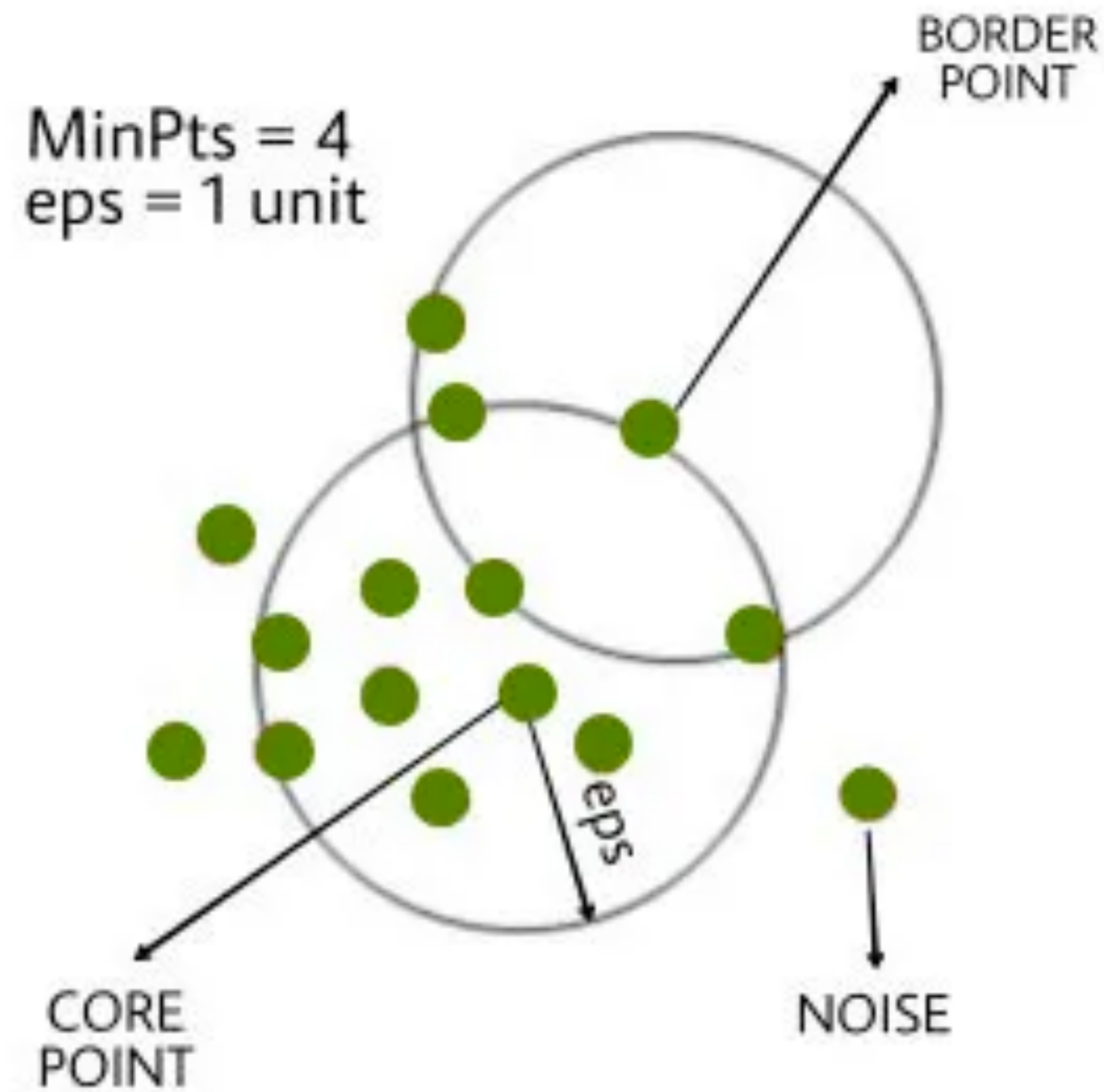


Isolation Forest



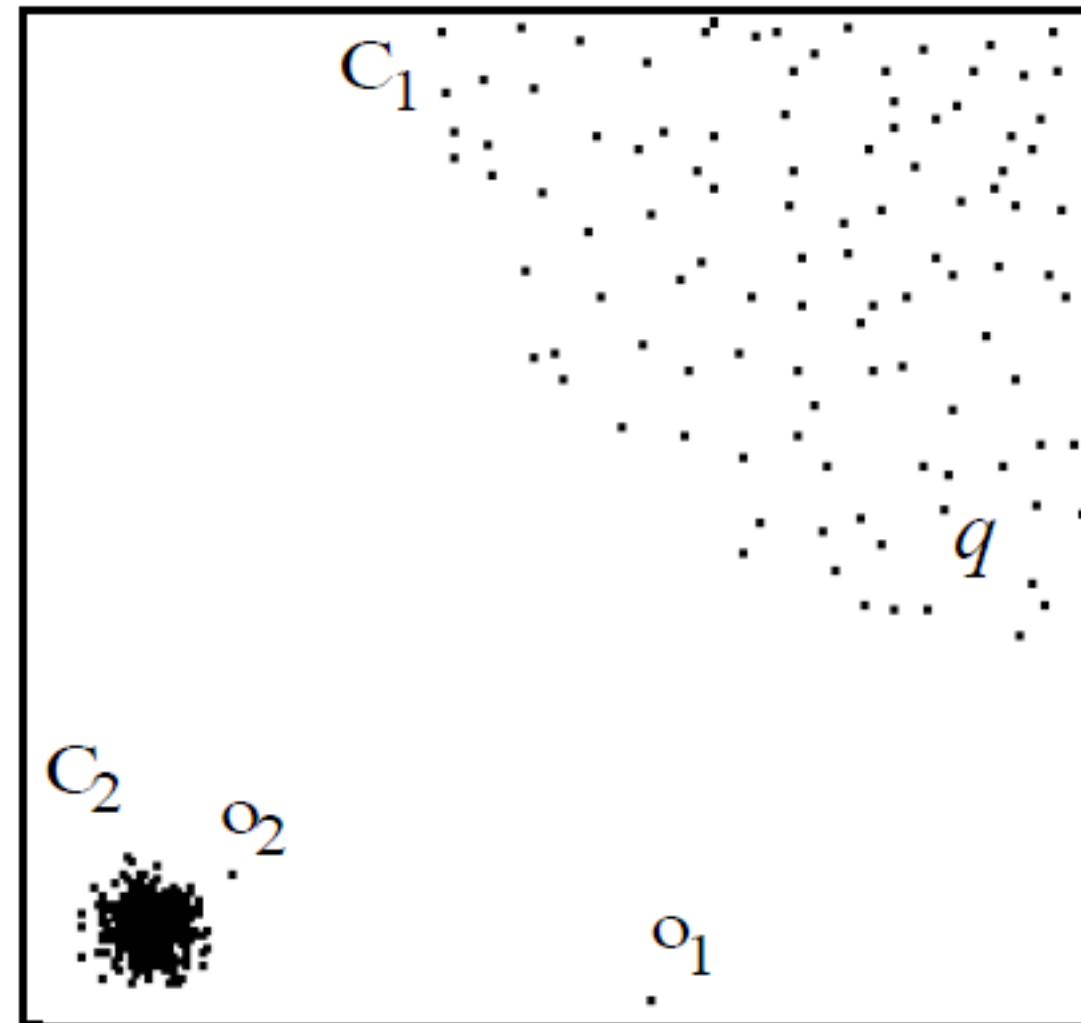
Isolation Forest는 데이터 포인트를 '고립'시키는 데 필요한 분할의 수를 기반으로 이상치를 식별
정상적인 관측치보다 이상치가 고립되는 데 적은 분할이 필요

DBSCAN(Density Based Spatial Clustering of Application with noise)

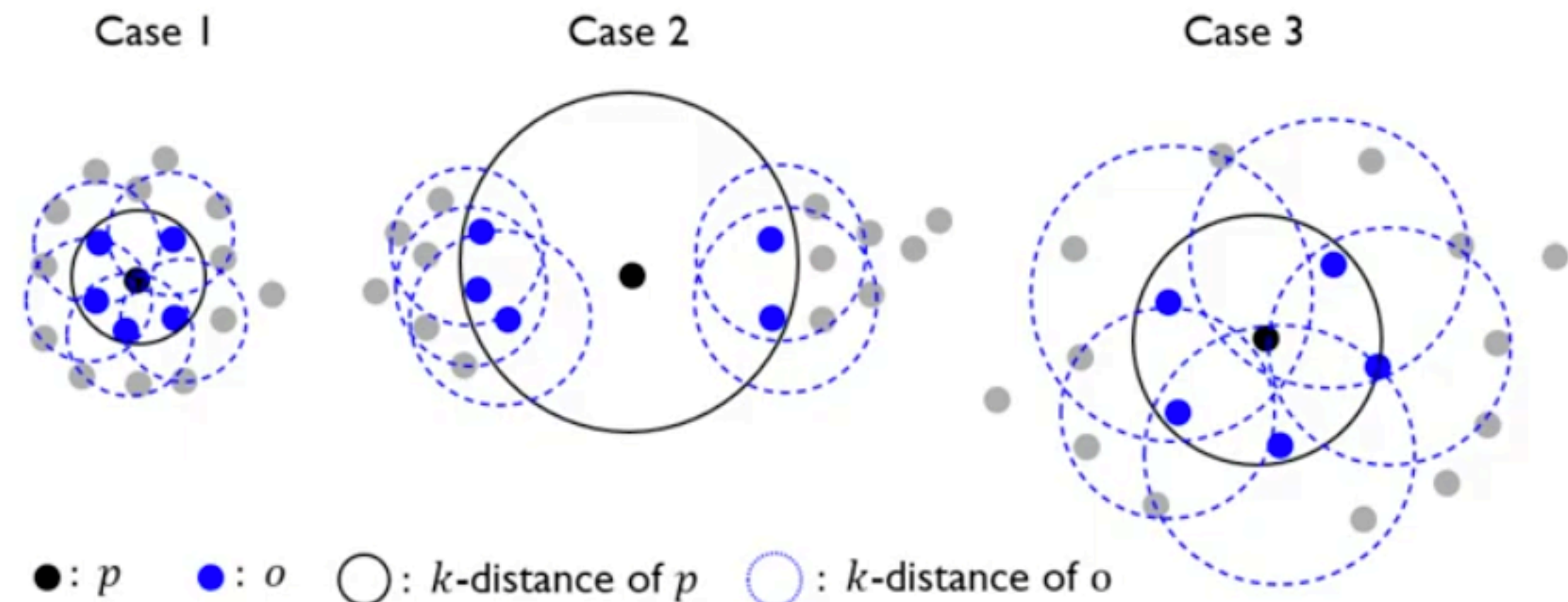


밀도 기반 클러스터링 알고리즘으로, 밀도가 높은 지역의 점들을 클러스터로 그룹화하고, 밀도가 낮은 지역의 점들을 이상치로 간주

LOF(Local Outlier Factor)



$$LOF_k(p) = \frac{\sum_{o \in N_k(p)} \frac{lrd_k(o)}{lrd_k(p)}}{|N_k(p)|} = \frac{1}{lrd_k(p)} \frac{\sum_{o \in N_k(p)} lrd_k(o)}{|N_k(p)|}$$



Case	$lrd_k(p)$	$lrd_k(o)$	$LOF_k(p)$
Case 1	Large	Large	Small
Case 2	Small	Large	Large
Case 3	Small	Small	Small

밀도기반으로 이상값 찾는 알고리즘, 다차원 공간 안에서 이상값들은 주위 밀도에 비해 밀도가 상당히 낮다는 점을 통해 시작

슈하르트 (시계열 데이터 이상치 탐지)

Shewart Control Charts

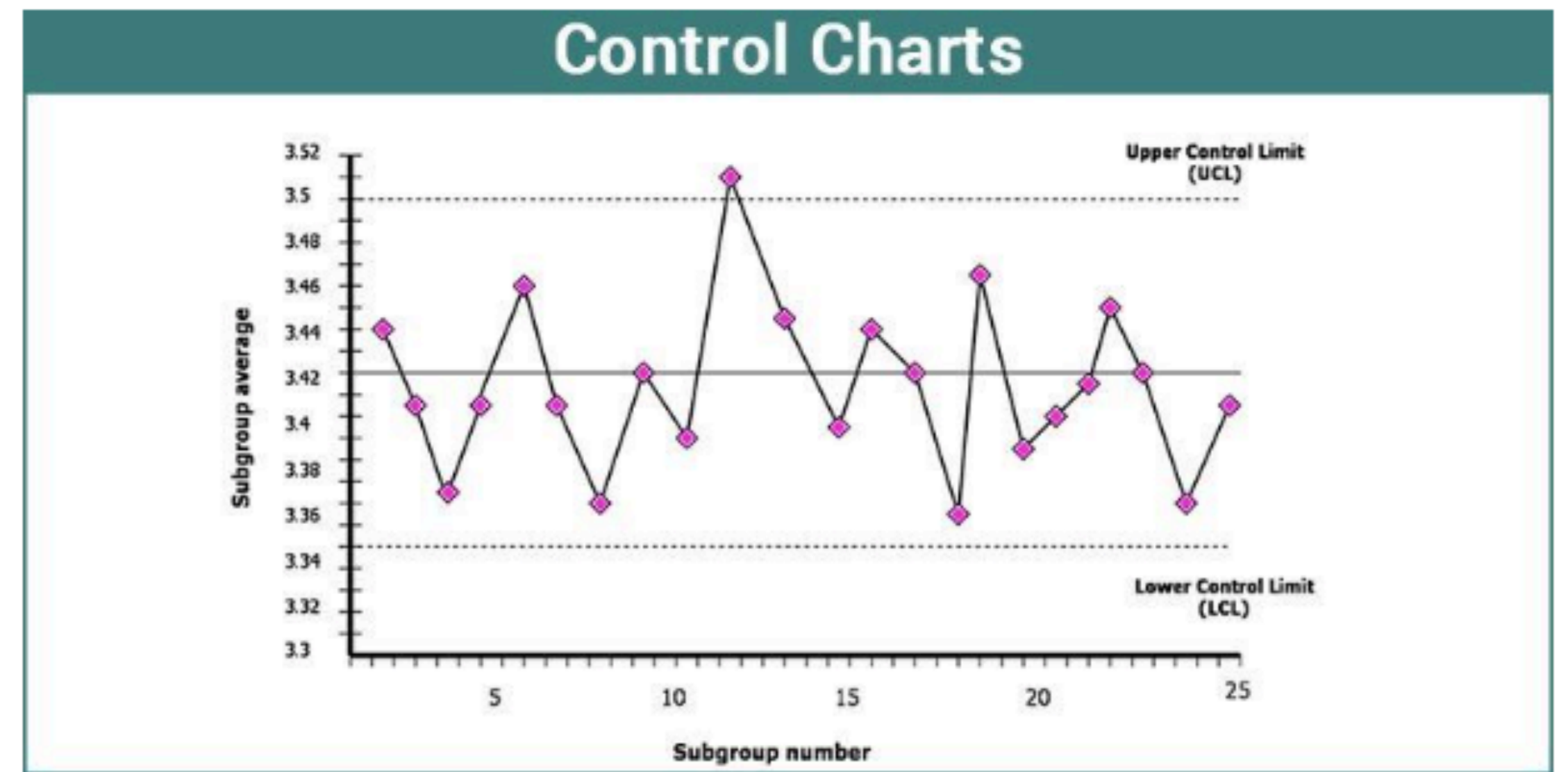
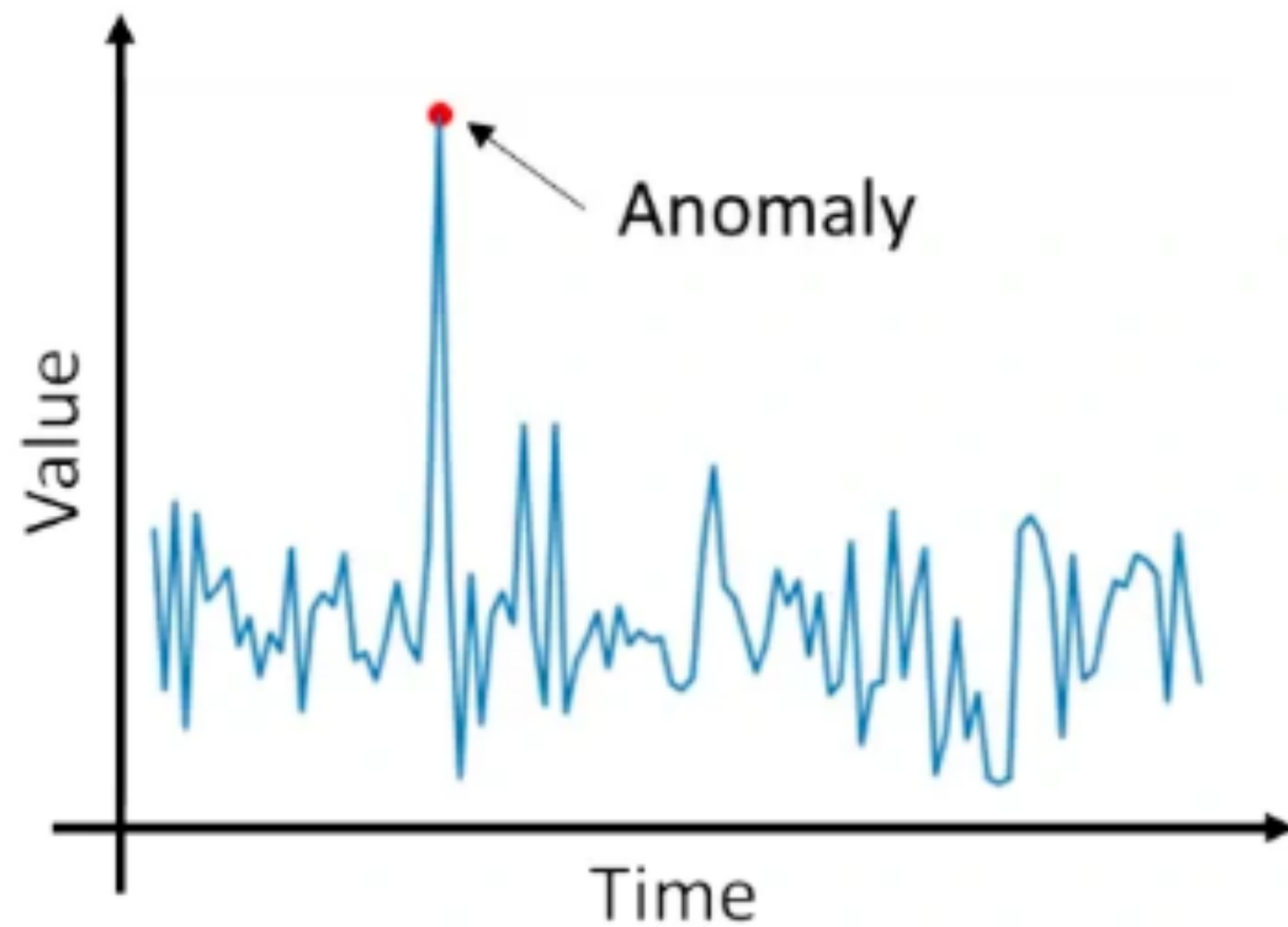


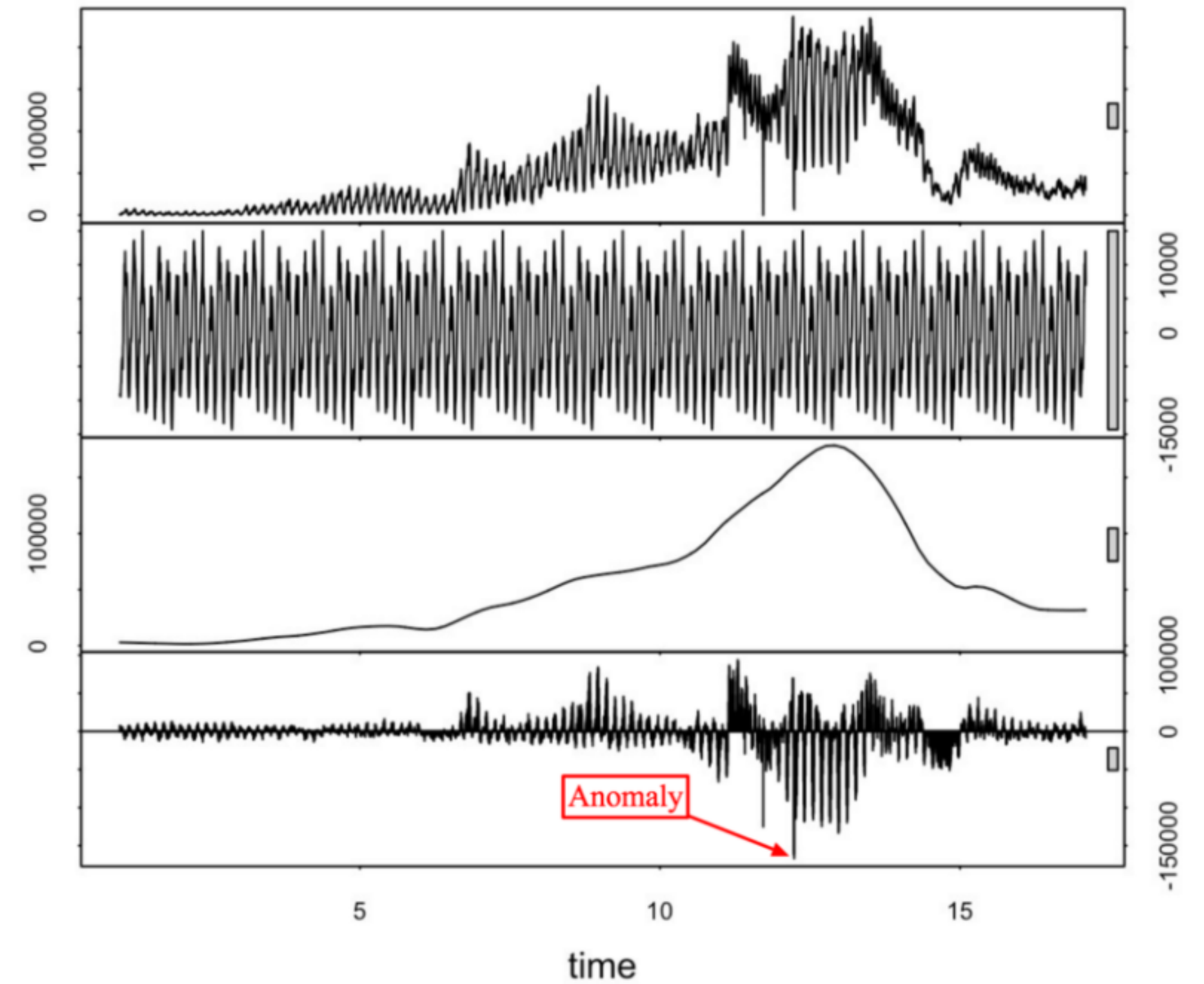
그림1 Control Chart

통계 기반 이상치 탐지 방법론

STL Decomposition

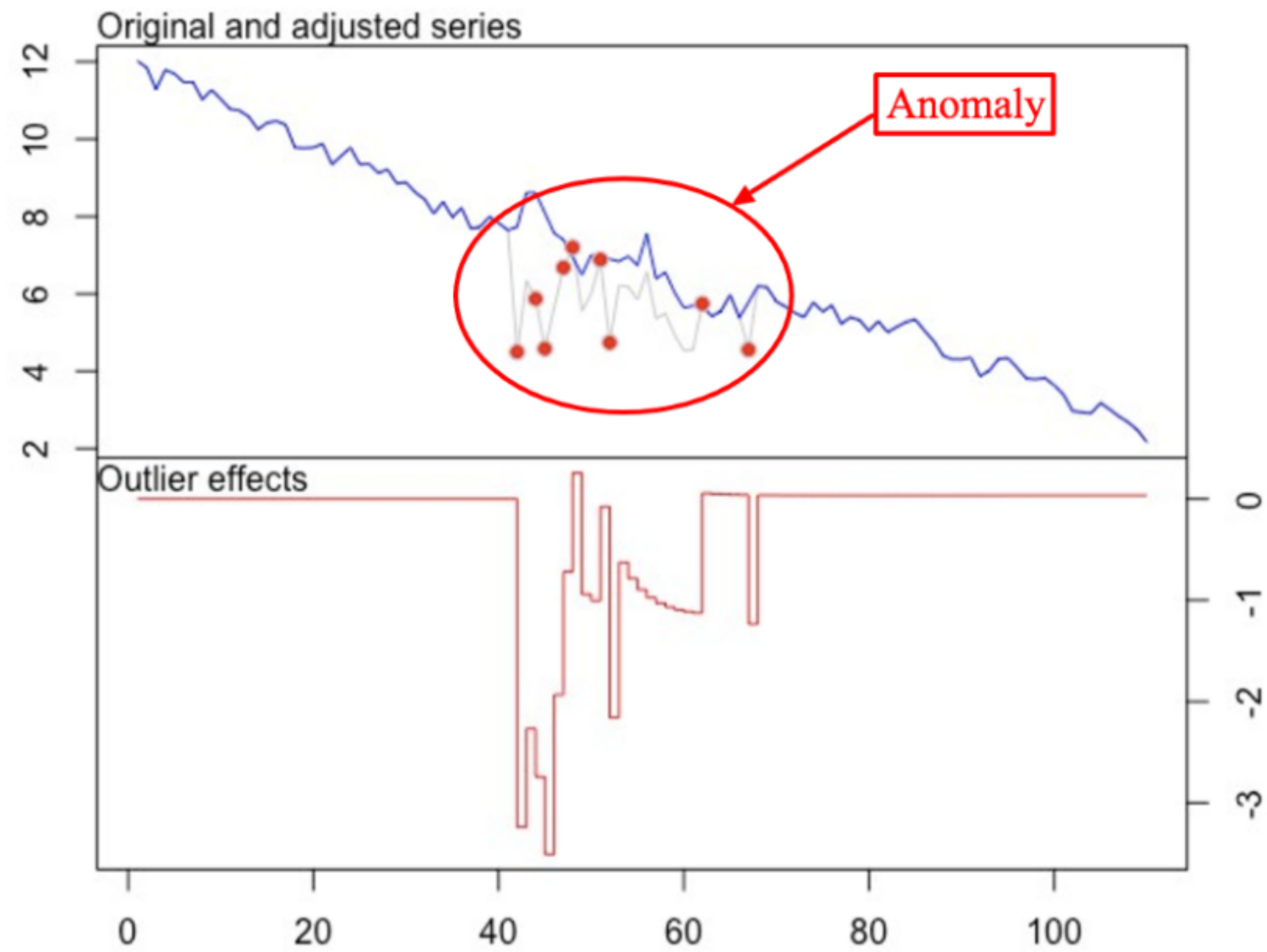
순서대로 original time series, seasonal, trend, residue

맨 위의 original time series를 아래의 3 개의 그래프로 분해



From top to bottom: original time series, seasonal, trend and residue parts retrieved using STL decomposition.

ARIMA



Two time series built using original ARIMA model and adjusted for outliers ARIMA model.

데이터 분석가의 도메인 지식을 통한 Outlier 선정