

EcoAI

랩세미나/논문 리뷰

이상치 정의에 관하여
논문: A survey on outlier explanations

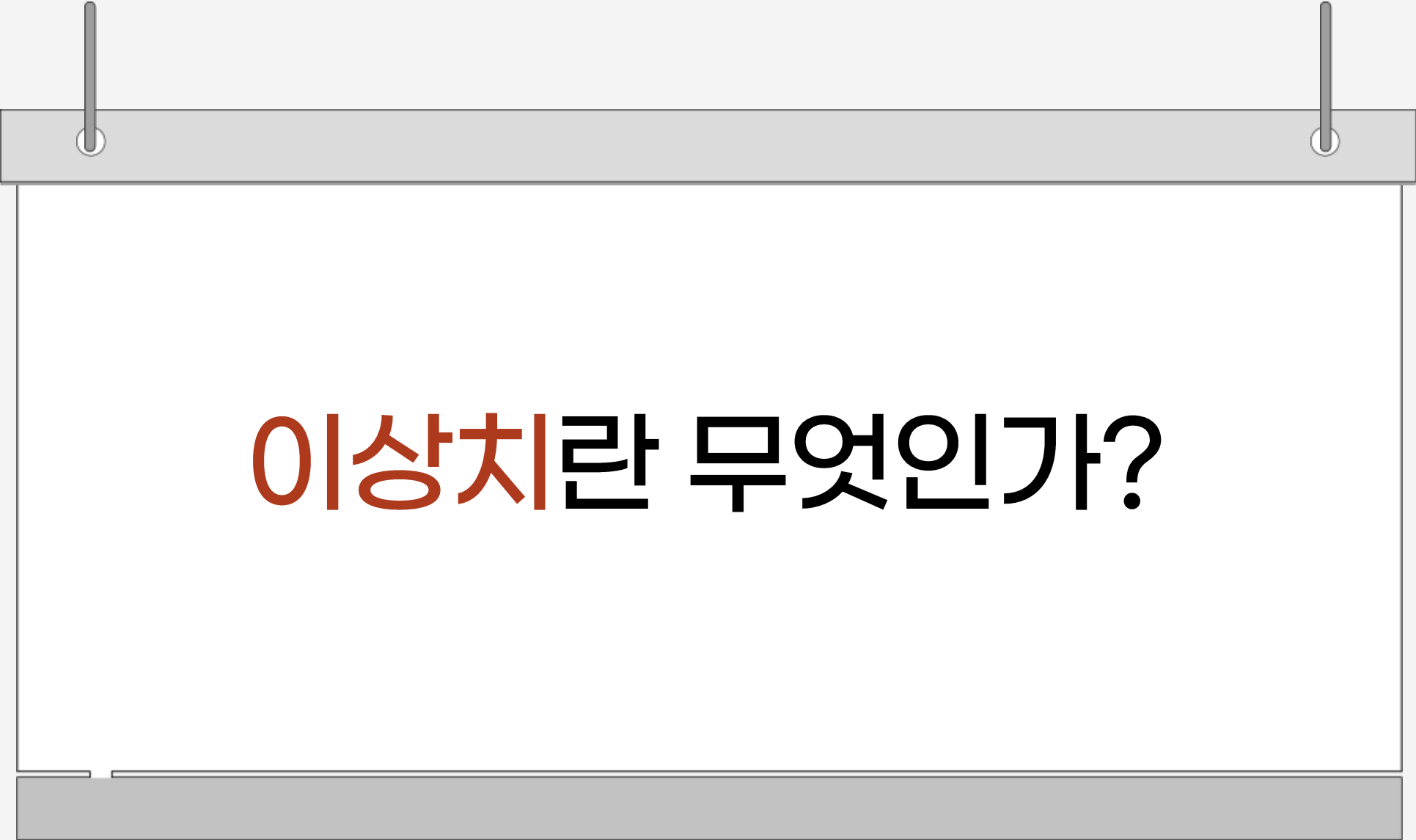
권우현, 정민성

#이상치 #정의

CONTENTS

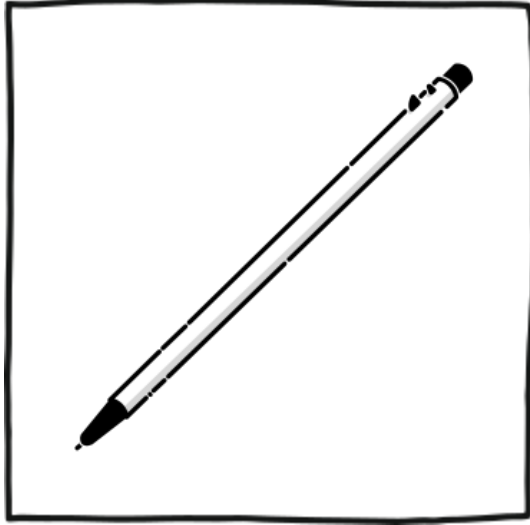


1	이상치란 무엇인가?
2	논문 리뷰
3	GitHub



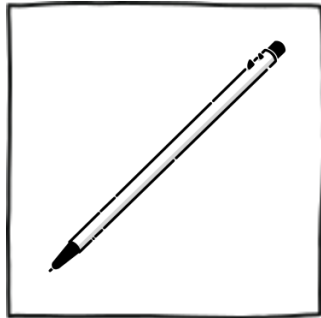
이상치란 무엇인가?

이상치란 무엇인가?



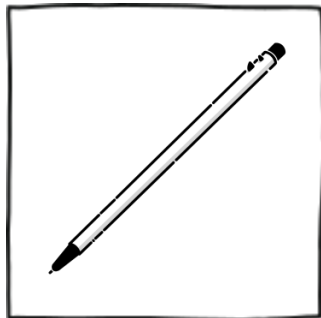
- A normal pen
- 0.5mm
- 모〇미

이상치란 무엇인가?



-A normal pen
-0.5mm
-모〇미

이상치란 무엇인가?

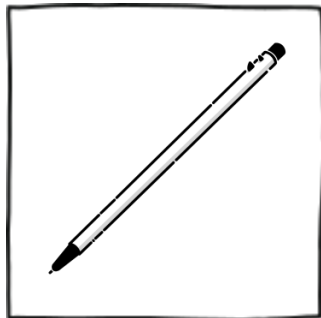


-A normal pen
-0.5mm
-모〇미



>평범한 공장

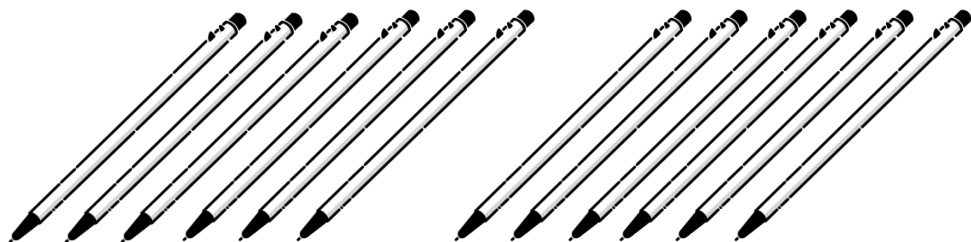
이상치란 무엇인가?



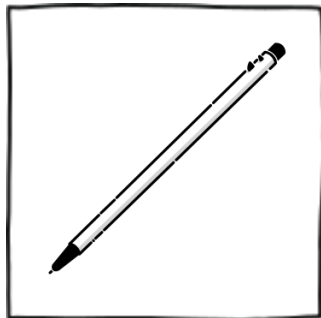
-A normal pen
-0.5mm
-모〇미



>평범한 공장



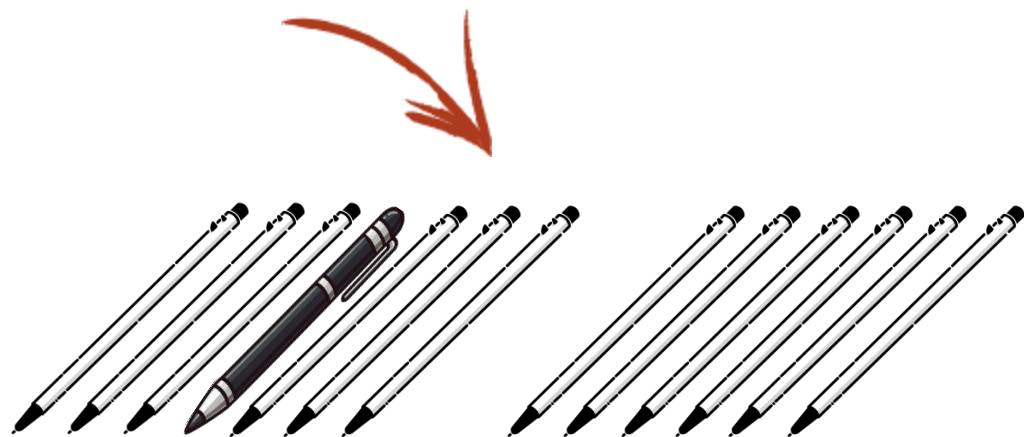
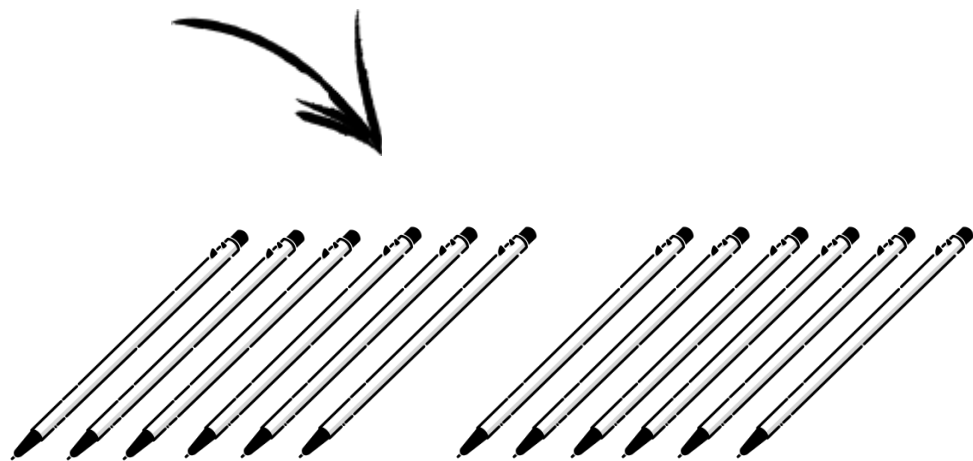
이상치란 무엇인가?



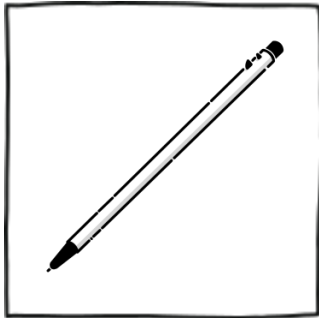
-A normal pen
-0.5mm
-모〇미



>평범한 공장



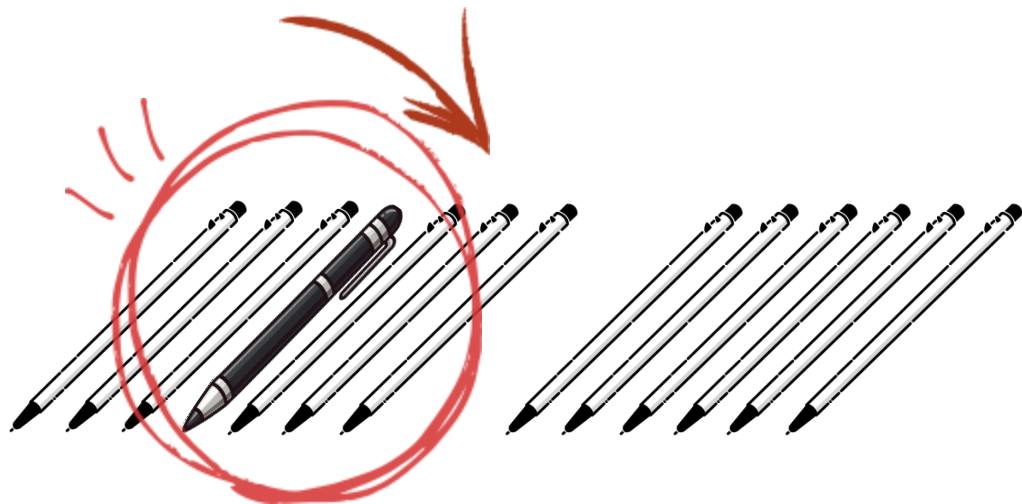
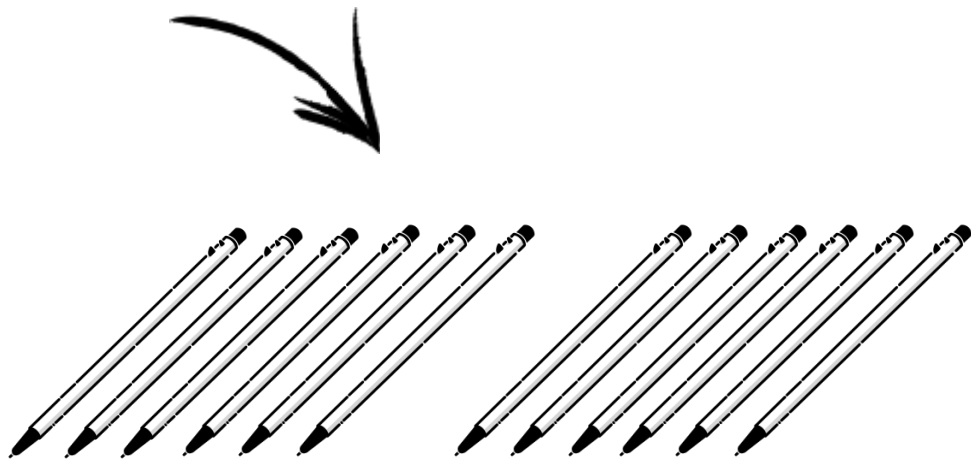
이상치란 무엇인가?



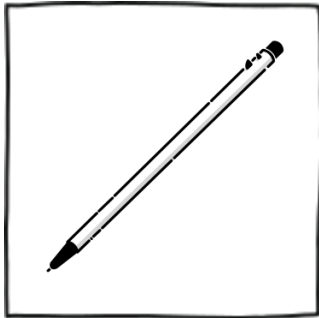
-A normal pen
-0.5mm
-모〇미



>평범한 공장



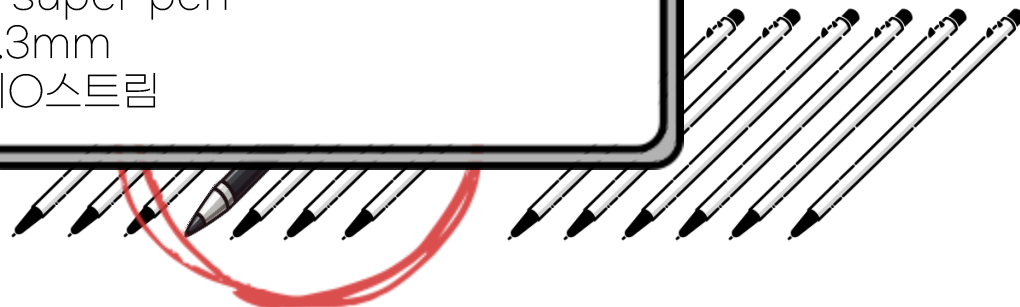
이상치란 무엇인가?



>평범한 공장

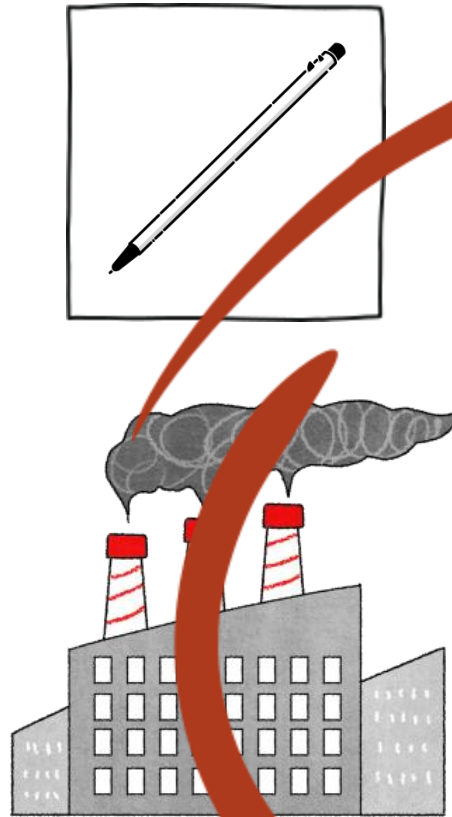


-A super pen
-0.3mm
-제O스트림



이상치란 무엇인가?

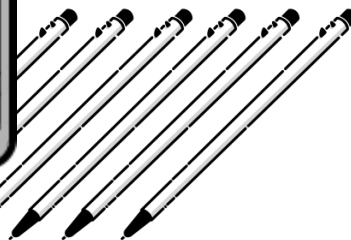
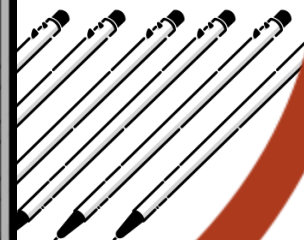
이상 볼펜!



>평범한 공장



-A super pen
-0.3mm
-제0스트림



| EcoAI |

이상치란 무엇인가?

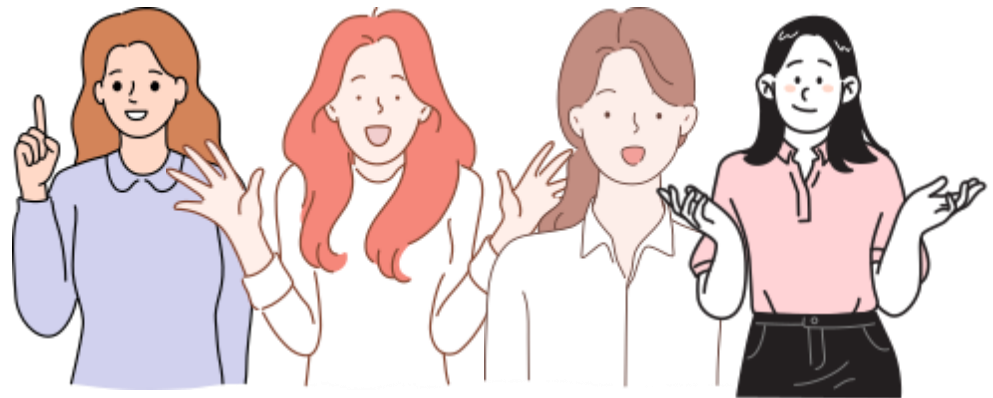


평범한 대학생A

이상치란 무엇인가?

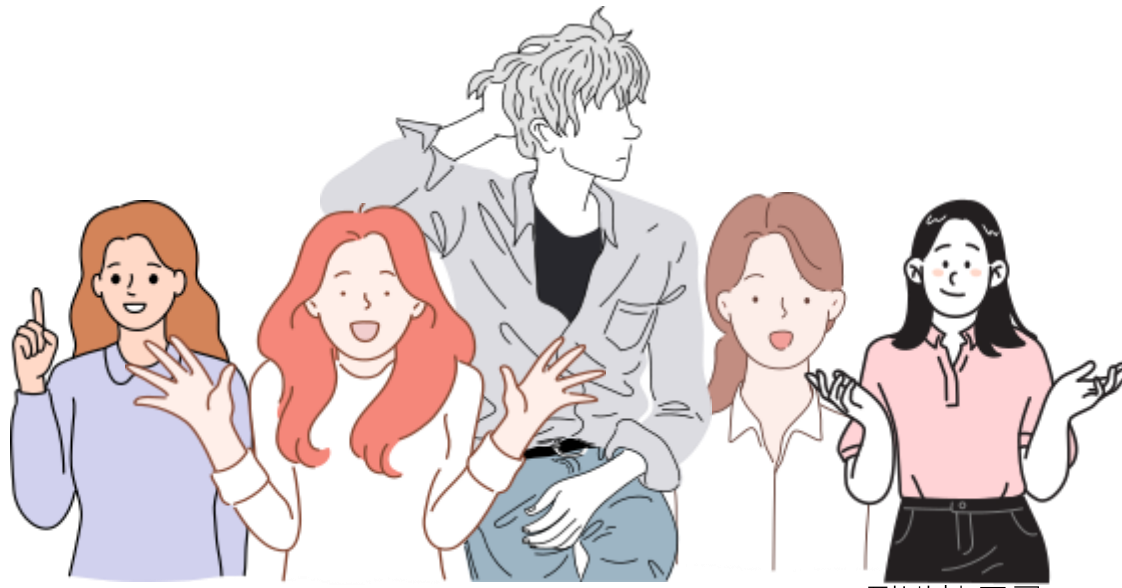


평범한 대학생A



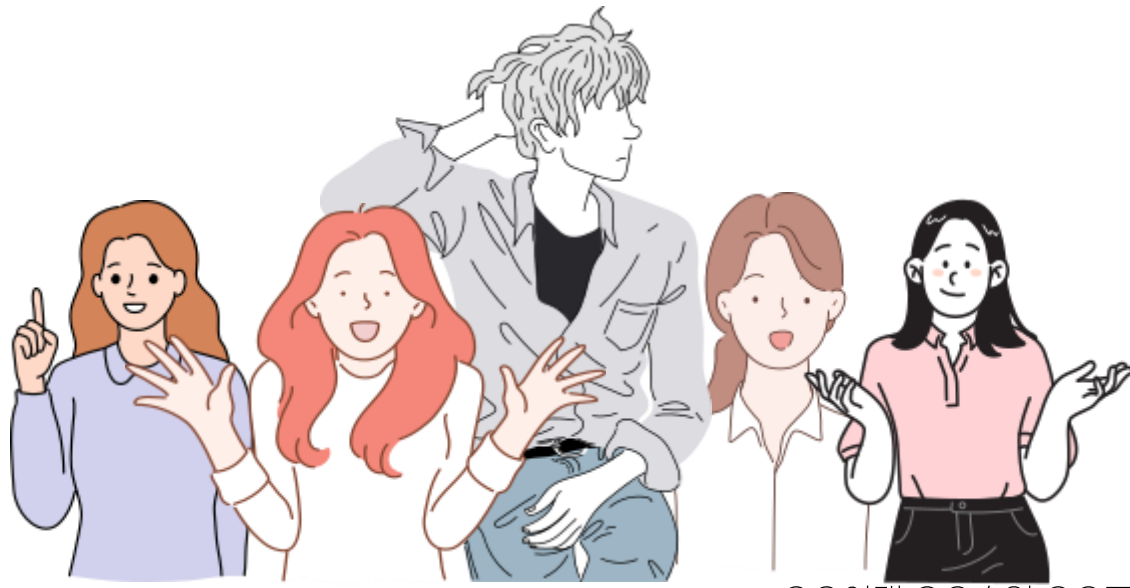
평범한 여학생 그룹

이상치란 무엇인가?



평범한 그룹

이상치란 무엇인가?

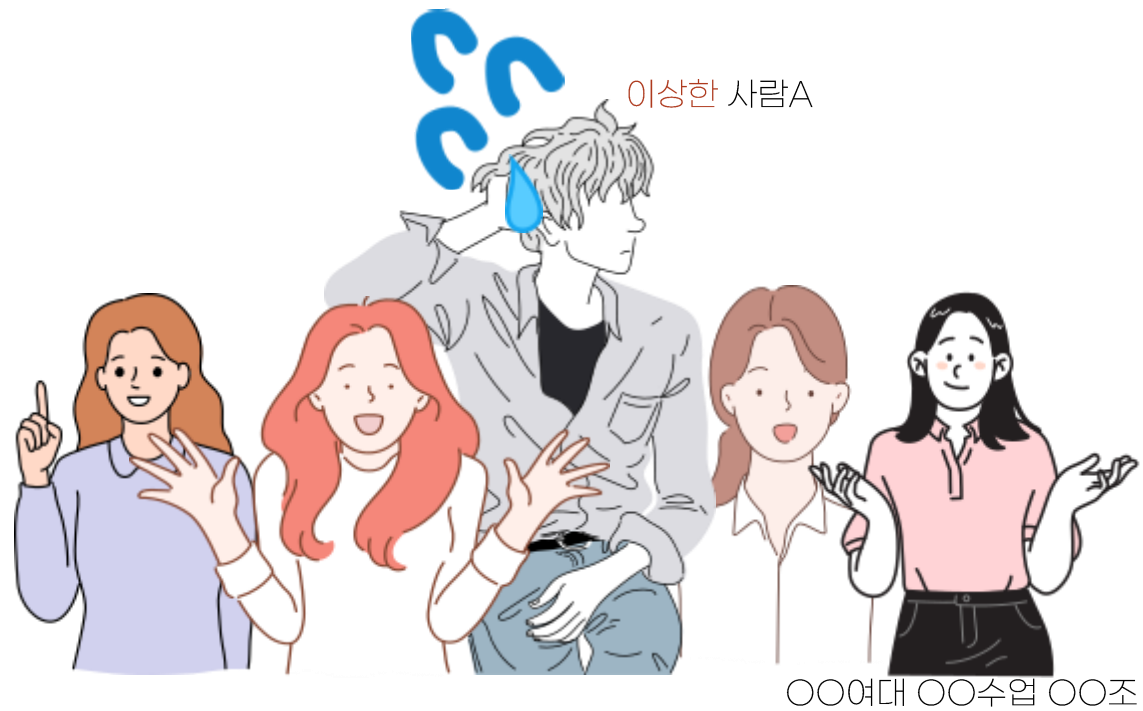


〇〇여대 〇〇수업 〇〇조

이상치란 무엇인가?



이상치란 무엇인가?

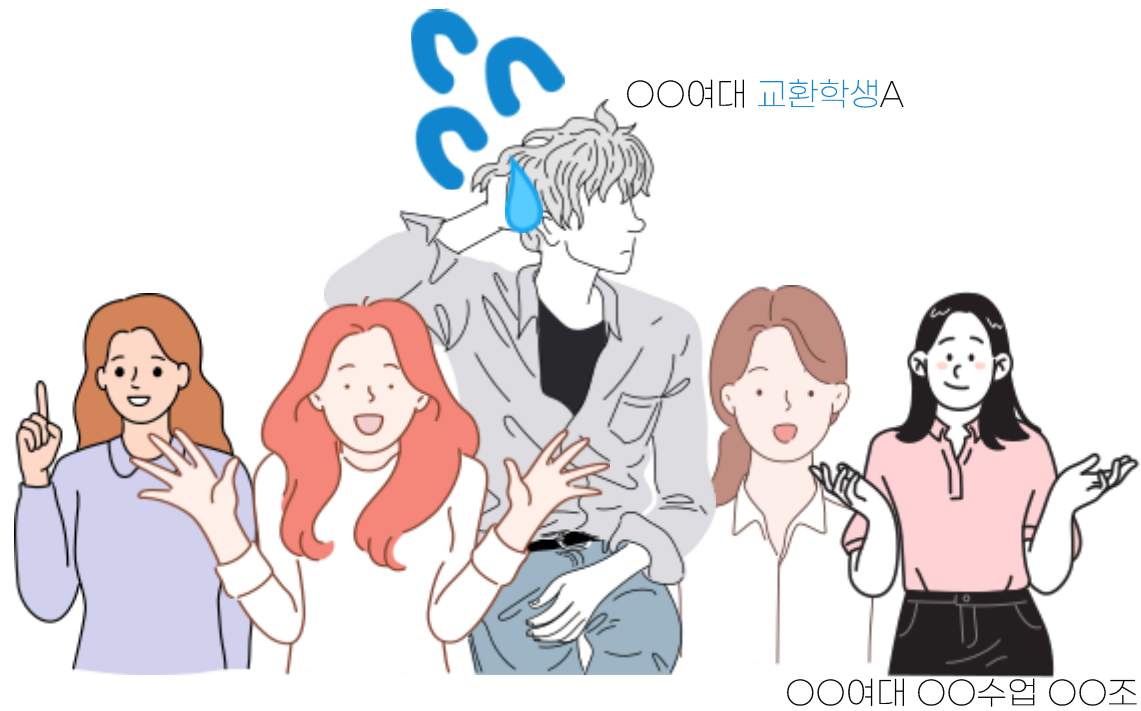


-따라서 상황, 환경에 따라 **이상치** 정의가 달라짐

ex) OO여대 발표 조에 낀 남학생

-하지만 **인지적 관점**에 따라 **이상치**로 보지 않을 수도 있음

이상치란 무엇인가?

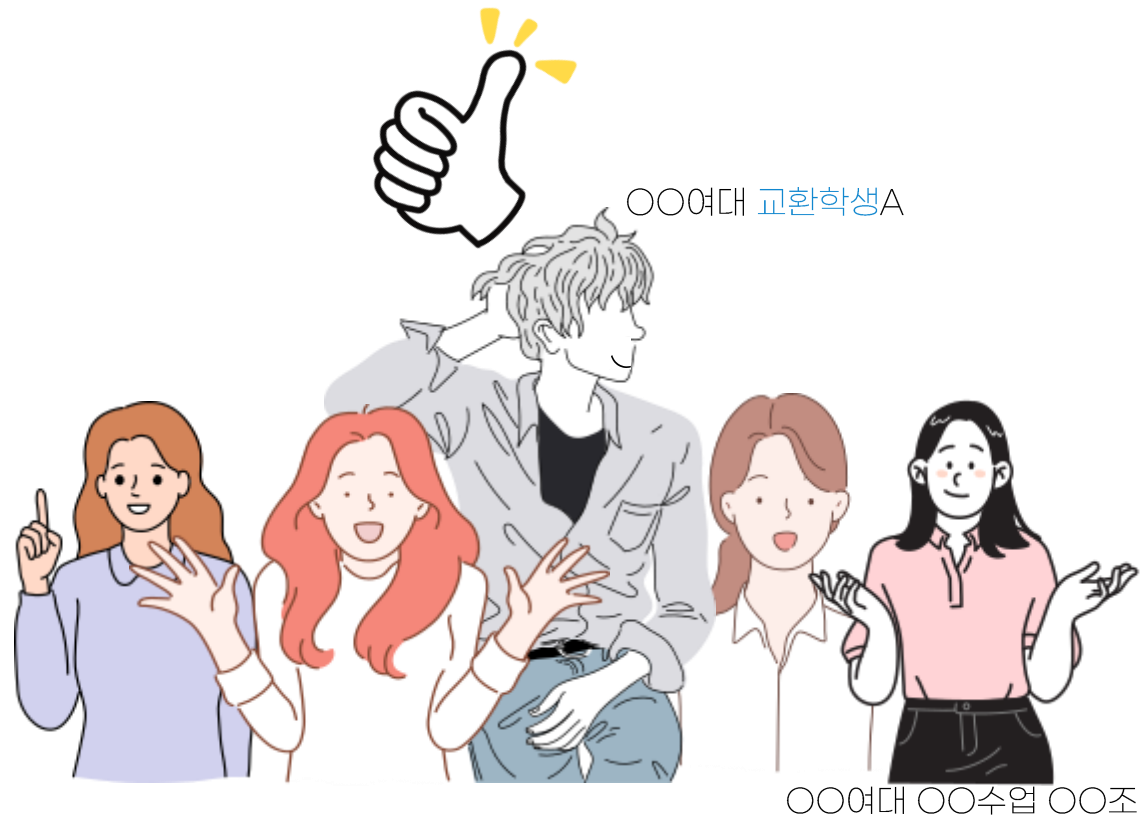


-따라서 상황, 환경에 따라 **이상치** 정의가 달라짐

ex)OO여대 발표 조에 낀 남학생

-하지만 **인지적 관점**에 따라 **이상치**로 보지 않을 수도 있음

이상치란 무엇인가?



-따라서 상황, 환경에 따라 **이상치** 정의가 달라짐

ex)OO여대 발표 조에 낀 남학생

-하지만 **인지적 관점**에 따라 **이상치**로 보지 않을 수도 있음

이상치란 무엇인가?



-따라서 상황, 환경에 따라 **이상치** 정의가 달라짐

ex)OO여대 발표 조에 낀 남학생

-하지만 **인지적 관점**에 따라 **이상치**로 보지 않을 수도 있음

>우리가 가진 **정보**가 부족했기 때문

| EcoAI |

이상치란 무엇인가?

[10, 13, 1000, 50, 15, 20]

| EcoAI |

이상치란 무엇인가?

[10, 13, 1000, 50, 15, 20]

| EcoAI |

이상치란 무엇인가?

[10, 13, 50, 15, 20]

이상치란 무엇인가?

[10, 13, 50, 15, 20]

이상치란 무엇인가?

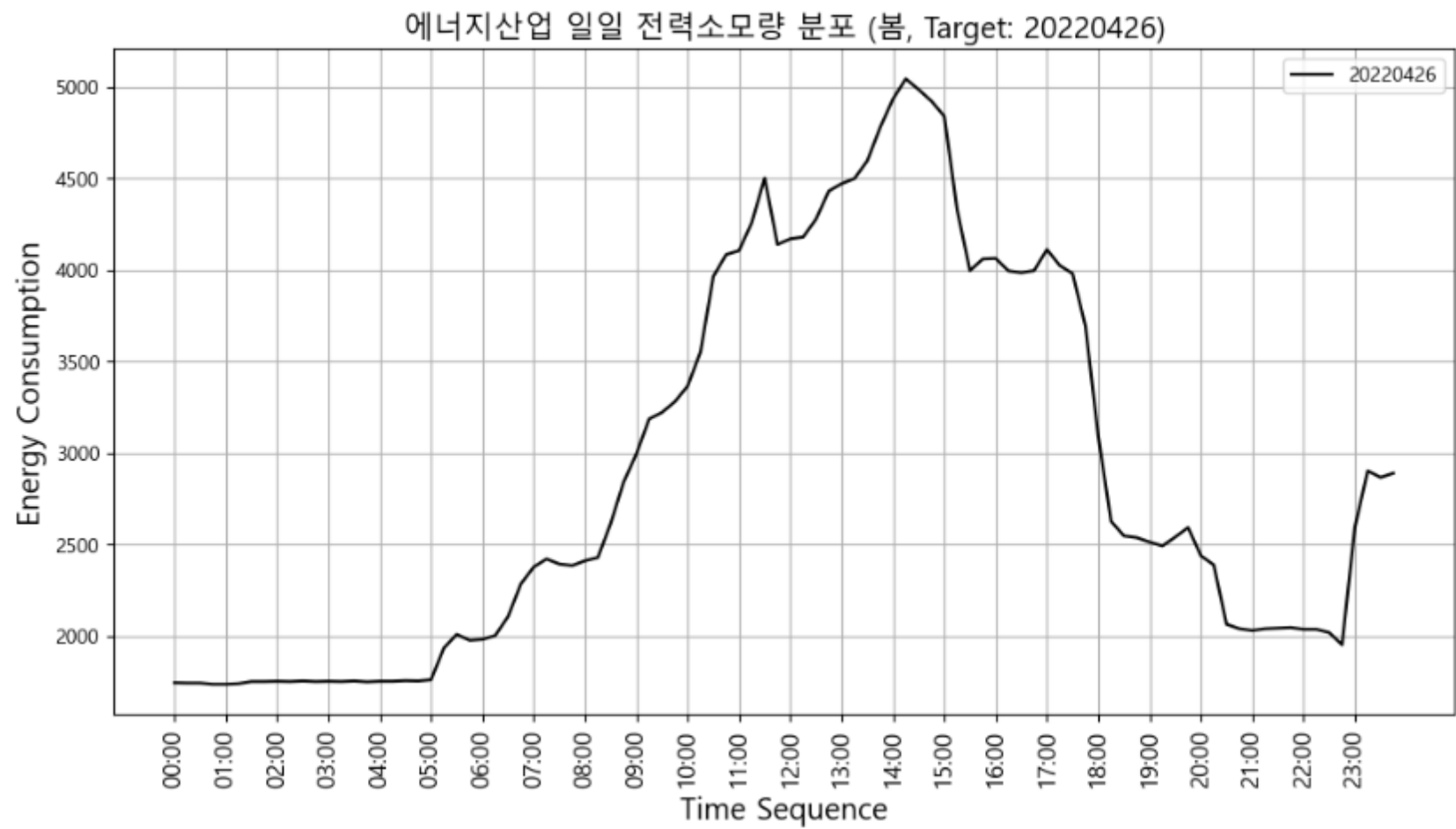
[10, 13, 50, 15, 20]

→ 국소 영역에 존재하는 이상치 판별은 어떻게?

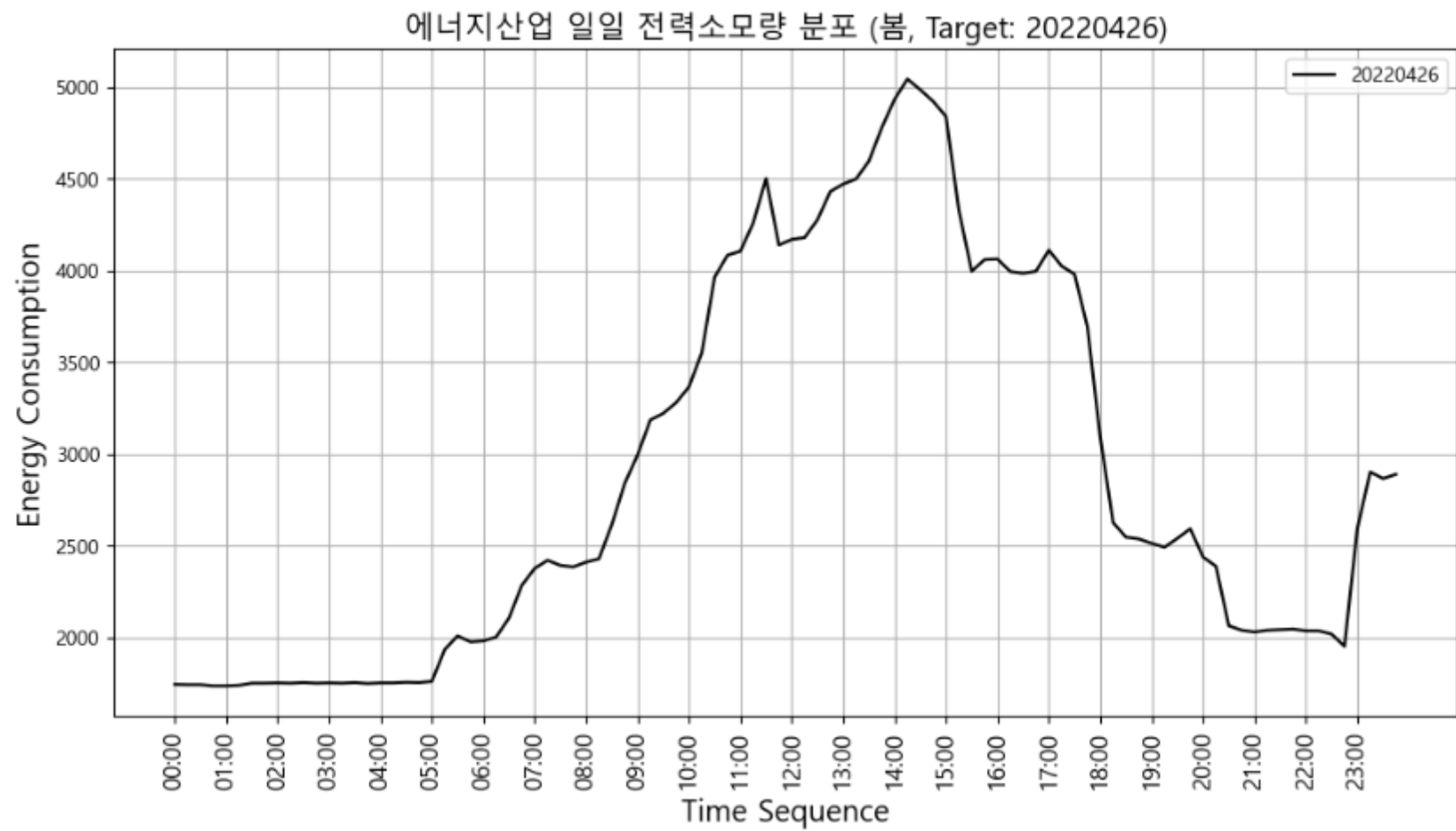
이상치란 무엇인가?

[1744, 1742, 1742, 1736, 1736, 1739, 1750, 1750, 1752, 1750, 1753, 1750, 1752, 1750, 1753, 1749, 1752, 1752, 1755, 1753, 1760, 1933, 2008, 1976, 1980, 2002, 2105, 2284, 2376, 2420, 2392, 2384, 2410, 2428, 2614, 2838, 2994, 3185, 3221, 3279, 3365, 3551, 3963, 4083, 4104, 4258, 4501, 4138, 4169, 4179, 4276, 4430, 4470, 4498, 4595, 4777, 4932, 5043, 4984, 4921, 4838, 4328, 3995, 4059, 4062, 3994, 3984, 3995, 4110, 4023, 3979, 3692, 3091, 2626, 2546, 2536, 2512, 2491, 2540, 2592, 2438, 2387, 2063, 2039, 2029, 2039, 2042, 2045, 2036, 2036, 2018, 1951, 2593, 2901, 2866, 2888]

이상치란 무엇인가?

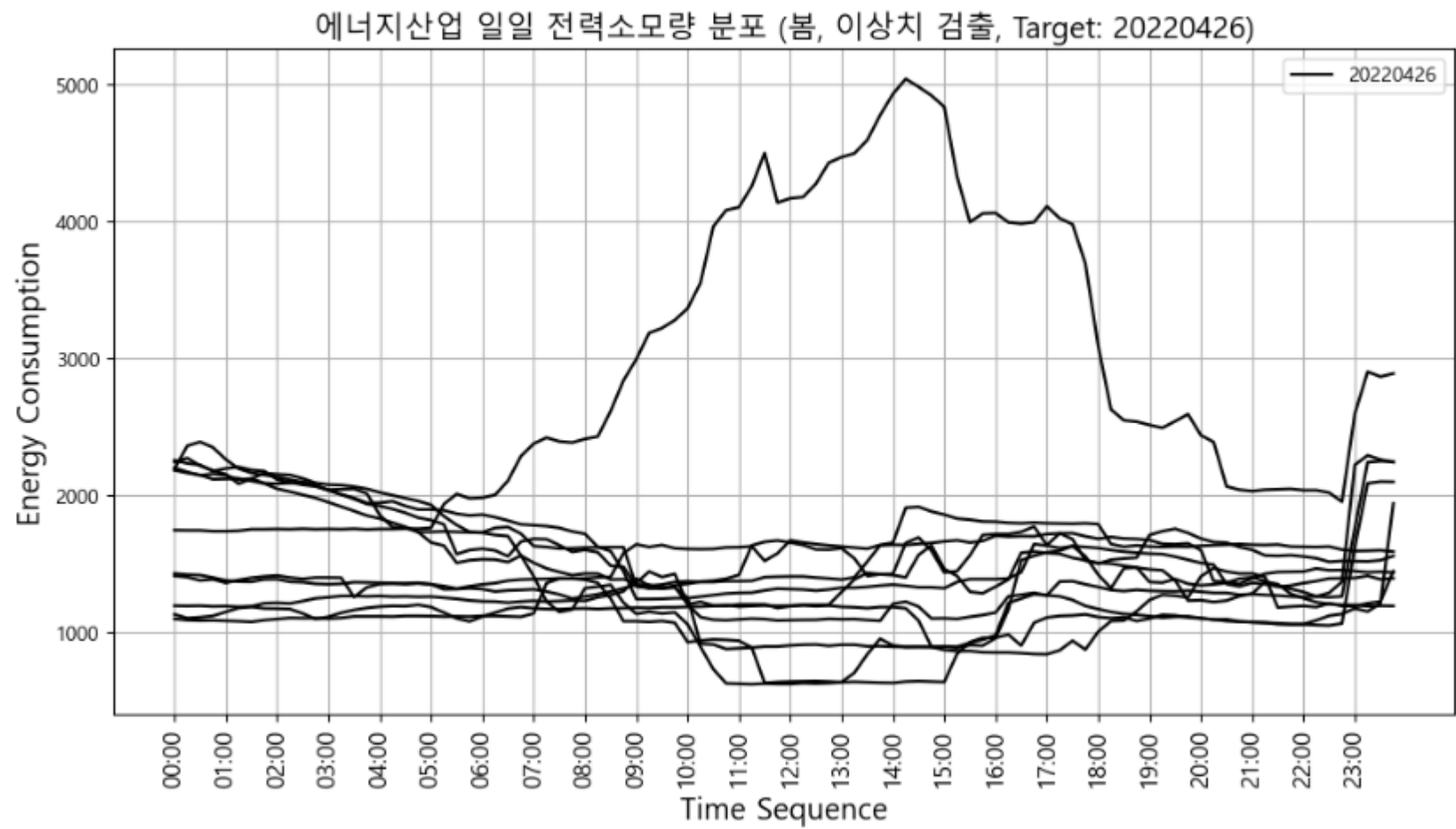


이상치란 무엇인가?

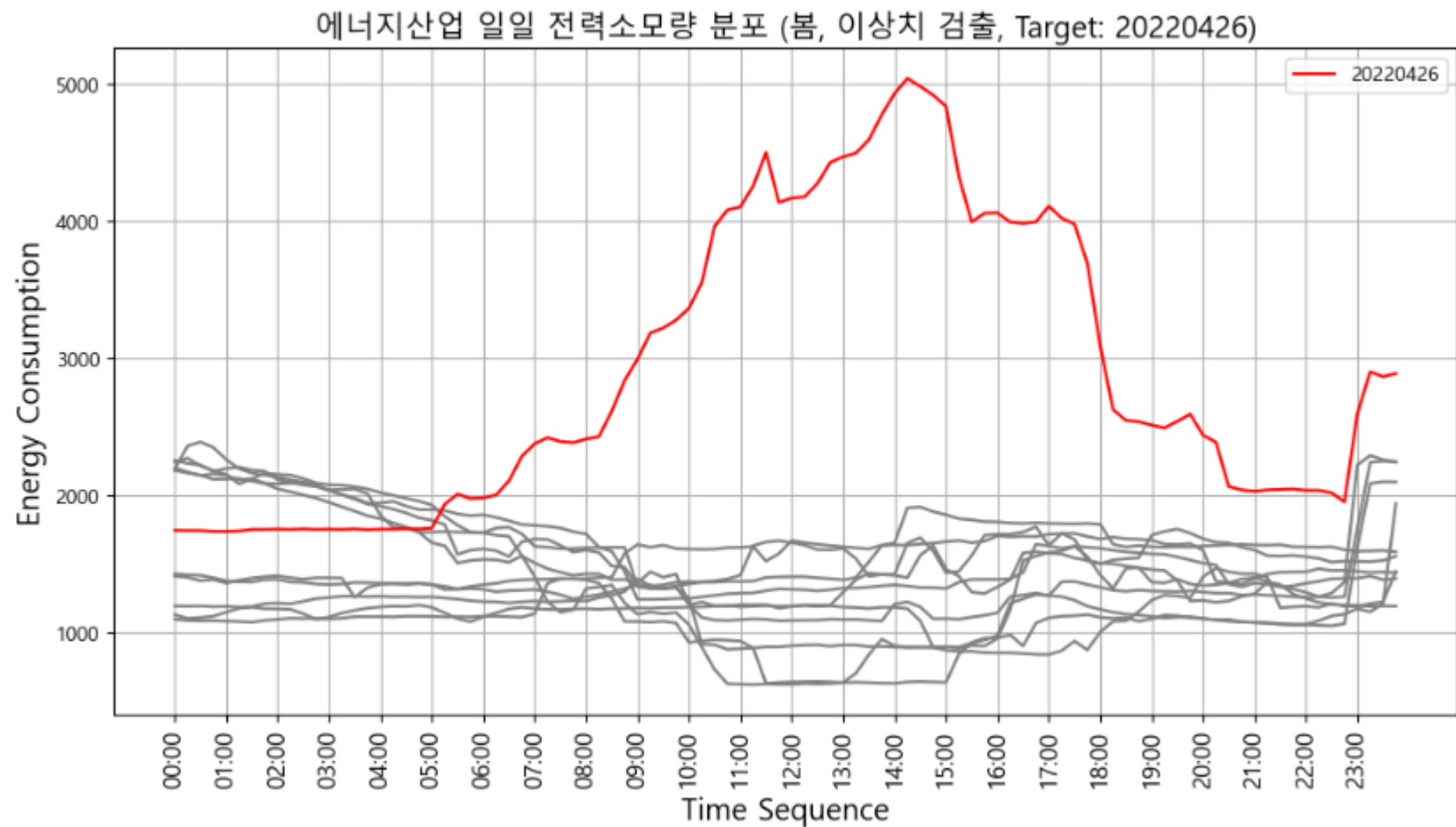


→ 정보의 부족

이상치란 무엇인가?



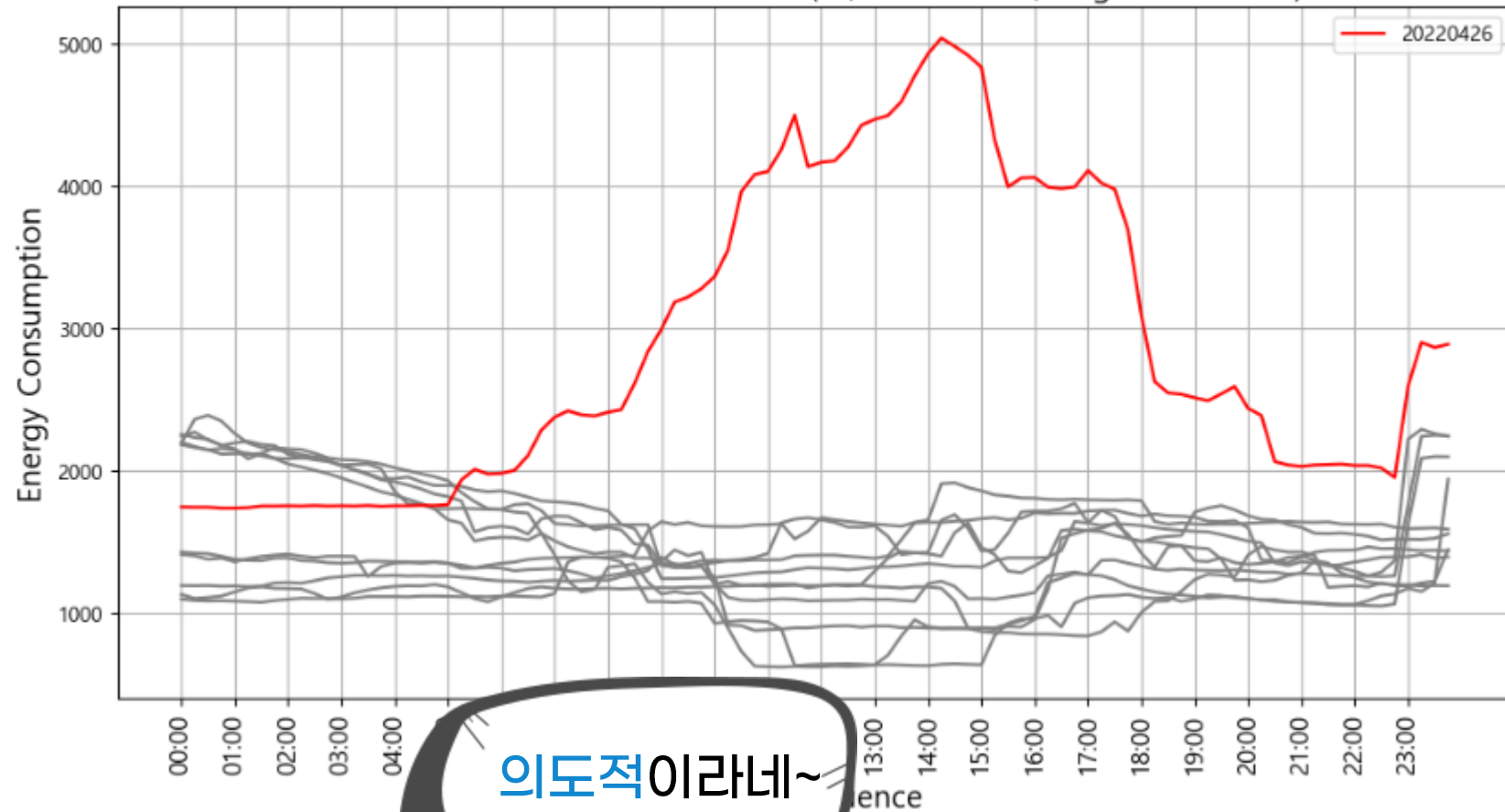
이상치란 무엇인가?



시그마-3에 따라 데이터 분포를 이상치로 규정

이상치란 무엇인가?

에너지산업 일일 전력소모량 분포 (봄, 이상치 검출, Target: 20220426)



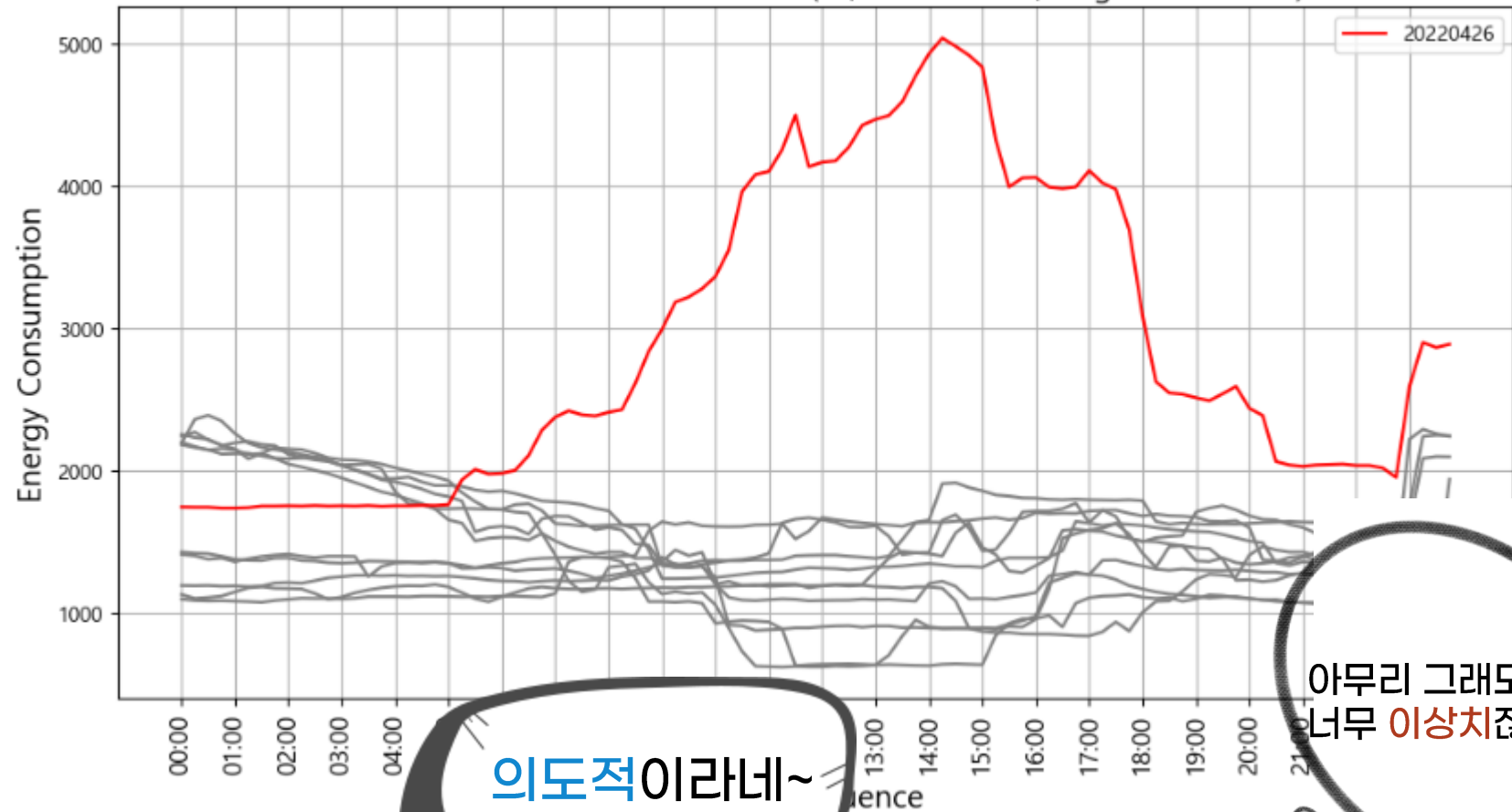
의도적이라네~

시그마-3에 따라 이상치로 규정



이상치란 무엇인가?

에너지산업 일일 전력소모량 분포 (봄, 이상치 검출, Target: 20220426)



아무리 그래도 저건
너무 이상치잖아..

시그마-3에 따라 이상치로 규정



이상치란 무엇인가?

따라서, 사용자와 관리자 간의 관점이 차이가 발생할 수 있음

→ 이상치 판별은 연구자의 몫

- 모호한 데이터 분포 (국소 영역에서의 데이터는 이상치?)

- 인지적 관점의 차이

- 단순 통계로 이상치를 정의하는 것의 어려움

Energy Consumption
50k
40k
30k
20k
10k

도 저건
잖아..



논문 리뷰

논문 리뷰

SPECIAL ISSUE PAPER



A survey on outlier explanations

Egawati Panjei¹ · Le Gruenwald¹ · Eleazar Leal² · Christopher Nguyen¹ · Shejuti Silvia¹

Received: 21 February 2021 / Accepted: 3 December 2021 / Published online: 26 January 2022
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2021

Abstract

While many techniques for outlier detection have been proposed in the literature, the interpretation of detected outliers is often left to users. As a result, it is difficult for users to promptly take appropriate actions concerning the detected outliers. To lessen this difficulty, when outliers are identified, they should be presented together with their explanations. There are survey papers on outlier detection, but none exists for outlier explanations. To fill this gap, in this paper, we present a survey on outlier explanations in which meaningful knowledge is mined from anomalous data to explain them. We define different types of outlier explanations and discuss the challenges in generating each type. We review the existing outlier explanation techniques and discuss how they address the challenges. We also discuss the applications of outlier explanations and review the existing methods used to evaluate outlier explanations. Furthermore, we discuss possible future research directions.

Keywords Outlier explanation · Outlier interpretation · Outlier description · Outlier detection · Anomaly analysis

1 Introduction

Hawkins [40] defines an outlier as “an observation which deviates so much from other observations as to arouse suspicions that it was generated by a different mechanism.” Outliers are also called anomalies, abnormalities, aberrations, contaminants, deviants, discordant observations, exceptions, peculiarities, or surprises in some applications [4,20]. Outlier detection plays an important role in many applications. Identified outliers reveal meaningful information about abnormal behavior in a system. For example, using outlier detection algorithms, medical and public health researchers can identify unusual patient symptoms that can be indicative of medical errors or unusual outcomes [96]. Outlier

For applications to benefit more from the results of the outlier detection process, the results should be explainable. To this end, the process should include two tasks: *outlier detection* and *outlier explanation*. The Merriam-Webster dictionary defines explanation as “the act or process of explaining.” To explain is “to make known” or “to give the reason for or cause of” or “to make something plain or understandable.” In [70], Miller argues that “explainable artificial intelligence can benefit from existing models of how people define, generate, select, present, and evaluate explanation.” In the context of outlier detection, the outlier explanation task provides guidance for users in investigating detected outliers.

Explanations will enhance the users' understanding of outliers and can be used to improve the outlier detection task

저자: Egawati Panjei · Le Gruenwald · Eleazar Leal · Christopher Nguyen · Shejuti Silvia

논문: A survey on outlier explanations

-VLDB Journal(2022.09)등재, 총 32페이지

-데이터를 **사용자** 관점에서 이상치 설명

-각 유형에 대한 **기술적 어려움**에 대한 내용

논문 리뷰

이상치 분석 유형

Importance Levels of Outliers

- 이상치의 우선 순위 또는 중요도 판별
- 조사할 데이터 포인트의 우선순위를 선정 가능
- 점수 기반, 범주를 통한 분류 등의 기법 사용

Casual Interactions Among Outliers

- 어떤 이상치가 추가적인 이상치를 유발했는지를 분석
- 이상치 발생의 원인-결과 구조 설명 가능
- 시간 기반 인과 관계의 연계를 통해 사용

Outlying Attributes of Outliers

- 이상치 발생에 핵심적인 영향을 미친 특성 파악
- 빠른 원인 파악 및 탐색 가능
- 개별 또는 그룹 속성으로 나누어 분석해 설명

논문 리뷰

Importance Level of Outliers

Numerical Ranking of Outliers

- 정상값과 이상치의 편차를 통해 점수 산출
- 점수가 높은 이상치부터 조사

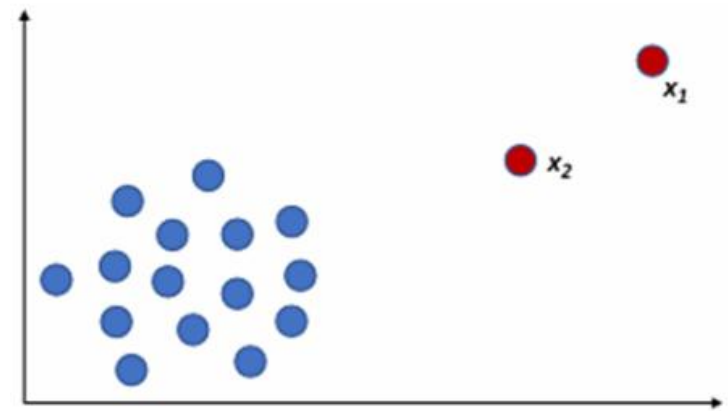


Fig. 1 Example of two outliers in a dataset

Categorical Ranking

- 중요도 혹은 우선 순위에 따른 기준 선정 후 그룹으로 분류
- 분류에 따라 중요도가 높은 순으로 조사

논문 리뷰

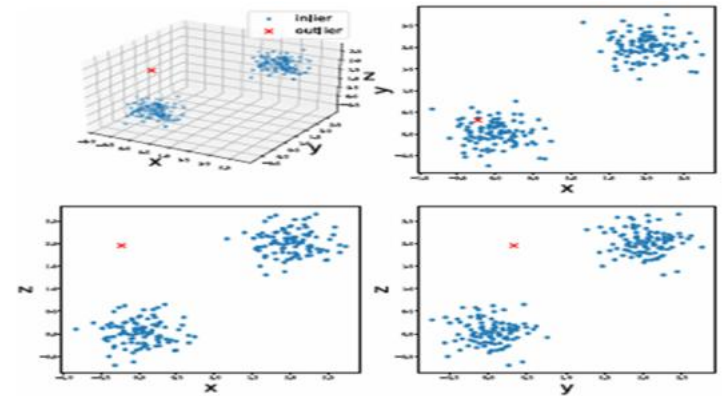
Casual Interactions Among Outliers

Casual Interactions Among Outliers

- 이상치 집합 $O = \{o_i, \dots, o_n\}$ 과
상응하는 timestamp $T = \{t_1, \dots, t_n\}$ 가 필요
- 이상치 o_i 가 o_j 를 유발시켰다고 했을 때 다음과 같은 조건 만족 필요
- o_i 의 timestamp가 o_j 보다 오래되었을 것
- o_i 를 집합 O 에서 제거할 시 o_j 도 제거될 것

논문 리뷰

Outlying Attributes of Outliers



Individual Outliers

- 특정 이상치가 속한 **고차원 공간**에서 이상치를 정상값과 명확히 구분짓는 **하위 속성**을 탐색
- 이상치 점수 기반 해석: 어떤 **속성 조합**에서 이상치가 튀는 지 파악
- 속성 기여도 기반 해석: 각 속성에 대해 **기여도 점수**를 계산 후 특정 임계값보다 큰 속성만을 이상 속성으로 판단

Group of Outliers

- 이상치 그룹에서, 공통적으로 튀는 속성 집합을 탐색
- 그룹 안의 모든 이상치에 대해, 해당 속성 하위집합에서의 **이상치 점수**가 **임계값**보다 크다면, 하위 집합은 이상치 그룹의 공통 이상 속성으로 판단

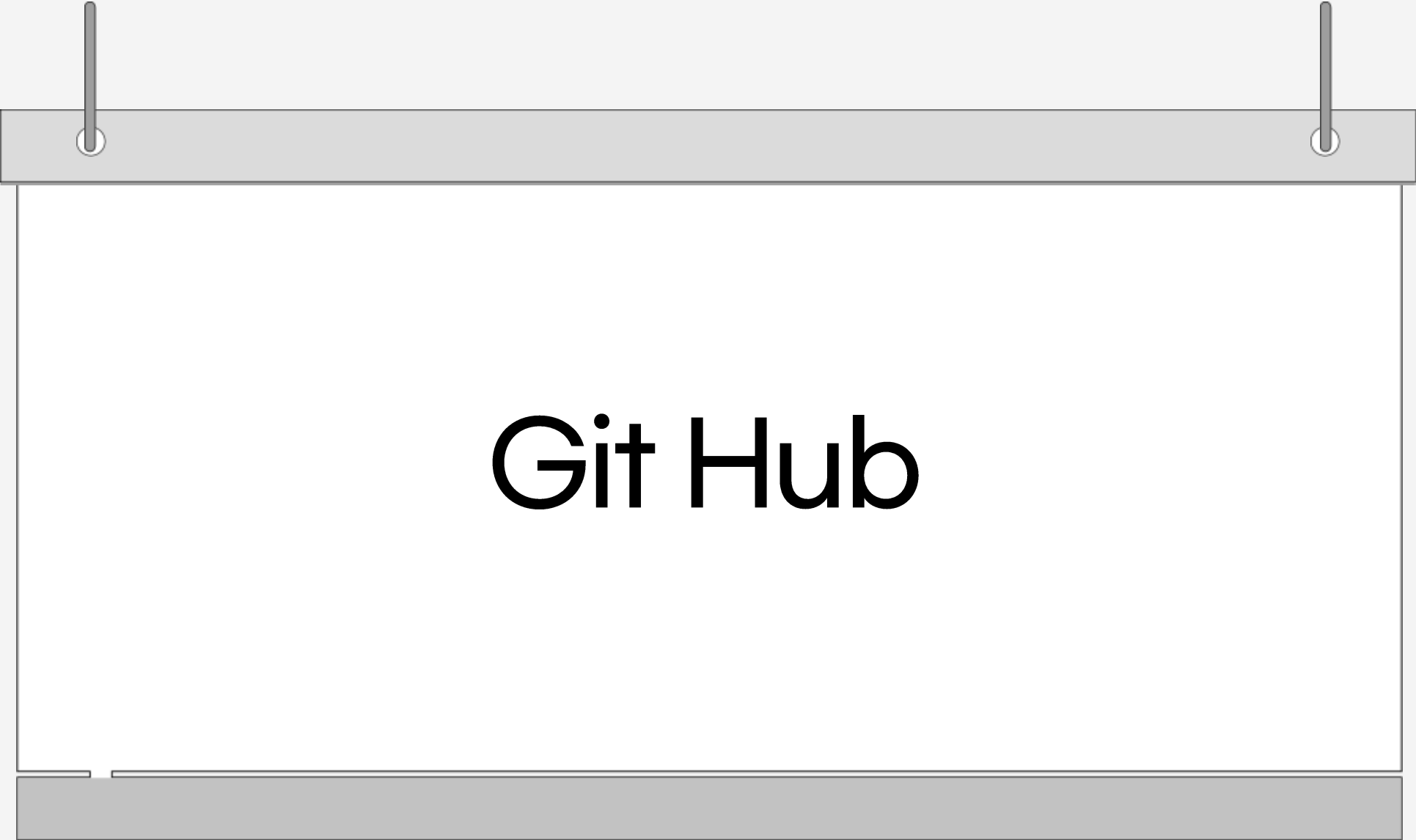
논문 리뷰

향후 연구에 연결 방안

“

논문에서 제안한 사용자 관점에서 적용했던 이상치 탐지
방법을 향후 연구 및 데이터분석에 적용하는 것을 시도해 볼
예정

”



Git Hub

Git Hub

요구사항 명세서

★ 개요

AI를 이용하여 이미지 및 비디오를 생성하는 모델을 구성하고, 웹사이트를 통해 클라이언트에게 해당 콘텐츠(이미지 및 비디오)를 제공하는 서비스

● 핵심 기능 흐름

1. AI 이미지 및 비디오 생성 모델 사용
2. 클라이언트는 API를 통해 POST 및 GET 메서드로 서비스 이용
 - 사용자는 웹사이트를 통해 모델을 선택하고 프론트엔드 입력
 - 주주 추가 기능으로 더 다양한 인터페이스 제공 예정

⚙ 개발 구성

1. AI 생성 모델 모듈화

- 이미지 및 비디오를 생성하는 AI 모델을 설계
- 각 모델은 모듈 단위로 분리하여 관리

2. 한/영 번역기 구성

- 사용자가 한글로 프론트엔드를 입력했을 때 영어로 번역
- 번역 로직은 별도 모듈(hwt.py)로 구성

3. 이미지 생성 모듈 Generate.py

- 모델과 번역기를 사용해 이미지를 생성
- 이미지 생성 함수는 재사용 가능하도록 모듈화하여 유지보수 용이성 확보

4. API 서버 모듈 Generate_service.py

- FastAPI 기반의 웹 서버
- 클라이언트에게 다음 메서드 제공:
 - POST /generate : 이미지 생성 요청
 - GET /models : 사용 가능한 모델 목록 조회
 - GET /image/{image_id} : 생성된 이미지 파일 반환

5. 웹사이트 연동

- 웹사이트에서 클라이언트는 프론트엔드 입력 및 모델 선택 UI를 통해 API 호출
- 백엔드 서버와 연동되어 이미지가 사용자에게 반환됨

젠토 인턴십 프로젝트 일지

프로젝트 개요

FastAPI를 활용한 AI 기반 이미지 및 영상 생성 서비스 테스트용 개발

일별 진행 상황

1일차 - 환경 구성 및 기초 학습

주요 작업

- 개발 환경 구성
 - PyTorch, CUDA, Anaconda, VSCode, Python 설치
- FastAPI 사용법 숙지
- HuggingFace 모델 활용 학습
 - Text-to-Video: Wan2.1
 - 번역: NLLB, opus-mt-ko-en
- 이미지 생성 API 구현
 - 번역 모델: opus-mt-ko-en
 - 이미지 생성: SD v1.5, v3.5M

어려웠던 점 및 느낀 점

학교에서 진행했던 프로젝트들과 다르게 처음부터 환경 구성을 진행해야했음. 허깅페이스와 모델 사용에 대한 경험 미숙으로 제대로 갈피를 잡지 못함. FastAPI 같은 사용자의 요청에 따라 상호작용하는 서버의 작동방식에 대한 이해도가 부족했음.

민성

우현



감사합니다