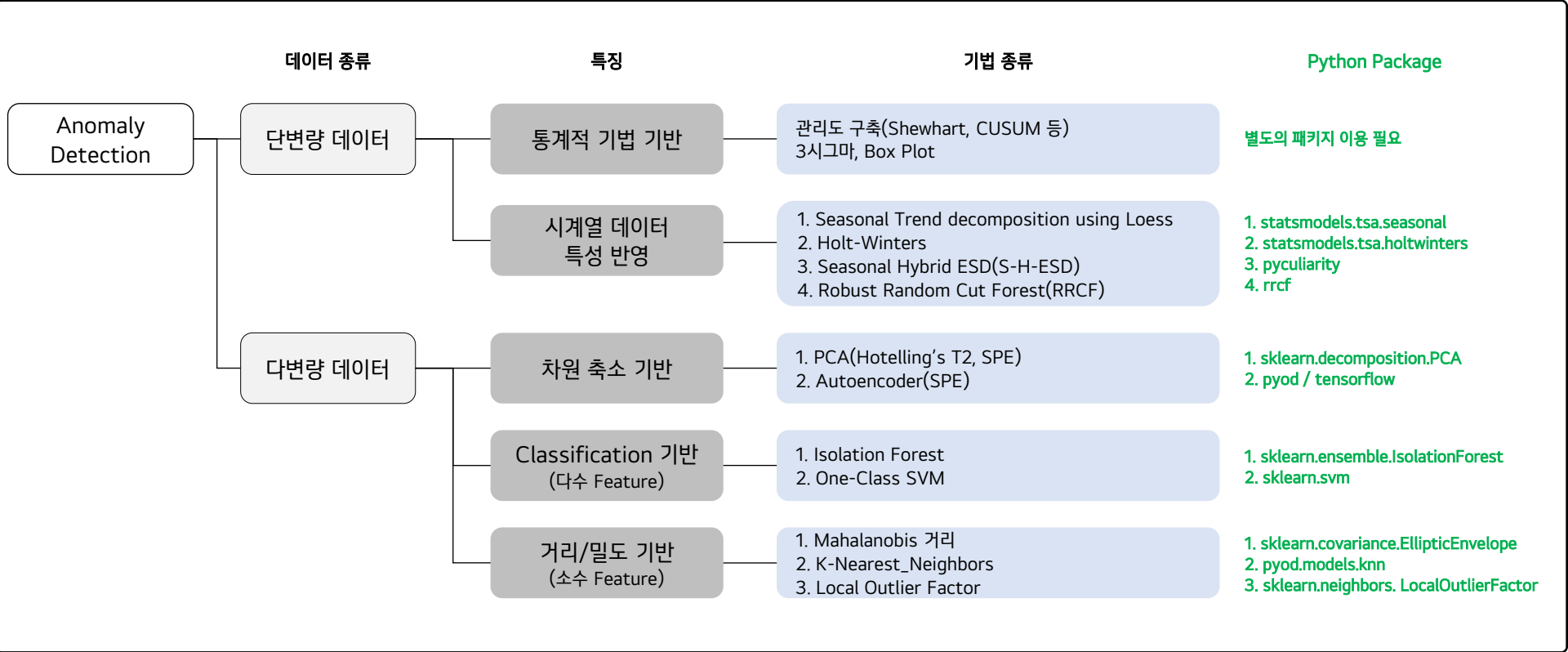


# Anomaly Detection Algorithm Guide

데이터 및 문제 해결 상황에 맞는 이상 탐지(Anomaly Detection) 알고리즘 선택가이드임.  
모든 데이터에 완벽한 이상 탐지 알고리즘은 없으며 알고리즘 별 주요 Hyper parameter 설정을 통해 데이터에 Fit시켜야 함.

## 선택 가이드 구성도



# Anomaly Detection Algorithm Guide

## 선택 가이드 구성도

기법	장점	단점
3 sigma	<ul style="list-style-type: none"> <li>구현이 쉽다.</li> <li>데이터 분포가 변하지 않는 이상 재학습이 필요 없음</li> </ul>	<ul style="list-style-type: none"> <li>정규 분포가 가정되어야 함.</li> <li>Feature간 상관관계 파악이 어려움</li> <li>데이터의 양이 충분하지 않다면, Outlier가 통계치에 영향을 미칠 수도 있음.</li> </ul>
Box Plot	<ul style="list-style-type: none"> <li>데이터를 눈으로 확인하기 어려울 때, 그림을 이용해 데이터의 범위를 빠르게 파악할 수 있음.</li> </ul>	<ul style="list-style-type: none"> <li>가운데 선은 평균이 아님.</li> <li>오해 소지가 있음(Median과 Mean은 다름.)</li> </ul>
STL	<ul style="list-style-type: none"> <li>시계열 데이터에 대해 분기별, 월별, 일별 분해 모두 가능</li> <li>MA(Moving Avg) 방식이 아니기 때문에 데이터 유실 없음.</li> <li>돌발스런 이상치에 대해 추세, 주기에 영향을 미치지 않음.</li> </ul>	<ul style="list-style-type: none"> <li>시간 데이터 전처리 필수적임.</li> <li>덧셈 분해 기능만 제공</li> </ul>
Holt-winter	<ul style="list-style-type: none"> <li>연산량이 적음, 큰 데이터 세트에 대해서 리소스 절약, 자동화 가능</li> </ul>	<ul style="list-style-type: none"> <li>변량 데이터에 대해서만 적용 가능, 상관관계 고려 X</li> <li>계절성이 없는 데이터에 대해서는 성능 저조</li> <li>변동이 적은 계절성 데이터에 대해서는 민감하게 탐지할 우려</li> </ul>
S-H-ESD	<ul style="list-style-type: none"> <li>데이터 내의 노이즈에 어느 정도 대응할 수 있음.</li> <li>급작스런 상승하는 이상치를 탐지할 수 있음</li> </ul>	<ul style="list-style-type: none"> <li>이상 탐지가 안 되는 경우가 아래와 같이 존재함. <ul style="list-style-type: none"> <li>- 점진적 증가 신호(seasonal grow)</li> <li>- 점진적 증가하는 신호에서의 음의방향 이상치 (Negative seasonal anomaly)</li> <li>- 평면적 신호 (Flat signal)</li> </ul> </li> </ul>
RRCF	<ul style="list-style-type: none"> <li>S-H-ESD에서 탐지되지 않았던 경우에 어느 정도 대응 가능함.</li> <li>Batch 및 Streaming data 모두 활용 가능함.</li> <li>Subsampling을 통한 적은 연산량</li> </ul>	<ul style="list-style-type: none"> <li>분리를 위한 선을 수직과 수평으로만 자르기 때문에 잘못된 scoring이 발생할 수 있음</li> </ul>
PCA-Hotellings'T2	<ul style="list-style-type: none"> <li>고차원에 데이터의 특징(잠재변수)을 추출할 수 있음.</li> <li>선택한 변수들의 해석이 용이함.</li> <li>PCA의 잠재변수 축을 이용하여 거리기반의 이상탐지를 실시하기 때문에 직관적임.</li> </ul>	<ul style="list-style-type: none"> <li>Threshold에 따라 민감함.</li> </ul>
Autoencoder Reconstruction	<ul style="list-style-type: none"> <li>데이터 Label이 존재하지 않아도 사용 가능</li> <li>고차원에 데이터의 특징(잠재변수)을 추출할 수 있음.</li> <li>Auto encoder를 기반으로 다양한 알고리즘 존재</li> </ul>	<ul style="list-style-type: none"> <li>Hyper parameter (※ hidden layer) 설정이 어려움.</li> <li>Loss(Reconstruction Error)에 대한 threshold 설정이 어려움.</li> </ul>

# Anomaly Detection Algorithm Guide

## 선택 가이드 구성도

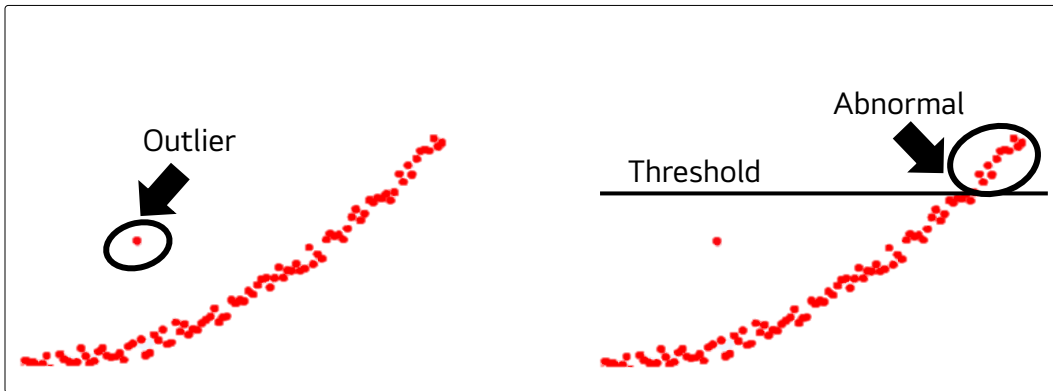
기법	장점	단점
Isolation Forest	<ul style="list-style-type: none"> <li>• 군집기반 이상탐지 알고리즘에 비해 계산량이 매우 적음. (※ Sampling 사용 Tree 생성)</li> <li>• Anomaly Detection 성능 우수</li> <li>• Train Data에 이상치가 없어도 Test Data에서 우수한 성능을 보임.</li> </ul>	<ul style="list-style-type: none"> <li>• 분리를 위한 선을 수직과 수평으로만 자르기 때문에 잘못된 scoring이 발생할 수 있음 (※ 대안 방법: Extended Isolation Forest)</li> </ul>
One-Class SVM	<ul style="list-style-type: none"> <li>• 데이터 Label이 존재하지 않아도 사용 가능</li> <li>• 저차원이나 고차원의 적은 데이터에서 일반화 능력이 좋음 (데이터 특성이 적어도 성능이 좋게 나오는 편, Robust 함)</li> </ul>	<ul style="list-style-type: none"> <li>• Kernel 기반의 방법론, 데이터가 늘어날 수록 연산량이 크게 증가함.</li> <li>• Scaling과 Hyper parameter에 민감함.</li> </ul>
Mahalanobis 거리	<ul style="list-style-type: none"> <li>• 비선형 관계의 데이터에 활용 가능</li> <li>• 데이터에 자체에 대한 가정이 필요 없음</li> </ul>	<ul style="list-style-type: none"> <li>• 변수 간의 관계가 모두 독립이라면 유클리드 거리와 같은 개념</li> <li>• 변수 간의 상관성이 명확히 알려져 있지 않은 경우 적용하기 어려움</li> </ul>
K-Nearest_Neighbors	<ul style="list-style-type: none"> <li>• 단순하고 효율적임</li> <li>• 기존 분류 체계 값을 모두 검사하여 비교하므로 높은 정확도를 보임.</li> <li>• 수치 기반 데이터 분류 작업에서 성능 우수함.</li> <li>• 기존 데이터를 기반으로 하기 때문에 데이터에 대한 가정이 없음.</li> </ul>	<ul style="list-style-type: none"> <li>• 기존의 모든 데이터를 비교해야 하기 때문에 데이터가 많으면 많을 수록 처리 시간이 증가</li> <li>• 특징과 클래스간 관계를 이해하는데 제한적</li> <li>• 카테고리컬 데이터를 위한 추가 처리가 필요</li> </ul>
Local Outlier Factor	<ul style="list-style-type: none"> <li>• 굉장히 밀집한 클러스터에서 조금만 떨어져 있어도 이상치로 탐지</li> <li>• KNN과 다르게 특별한 라벨링이 없어도 사용할 수 있음</li> <li>• Local Outlier를 탐지할 수 있음</li> </ul>	<ul style="list-style-type: none"> <li>• 데이터의 차원수가 증가할 수록 연산량 증가</li> <li>• 이상치 판단 기준 설정 어려움 (밀집도가 다른 여러 클러스터가 존재한다면 민감하게 반응함)</li> </ul>

# Outlier vs Abnormal vs Novelty

이상 탐지를 학습 데이터에 따라 아래와 같이 나눠 볼 수 있으며 학습할 데이터를 어떻게 정의하는가에 따라서 문제의 성격과 해결 방법이 다름.  
새로운 관측치가 기존 분포에 속하는지 기존 분포를 벗어나는지 구분함

- Outlier: 관측된 데이터의 범위에서 많이 벗어난 아주 작은 값이나 큰 값으로 분석 결과 해석 시 오해를 발생시킬 수 있기 때문에 사전 제거 필요
- Abnormal: Domain-Knowledge 기반의 문제 상황 범주에 속한 데이터
- Novelty: 학습된 데이터 외의 새로운 패턴의 데이터

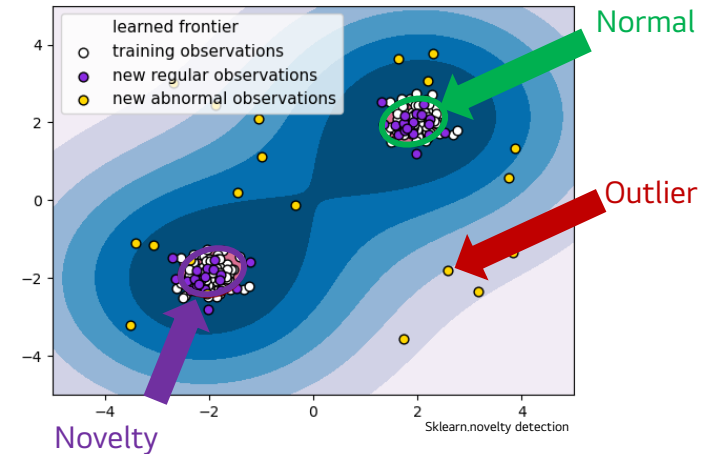
## Outlier vs Abnormal



## Outlier/Abnormal Detection

학습 데이터를 통해 정상 데이터의 범위를 결정하고, 이를 초과할 경우 Abnormal로 간주  
훈련 데이터 셋에 정상 샘플과 이상치 샘플을 모두 포함하고 있음.

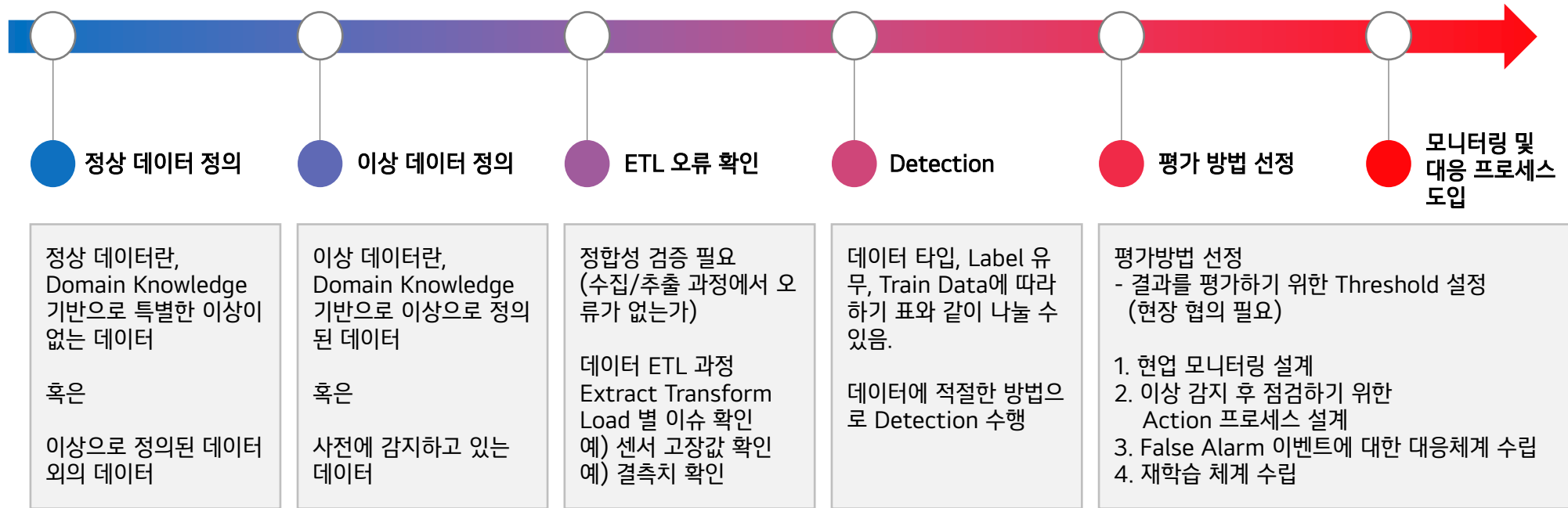
## Outlier vs Abnormal vs Novelty



## Novelty Detection

학습 데이터 내에 존재하지 않는 패턴의 데이터를 탐지  
탐지 목표 1: In-distribution Test Set의 정확한 예측  
탐지 목표 2: Out-of-distribution 데이터 셋은 걸러내는 것

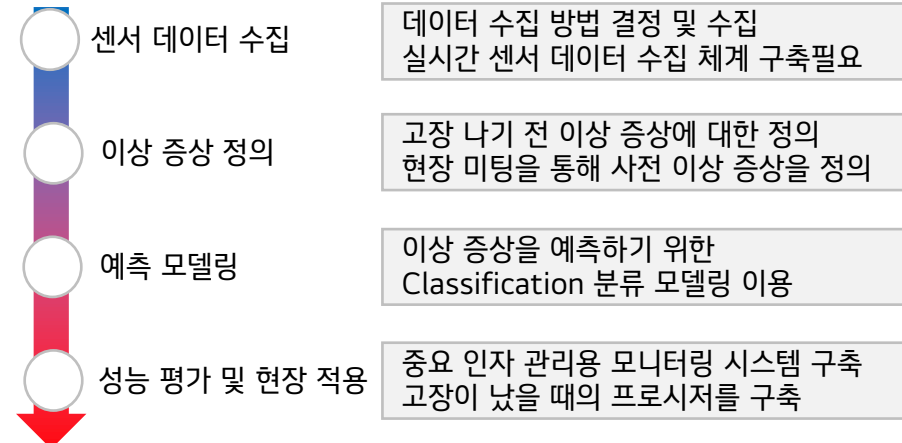
# Anomaly Detection Process (1/2)



## Anomaly Detection의 종류

학습 유형	패턴 유형	데이터 유형	데이터 종류
지도 학습 (Supervised)	이상 탐지 (Outlier/Abnormal Detection)	단변량 데이터 (Univariate)	시계열 데이터 (Time-Series)
반 지도 학습 (SemiSupervised)	신규성 이상 탐지 (Novelty Detection)	다변량 데이터 (Multivariate)	비 시계열 데이터 (Non Time- Series)
비지도 학습 (Unsupervised)			

## 예시 프로세스 - 설비 이상 탐지



# Anomaly Detection Process (2/2)

## 예시 프로세스 – 품질 이상 탐지 모니터링



## 기대 효과 산출

주요 공정 변수 모니터링 전/후 생산량 비교

주요 공정 변수 모니터링 전/후 불량률 비교

주요 공정 변수 모니터링 전/후 폐기 비용 비교

# Anomaly Detection Project

프로젝트 이해관계를 알고 관계 간 적절한 관리가 필요함.

널리 알려진 시행착오와 프로젝트 진행을 지연/방해하게 되는 요소들을 이해하고 사전에 해당 내용을 점검, 프로젝트의 불확실성을 줄여야 함.

프로젝트 일정 수립 시 머신러닝 프로젝트의 불확실성\*을 고려하여 범위의 확정, 변경되지 않는 기획이 필요함.

머신러닝 프로젝트의 불확실성: 새로운 아이디어, 반복적인 실험, 예측할 수 없는 결과로 가시적인 결과를 예상하기 어려움

프로젝트 이해관계	비즈니스 관점의 점검 사항	데이터 관점의 점검 사항
1. 이상 탐지 필요 부서(제조 공정 현장) 2. 데이터 사이언스/TF팀 3. IT/외주개발 4. CEO/임원	1. 머신러닝 문제 정의 2. 요구사항 점검(KPI, 예측 주기 등) 3. 가용 가능한 자원/환경 점검 4. 선행 모델의 장/단점 파악	1. 데이터의 위치/권한 확인 2. 이상 감지 여부 확인 - 이상 데이터 포함 여부/발생주기 - 이상에 영향을 주는 Feature 유무

## 프로젝트 일정 분배 예시

