

# 적대적 생성 신경망을 활용한 비지도 학습 기반의 대기 자료 이상 탐지 알고리즘 연구

양호준<sup>1</sup>, 이선우<sup>1</sup>, 이문형<sup>1</sup>, 김종구<sup>1</sup>, 최정무<sup>2</sup>, 신유미<sup>2</sup>, 이석채<sup>3</sup>, 권장우<sup>4\*</sup>, 박지훈<sup>5</sup>, 정동희<sup>5</sup>, 신혜정<sup>5</sup>  
<sup>1</sup>인하대학교 전기컴퓨터공학과 학생, <sup>2</sup>인하대학교 컴퓨터공학과 학생, <sup>3</sup>인하대학교 행정학과 학생,  
<sup>4</sup>인하대학교 컴퓨터공학과 교수, <sup>5</sup>국립환경과학원 대기환경연구과 환경연구원

## A Study on Atmospheric Data Anomaly Detection Algorithm based on Unsupervised Learning Using Adversarial Generative Neural Network

Ho-Jun Yang<sup>1</sup>, Seon-Woo Lee<sup>1</sup>, Mun-Hyung Lee<sup>1</sup>, Jong-Gu Kim<sup>1</sup>, Jung-Mu Choi<sup>2</sup>, Yu-mi Shin<sup>2</sup>,  
Seok-Chae Lee<sup>3</sup>, Jang-Woo Kwon<sup>4\*</sup>, Ji-Hoon Park<sup>5</sup>, Dong-Hee Jung<sup>5</sup>, Hye-Jung Shin<sup>5</sup>

<sup>1</sup>Student, Department of Electric Computer Engineering, Inha University

<sup>2</sup>Student, Department of Computer Engineering, Inha University

<sup>3</sup>Student, Department of Public Administration, Inha University

<sup>4</sup>Professor, Department of Computer Engineering, Inha University

<sup>5</sup>Researcher, Air Quality Research Department, Air Quality Research Division

**요약** 본 논문에서는 기존에 전문가에 의해서 이루어지던 국가 대기오염 측정망 데이터들의 이상 탐지 작업을 인공지능을 통해 자동화하고자 심층 신경망을 이용한 이상 탐지 모델을 제안하였다. 환경과학원에서 제공받은 기상자료 데이터의 결측치 및 이상치를 분석하여 학습데이터를 생성하였으며 비지도 학습 방식의 BeatGAN 모델에 기반하여 커널 구조 변경과 합성곱 필터층 및 전치 합성곱 필터층의 추가를 통해 새로운 모델을 제안하여 이상 탐지 성능을 높이고자 하였다. 또한 제안하는 모델의 생성적 특징을 활용하여 새로운 데이터를 생성하고 이를 학습에 사용하는 재학습 알고리즘을 구현 및 적용하여 기존 BeatGAN 모델뿐 아니라 다른 비지도 학습 모델인 Iforest, One Class SVM과 비교하였을 때 제안모델의 성능이 가장 높았음을 확인할 수 있었다. 본 연구를 통해 실제 산업현장에서 센서의 이상, 점검 등의 여러 요인으로 인해 학습 데이터가 부족한 상황에서 추가적인 비용없이 과적합을 피하며 제안하는 모델의 이상탐지 성능을 올릴 수 있는 방법을 제시할 수 있었다.

**주제어** : 대기질, 머신러닝, 딥러닝, 비지도 학습, 이상탐지

**Abstract** In this paper, We propose an anomaly detection model using deep neural network to automate the identification of outliers of the national air pollution measurement network data that is previously performed by experts. We generated training data by analyzing missing values and outliers of weather data provided by the Institute of Environmental Research and based on the BeatGAN model of the unsupervised learning method, we propose a new model by changing the kernel structure, adding the convolutional filter layer and the transposed convolutional filter layer to improve anomaly detection performance. In addition, by utilizing the generative features of the proposed model to implement and apply a retraining algorithm that generates new data and uses it for training, it was confirmed that the proposed model had the highest performance compared to the original BeatGAN models and other unsupervised learning model like Iforest and One Class SVM. Through this study, it was possible to suggest a method to improve the anomaly detection performance of proposed model while avoiding overfitting without additional cost in situations where training data are insufficient due to various factors such as sensor abnormalities and inspections in actual industrial sites.

**Key Words** : Air Quality, Machine Learning, Deep Learning, Unsupervised Learning, Anomaly Detection

\*This work was supported by a grant from the National Institute of Environmental Research (NIER), funded by the Ministry of Environment (MOE) of the Republic of Korea.

\*This research was supported by the BK21 Four Program funded by the Ministry of Education(MOE, Korea) and National Research Foundation of Korea(NRF).

\*Corresponding Author : Jang-Woo Kwon(jwkwon@inha.ac.kr)

Received February 20, 2022

Revised March 20, 2022

Accepted April 20, 2022

Published April 28, 2022

## 1. 서론

점점 더 심해지고 있는 지구 온난화, 환경오염, 미세먼지 증가를 막기 위해 유엔환경계획, 그린피스 등의 국제 환경단체는 환경문제 조정기능 및 촉매 기능 유지, 환경 상태 평가 및 환경관리, 환경 보호를 위한 지원조치 등을 진행하고 있으며, 국내에서도 환경운동연합, 녹색연합, 세계자연기금 한국본부 등의 단체들이 환경오염의 심각성과 경각심을 강화 시키기 위한 움직임을 보이고 있다. 이러한 활동들에 의해 국민들의 환경과 미세먼지 등에 대한 관심이 높아지게 되자, 전국의 국가대기오염측정망도 '10년 290개소에서 '20년 505개로 늘어났으며 계속해서 확장되고 있는 추세이다.

측정소가 확장되고 있는 만큼 내부 시설의 공사, 장치의 결합과 점검 혹은 주변지대의 개발, 기후 등 예측하기 힘든 경우가 증가될 수 있으며, 이러한 상황들은 측정소 내의 대기오염물질 센서값에 결측치 혹은 이상치와 같이 데이터의 신뢰도를 하락시킬 수 있는 문제의 원인이 될 수 있다.

이와 같은 사태를 미연에 방지하기 위해 국가대기오염측정망 측정소의 전문가들은 정기적으로 데이터의 결측치, 이상치에 해당하는 값 혹은 구간에 표시(Labeling)하는 작업을 수행하고 있다. 이러한 작업적 특성 때문에 표시가 끝난 데이터는 전문가 각각의 주관적 판단이 들어있을 가능성이 높아지게 된다. 기준값을 초과하거나, 동일값이 지속되는 구간이 있는 경우 단순한 알고리즘으로 자동화가 가능한 것처럼 보이나, 특정 측정소에서는 기준값이 달라질 수 있다는 점, 동일값 지속이 허용될 수 있다는 점 등을 주관적 판단사례의 예시로 들 수 있다. 여기서 발생할 수 있는 문제점은 측정소마다 상이할 수 있는 판단기준에 대한 명확한 근거가 없이 전문가들마다 그동안 쌓아온 경험에 의해 작업을 수행하게 된다는 점이다. 점점 늘어나고 있는 측정소로 인하여 검증해야 하는 데이터가 더욱 늘어나고 있으며 이에 따라 전문가들이 검토해야 하는 데이터의 양 또한 늘어나고 있다.

본 논문에서는 이러한 문제를 해결하기 위한 방법으로 심층 신경망을 활용한 기법중 하나인 이상탐지(Anomaly Detection)기법을 적용하는 것을 제안하고자 한다.

## 2. 관련 연구

시계열 데이터에 있어서 이상치는 데이터 전체의 신뢰도에 악영향을 주는 만큼 이상 탐지에 대한 연구는 상당한 중요도를 가지고 관심 있게 연구 되어왔다. Li, Izakian 등은[1]은 클러스터링 기반 이상치 탐지 기술을 이용하여 다변량 시계열 데이터(Multivariate time series) 이상 탐지를 수행하는 연구를 진행하였다. 다변량 시계열 데이터에 Sliding Window 기술을 이용하여 작은 부분으로 나눈 뒤 FCM 클러스터링을 이용하여 구조를 나타내고 나타난 클러스터를 기준으로 이상치에 확률을 부여하여 이상 탐지를 하는 방법을 제안하였다.

Deng, Hooi는 시계열 데이터 이상 탐지를 위해 GNN(Graph Neural Network) 모델과 임베딩 벡터를 이용하여 정상 데이터로만 이루어진 학습 데이터 셋의 특성을 포착한 뒤 각 데이터 간의 의존 관계를 분석하는 그래프 구조를 학습하여 학습한 의존 관계에 부합하지 않는 데이터를 이상 데이터로 진단하는 방법을 제안하였다.[2,3]

Ren, Xu는 돌출 감지(Spectral Residual)를 이용하여 1차 시계열 데이터에 대한 학습을 진행하고 CNN과 연계하여 이상 데이터를 찾아내는 방법을 제안하였다[4].

위에서 언급한 방법론들은 충분한 학습 데이터가 확보 되어야 제 성능이 발휘될 수 있다는 점을 전제로 하고 있다. 그러나 본 연구와 같이 실제 현상에서의 센서값을 데이터로 이용하는 경우, 기기의 이상 혹은 점검 등의 이유로 정상 데이터를 충분한 만큼 확보하기 어려운 경우가 있을 수 있으며 본 논문에서는 이러한 상황에서도 적정 성능 이상의 이상 탐지를 진행하기 위해 데이터를 스스로 생성해 내서 학습 데이터를 보충할 수 있는 생성 모델(Generative model) 중심의 딥러닝 방법론을 제시하고자 하였다.[5]

## 3. 실험 방법

### 3.1 데이터 분석

Table 1은 대기오염측정망을 통해 측정되는 항목들의 자동측정 주기와 자료생성 주기에 관한 표이며[6] 본 논문에서는 위의 8가지 항목(SO<sub>2</sub>, CO, O<sub>3</sub>, NO, NO<sub>2</sub>, NOX, PM<sub>10</sub>, PM<sub>2.5</sub>)에 대하여 2018년 4월 1일부터 2019년 12월 31일까지 총 1년 9개월동안 축적된 국립 환경 과학원의 이상치 판정 결과를 실험 데이터로 사용하였다.

Table 1. The measurement period of the automatic measurement item

Item	Cycle	
	Automatic measurement	Data generation
SO2 CO O3 NO NO2 NOX	5 min	Average value for an hour
PM10 PM25	1 hour	Concentration value for an hour

딥러닝 기반의 알고리즘 구현 및 실험을 진행한 컴퓨터 사양은 Table 2와 같다.

Table 2. Experiment environment

Components	Name
Operating system	Windows 10
CPU	Intel i9-9900K
GPU	NVIDIA GeForce RTX 3090
Memory size	64GB

입력 데이터에 의존하여 스스로 특징을 도출해가며 학습을 진행하는 딥러닝 알고리즘의 특성상 데이터의 품질 및 신뢰도는 학습 이후의 성능과 직결될 정도로 매우 중요하다. 따라서 원활한 학습을 진행하기 위한 적절한 전처리 방법 선택을 위해 제공받은 데이터의 이상치와 결측치의 비율 및 종류에 대해서 분석하였다.

이상치란 국립 환경 과학원에서 정상으로 간주하기 어려운 데이터를 동일값 지속, 과거 대비 이상 등의 9개 목록 중 하나에 속한다고 판단되는 값을 의미하며 결측치란 측정 센서의 점검, 기기 이상 등으로 제대로 된 측정이 진행되지 않아 값 자체가 비어있는 것을 의미한다.

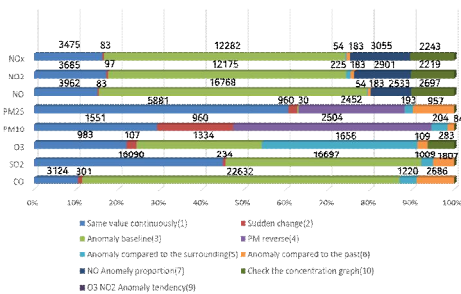


Fig 1. Analyzing data in terms of Anomaly

Fig. 1은 각 측정 항목별 이상치들의 수를 나타낸 그림이며 각 항목 대비 이상치 및 결측치 비율을 Table 3에서 확인할 수 있다.

Table 3. Rate of Anomaly data

Item	Proportion of anomaly (%)	Proportion of missing (%)
SO2	4.79	1.35
CO	5.12	1.46
O3	4.39	0.53
NO	4.04	0.72
NO2	4.05	0.71
NOX	4.01	0.72
PM10	5.48	0.31
PM25	9.25	3.17

### 3.2 학습 데이터 생성

학습 성능을 올리기 위한 전처리 과정으로 위에서 언급한 이상치 및 결측치를 제거하기 위해 Fig. 2 알고리즘을 적용하였다.

Algorithm 1 Data preprocessing

```

1: Get StepSize of data
2: for iteration = 1, 2, ..., DataSize do
3:   Extract Data according to the StepSize
4:   Check DataCode of same range of extracted Data
5:   if Whole DataCode equal to 0 then
6:     Save the Data in NormalDataArray
7:   else
8:     Save the Data in AnormalDataArray
9:   end if
10: end for
    
```

Fig 2. Applied preprocessing algorithm

센서값 데이터와 같은 범위를 공유하는 이상 코드 (WrongCode) 데이터에서 한 구간이라도 정상이 아닐 시 해당 데이터를 비정상적으로 간주하여 학습 데이터에서 제외했다.

부족한 데이터셋을 최대한 활용하기 위해 결측치 값들은 결측치를 중심으로 일정 구간 까지의 평균값을 구하여 대체하는 방식을 사용하였으며, 주어진 시계열 데이터에 대해 딥러닝 모델이 빠른 속도로 최적의 해 (Optimal Solution)를 찾게 하고, 인공지능 모델이 생성한 값과의 비교를 통한 이상치 점수 (Anomaly Score)를 구하기 위해 수식 1 형태로 데이터가 최소 - 1에서 최대 1 값을 갖도록 최소·최대 정규화 (Min-Max Normalization) [7]를 적용하였다

$$X_{\neq w} = 2 \times \frac{X - X_{\min}}{X_{\max} - X_{\min}} - 1 \quad (1)$$

### 3.3 제안하는 모델

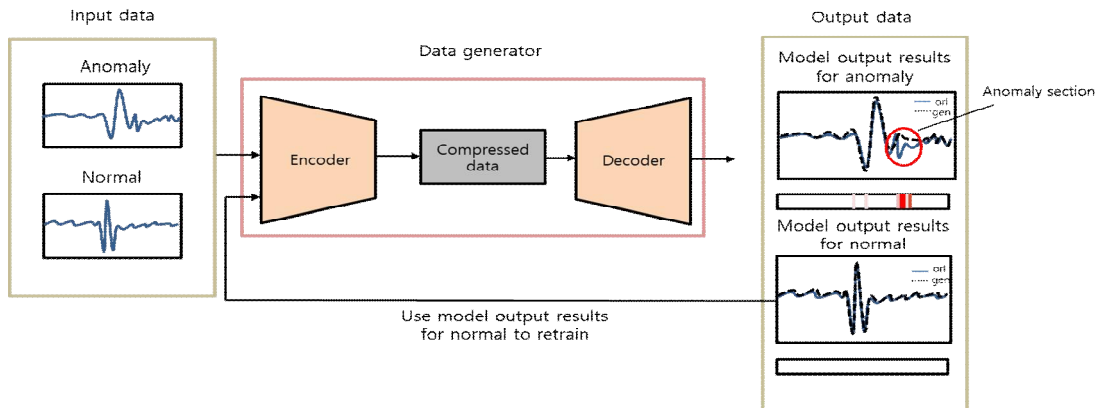
Fig. 3은 제안하는 모델의 전체 구성도로, 비지도 학습(Unsupervised) 방법을 적용한 BeatGAN[8] 모델을 기반으로 하여 국립 환경 과학원에서 제공한 데이터에 대한 성능 향상을 위해 모델의 구조 및 알고리즘의 흐름을 다음과 같이 변경하였다.

첫째, 국립 환경 과학원 전문가들이 이상 탐지 작업을 진행할 때 대략 2개월(1,440시간)의 변화 양상을 분석한다는 것에 의거하여 입력 데이터를 2개월 간격으로 받을 수 있도록 인코더 영역의 합성곱 필터 연산층을 기존 6개 층에서 8개 층으로 확장하였으며 디코더 영역의 전치 합성곱 연산층도 6개층에서 8개 층으로 확장하였다.[9]

구체적인 변경 사항은 Table 4에 커널 사이즈와 채널의 순서로 명시해 두었다.

**Table 4. Changes and additions to the convolution filter structure**

Area	Before change	After change
Encoder	conv4-32	conv4-32
	conv4-64	conv4-64
	conv4-128	conv4-128
	conv4-256	conv4-256
	conv4-512	conv4-512
	conv10-50	conv4-1024
	-	conv4-2048
	-	conv11-200



**Fig. 3. The architecture of the proposed model based on BeatGAN**

Area	Before change	After change
Decoder	convT10-512	convT11-2048
	convT4-256	convT4-1024
	convT4-128	convT4-512
	convT4-64	convT4-256
	convT4-32	convT4-128
	convT4-1	convT4-64
	-	convT4-32
	-	convT4-1

둘째, 정상데이터로의 학습이 끝난 모델이 새롭게 생성해 낸 데이터는 모델이 정상으로 판별할 고유 정보를 담고 있을 것이라고 판단하여 생성해낸 데이터를 학습에 반영할 수 있는 재학습 알고리즘을 구현하여 적용하였다.

### 3.4 평가 방법

제안하는 딥러닝 모델 알고리즘은 정상 데이터와 비정상 데이터를 분류하는 작업을 하므로 Fig. 4의 오차 행렬(Confusion Matrix) 표와 정밀도(Precision) 및 재현도(Recall)를 활용해 성능 평가를 진행하였다.[10]

Confusion Matrix		Predicted	
		Positive	Negative
True Condition	Positive	True Positive(TP)	False Negative(FN)
	Negative	False Positive(FP)	True Negative(TN)

**Fig. 4. Confusion matrix**

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$Recall = \frac{TP}{TP+FN} \quad (3)$$

수식 (2) 와 수식 (3)에서 확인할 수 있듯 정밀도란 모델이 참(Positive)으로 분류한 것 중 실제값이 참인 비율을 의미하며 재현도란 실제 값이 참인 것들 중에서 모델이 참으로 분류한 것을 의미한다.

이상 탐지에 대한 분류 성능을 정확히 측정하기 위해 정밀도와 재현도의 조화평균 값을 F1-Score로 수식 (4)를 통해 계산하였다

$$F1-Score = 2 \times \frac{Recall \times Precision}{Recall + Precision} \quad (4)$$

들어온 입력데이터에 대해 정상 데이터와 비정상 데이터로 판별하기 위해서 수식 (5)와 같이 실제 데이터와 모델이 예측한 데이터 값을 비교하여 일정 값 이상 차이가 날 경우 비정상(Anomaly)으로 판정하였으며 그렇지 않은 경우 정상으로 판정하였다.

$$Anomaly Score = \frac{X_{pred}^t - X_{true}^t}{Size\ of\ pred} \quad (5)$$

### 3.5 학습

국립 환경 과학원에서 제공받은 이상치 판정 데이터만 활용하여 학습을 진행한 것과, 정상 데이터로 학습이 끝난 모델이 생성해 낸 데이터로 새로 학습을 진행하는 것으로 Table 5와 같이 학습데이터, 검증데이터, 테스트 데이터를 나누어서 진행하였다. 모델이 생성한 데이터는 전문가를 통해 검증된 데이터가 아니기 때문

에 학습 데이터로만 사용하였으며 성능 평가의 공정성을 위해 검증 및 테스트 데이터로는 사용하지 않았다.

학습을 위해 사용한 주요 하이퍼 파라미터 정보는 Table 6과 같으며 배치사이즈와 입력채널은 기존의 모델 구조와 동일하게 유지하였고, 데이터 길이는 2개월에 해당하는 1,440시간, 잠재 벡터는 200으로 변경하였다. 잠재 벡터의 경우 입력 데이터의 길이가 길어져 학습해야 하는 특징 공간 (Feature space)이 늘어나기 때문에 효율적인 학습을 위해 기존보다 4배 늘려서 진행하였다. 사용한 손실함수는 판별자(Discriminator)의 경우 이진 크로스 엔트로피(Binary cross entropy)를 사용했으며, 생성자(Generator)의 경우 평균 제곱 오차(Mean Squared Error)를 사용하였고 최적화 함수로는 Adam을 적용하였다.

Table 6. Applied hyper parameters

Hyper parameter	Before	After
Batch size	64	
Input channel	1	
Epoch	100	
Data length	320	1440
Latent vector	50	200

### 3.6 실험 결과

실험은 기존 BeatGAN 모델과 제안하는 모델, 비지도 학습 기반 모델인 Iforest, One Class SVM의 F1-score, 정밀도 및 재현도를 비교하는 것으로 진행하였다[11,12].

각 모델별로 에폭, 입력 데이터별 길이, 재생성 데이터 사용 여부를 달리해서 8개의 원소별 F1-score, 정밀도 및 재현도를 산출하여 Table 7에 명시하였다.

Table 5. Information of train data

Length	Type	Original Data								Original Data + Generated Data							
	Measurement Item	SO2	CO	O3	NO	NO2	NOX	PM10	PM25	SO2	CO	O3	NO	NO2	NOX	PM10	PM25
320	Train set	9360	9465	9608	9272	9254	9712	8472	8106	15811	15916	16060	15724	15705	14924	14557	14557
	Validation set	1020	1054	1663	1155	1172	1726	857	1018	1020	1054	1663	1155	1172	1726	857	1018
	Test set	1329	1328	1329	1451	1352	1341	1365	1270	1329	1328	1329	1451	1352	1341	1365	1270
1440	Train set	1504	1564	1610	1468	1450	1597	1172	1610	2682	2793	1450	2595	1172	2826	2094	1904
	Validation set	306	302	209	384	402	258	168	209	306	302	209	384	402	258	168	209
	Test set	292	290	289	296	296	293	304	289	292	290	289	296	296	293	304	289

Table 7. Experiment result

Model	Length of each input data	Epoch	Whether to apply the re-training algorithm	F1 - SCORE for each measurement item							
				SO2	CO	O3	NO	NO2	NOX	PM10	PM25
BeatGAN original model	320	100	X	0.69	0.67	0.67	0.66	0.65	0.63	0.69	0.69
			O	0.71	0.68	0.75	0.68	0.71	0.69	0.75	0.71
BeatGAN proposed model	1440	100	X	0.69	0.67	0.73	0.66	0.65	0.63	0.71	0.68
			O	<b>0.83</b>	<b>0.88</b>	<b>0.79</b>	<b>0.82</b>	<b>0.89</b>	<b>0.82</b>	<b>0.85</b>	<b>0.86</b>
Isolation Forest	320	X	X	0.64	0.63	0.63	0.65	0.65	0.63	0.69	0.68
	1440			0.76	0.74	0.74	0.77	0.78	0.74	0.83	0.83
One Class SVM	320	X	X	0.64	0.65	0.62	0.65	0.65	0.62	0.70	0.68
	1440			0.75	0.74	0.75	0.77	0.78	0.74	0.83	0.83
				Precision for each measurement item							
Model	Length of each input data	Epoch	Whether to apply the re-training algorithm	SO2	CO	O3	NO	NO2	NOX	PM10	PM25
BeatGAN original model	320	100	X	0.54	0.51	0.57	0.50	0.48	0.47	0.53	0.53
			O	0.59	0.52	0.68	0.55	0.61	0.65	0.68	0.64
BeatGAN proposed model	1440	100	X	0.54	0.50	0.57	0.50	0.49	0.46	0.56	0.54
			O	<b>0.71</b>	<b>0.68</b>	<b>0.61</b>	<b>0.62</b>	<b>0.63</b>	<b>0.64</b>	<b>0.71</b>	<b>0.71</b>
Isolation Forest	320	X	X	0.46	0.46	0.46	0.49	0.48	0.46	0.53	0.52
	1440			0.62	0.59	0.56	0.63	0.62	0.59	0.72	0.75
One Class SVM	320	X	X	0.49	0.48	0.46	0.49	0.49	0.46	0.53	0.52
	1440			0.61	0.58	0.61	0.62	0.62	0.59	0.70	0.70
				Recall for each measurement item							
Model	Length of each input data	Epoch	Whether to apply the re-training algorithm	SO2	CO	O3	NO	NO2	NOX	PM10	PM25
BeatGAN original model	320	100	X	0.97	0.98	0.79	0.98	<b>1.00</b>	0.97	<b>1.00</b>	0.96
			O	0.89	0.98	0.84	0.90	0.84	0.72	0.82	0.83
BeatGAN proposed model	1440	100	X	0.96	0.98	0.79	0.97	0.99	<b>1.00</b>	0.97	0.94
			O	0.89	0.90	0.97	<b>1.00</b>	0.99	0.92	<b>1.00</b>	<b>1.00</b>
Isolation Forest	320	X	X	0.83	0.99	0.99	0.99	0.98	0.98	0.98	0.98
	1440			<b>1.00</b>	<b>1.00</b>	0.86	0.99	0.91	0.95	0.93	0.85
One Class SVM	320	X	X	0.96	0.98	0.95	0.99	0.99	0.99	0.98	<b>1.00</b>
	1440			0.94	0.89	0.98	0.91	0.89	0.99	0.89	0.54

본 논문에서 제안한 입력 데이터 길이 확장, 구조 변형과 재학습 알고리즘을 적용한 모델의 분류 성능이 기존 BeatGAN 모델뿐 아니라 Iforest, One Class SVM 모델들과 비교하였을 때 모든 측정 항목에서 가장 높은 F1-score와 정밀도를 달성하였음을 확인할 수 있었다.

재현도에서는 다른 모델들과 점수가 비슷하거나 떨어지는 경향을 보였으나 정밀도와 재현도가 트레이드 오프(Trade-Off) 관계에 있으며, 이 둘을 동시에 고려한

F1-score가 모든 항목에서 가장 높았다는 점에서 본 실험결과는 유의미한 결과를 산출한 것으로 생각된다.

또한 주목할 만한 점으로 BeatGAN 모델을 활용한 실험의 경우 재학습 알고리즘을 적용하고 100 에폭으로 학습한 모델과 재학습 알고리즘을 적용하지 않고 200 에폭으로만 학습한 모델을 비교해 봤을 때 200 에폭으로 학습한 모델이 과적합의 영향으로 몇 가지 항목에 대한 분류 성능이 감소하는 경우가 발생하였으며,

100 에폭과 재학습 알고리즘을 적용한 모델의 성능이 오히려 더 높았다는 점이다. 이는 모델의 성능을 높이기 위해 에폭을 늘렸을 때 발생할 수 있는 과적합 문제를 새로운 데이터를 만들 수 있는 생성 모델(Generative model)의 특성을 활용해 해결했다는 점에서 의미가 있다고 생각된다.

Fig. 5는 제안하는 모델의 학습 이후 이상진단 결과의 일부를 출력한 그림이다. Fig. 5의 주황색 영역은 제안 모델이 이상 데이터로 판정한 데이터의 이미지이며 빨간색 박스로 실제 이상치 구간을 표시하였다. 실제 이상치 구간에 대해 실제 데이터와 제안 모델이 생성해낸 데이터의 차이를 히트맵(Heatmap) 형태로 표시하였으며 실제 이상구간과 제안 모델이 생성해낸 데이터의 동일 구간 데이터 패턴이 크게 달라서 성공적으로 이상데이터로 검출하는 모습을 확인할 수 있었다. 초록색 영역은 제안 모델이 정상데이터로 판정한 데이터의 이미지이며 일부 구간에 대해 실제 데이터와 다른 패턴을 생성하는 모습이 나타나지만 미리 정해둔 기준치를 초과하지는 않아서 정상 데이터로 분류하는 모습을 확인할 수 있었다.

#### 4. 결론

본 논문에서는 국립 환경 과학원 전문가들이 진행하고 있는 대기자료 이상 탐지 작업을 비지도 학습 기반

의 딥러닝 모델을 활용하여 자동화 할 수 있는 방안을 제안하고자 하였다.

측정된 센서 데이터를 통해 정상 데이터와 비정상 데이터를 구분 짓기 위해서는 지역적 특성, 과거의 데이터 추세 등 다양한 요소를 종합적으로 고려해야 하기 때문에 단순한 알고리즘으로는 원활한 분류 작업이 어렵게 된다. 또한 여러 요소를 동시에 고려해야 한다는 것은 데이터 레이블링의 난이도가 올라갈 수 있다는 것을 의미하며 그럴수록 전문가들의 의견이 상이해질 가능성이 커지게 된다.

같은 현상의 비정상 데이터에 대해 상이한 레이블을 학습함으로 인해 모델의 성능이 하락되는 것을 방지하고자 비지도 학습 기반 방식을 통해 정상 데이터만을 학습 시키고자 하였다.

학습 데이터 생성을 위해 국립 환경 과학원에서 제공받은 기상 데이터의 결측치를 분석하여 부족한 데이터를 최대한 활용하기 위해 결측값을 중심으로 일정 구간에 대한 평균값을 산출하여 결측값을 대체하였다. 또한 다양한 원인의 이상치는 비정상 데이터로 두어 학습 과정에 포함하지 않고 검증 및 테스트 데이터로만 활용하였다.

본 논문에서는 BeatGAN 모델에 기반하여 이상 자료 탐지의 자동화 작업 및 탐지 성능을 높이기 위해 다음과 같은 방법 및 모델 구조를 제안하였다.

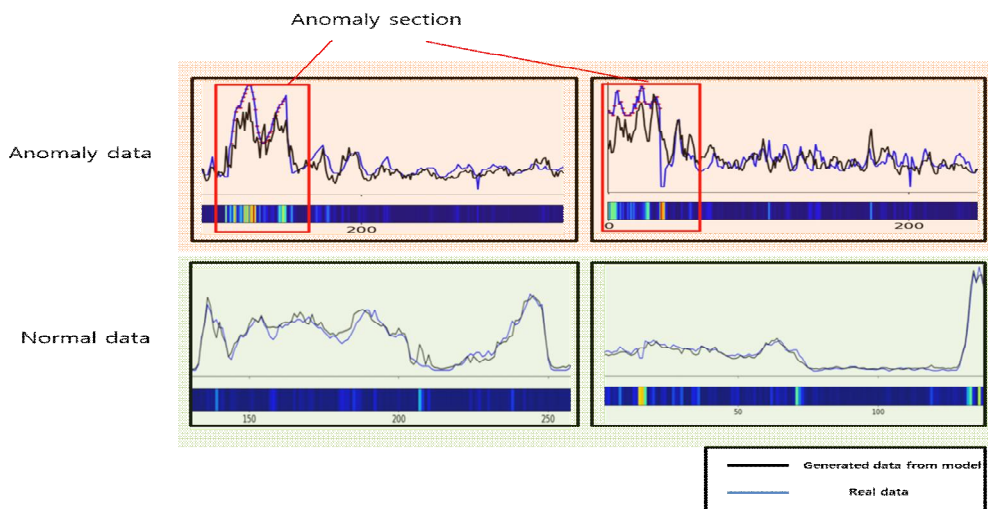


Fig. 5. An example of actual application of anomaly detection through proposed model

첫째, 전문가들이 기상 데이터의 이상 자료를 탐지하기 위해 활용하는 구간 길이가 대략 2개월 단위라는 점에 착안하여 데이터의 입력 구간을 320시간에서 1,440시간으로 변경하였으며 입력 데이터의 길이가 커진만큼 모델의 복잡성을 올리기 위해 BeatGAN 모델 합성곱 필터계층의 커널 구조를 변경하였다.

둘째, 제안하는 모델이 생성 모델(Generative Model)적 특징을 갖고 있다는 점을 활용하여 학습이 끝난 모델이 새로운 데이터를 생성하게 하고 이를 정상 데이터로 간주하여 학습에 다시 활용하는 재학습 알고리즘을 구현 및 적용하였다.

위와 같은 방법들을 적용하여 8개의 기상 자료에 대해 다른 비교 모델 대비 가장 우수한 이상 탐지 성능을 가짐을 F1-score 점수 산출을 통해 확인하였으며, 특히 100 에폭을 학습하고 재학습 알고리즘을 적용한 모델이 재학습 알고리즘의 적용없이 200 에폭을 학습한 모델에 비해 더 높은 성능을 가짐을 확인할 수 있었으며 이는 모델 학습 시 과적합을 피하면서 별다른 추가 비용 없이 모델의 이상 탐지 성능을 올릴 수 있음을 의미한다.

하지만 언급한 방법들의 적용을 통해서도 입력 데이터의 시간 대비 비정상 데이터가 차지하는 시간이 매우 작을 경우 비정상적으로 입력 데이터를 분류 해내는 작업이 원활하지 못한 것을 확인할 수 있었다. 또한 정상 패턴을 학습한 모델이 생성한 데이터를 재학습 시켰을 때 성능이 올라간다는 점은 모델이 다양한 상황을 학습할 수 있도록 하는 실제 데이터의 다양성이 부족한 것으로 보인다.

이와 같은 문제를 해결하는 방법으로 입력 데이터 구간 전체를 정상 데이터와 비정상 데이터로 나누는 것이 아니라 시간별로 정상과 비정상을 분류해 낼 수 있는 방안 및 평가법을 적용하는 것과 전문가들을 통해 선별된 실제 학습 데이터를 더욱 확보하는 것이 성능을 올리기 위해 필요한 부분이라고 판단되며 향후 연구 과정에서는 이를 제안하는 딥러닝 모델의 학습 및 예측 알고리즘에 적용하여 더욱 개선된 결과를 얻고자 한다.

## REFERENCES

- [1] J. Li, H. Izakian, W. Pedrycz & I. Jamal. (2020). Clustering-based anomaly detection in multivariate time series data. *Applied Soft Computing*, 100, 106919. DOI : 10.1016/j.asoc.2020.106919
- [2] A. Deng & B. Hooi. (2021). Graph Neural Network-Based Anomaly Detection in Multivariate Time Series. *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 35, No. 5, pp. 4027-4035). ArXiv:2106.06947 DOI : 10.48550/arXiv.2106.06947
- [3] J. Zhou et al. (2020). Graph neural networks: A review of methods and applications. *AI Open*, 1, 57-81. DOI : 10.1016/j.aiopen.2021.01.001
- [4] H. Ren et al. (2019). Time-Series Anomaly Detection Service at Microsoft. *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, (pp.3009-3017). DOI : 10.1145/3292500.3330680
- [5] L. Ruthotto & E. Haber. (2021). An Introduction to Deep Generative Modeling. *GAMM-Mitteilungen*, 44(2), e202100008. ArXiv:2103.05180 DOI : 10.1002/gamm.202100008
- [6] Air Korea. (2021). *Annual report of the Atmospheric Environment 2020*(Online). [https://www.airkorea.or.kr/web/detailViewDown?pMENU\\_NO=125](https://www.airkorea.or.kr/web/detailViewDown?pMENU_NO=125)
- [7] S. G. K. Patro & K. K. sahu. (2015). Normalization: A Preprocessing Stage. *IARJSET*, 2(3), 20-22. DOI : 10.17148/IARJSET.2015.2305
- [8] B. Zhou, S. Liu, B. Hooi, X. Cheng & J. Ye. (2019). BeatGAN: Anomalous Rhythm Detection using Adversarially Generated Time Series. *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, (pp. 4433-4439). DOI : 10.24963/ijcai.2019/616
- [9] K. Cho et al. (2014). Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. *ArXiv:1406.1078* DOI : 10.48550/arXiv.1406.1078
- [10] J. Davis & M. Goadrich. (2006). The relationship between Precision-Recall and ROC curves. *Proceedings of the 23rd International Conference on Machine Learning - ICML '06*, 233-240. DOI : 10.1145/1143844.1143874
- [11] F. T. Liu, K. M. Ting & Z.-H. Zhou. (2008). Isolation Forest. *2008 Eighth IEEE International Conference on Data Mining*, (pp. 413-422). DOI : 10.1109/ICDM.2008.17
- [12] L. M. Manevitz & M. Yousef. (2001). One-Class SVMs for Document Classification. *Journal of Machine Learning Research*, 2(Dec), 139-154.



양 호 준(Ho-Jun Yang)

[학생회원]



- 2021년 2월 : 인하대학교 컴퓨터공학과 (공학사)
- 2021년 3월 ~ 현재 : 인하대학교 전기컴퓨터공학과 석사과정
- 관심분야 : 머신러닝, 임베디드 인공지능, 데이터 분석
- E-Mail : hjyang@inha.edu

신 유 미(Yu-Mi Shin)

[학생회원]



- 2020년 3월 ~ 현재 : 인하대학교 컴퓨터공학과 학사과정
- 관심분야 : 인공지능, 데이터 분석
- E-Mail : yumishin43@gmail.com

이 선 우(Seon-Woo Lee)

[정회원]



- 2017년 2월 : 인하대학교 컴퓨터정보학과 (공학사)
- 2019년 2월 : 인하대학교 컴퓨터공학과 (공학석사)
- 2019년 3월 ~ 현재 : 인하대학교 전기컴퓨터공학과 박사과정

- 관심분야 : 머신러닝, 딥러닝, 이상진단, 데이터 분석
- E-Mail : x21999@inha.edu

이 석 채(Seok-Chae Lee)

[학생회원]



- 2017년 3월 ~ 현재 : 인하대학교 행정학과 학사과정
- 관심분야 : 머신러닝, 인공지능, 임베디드 시스템, 데이터 분석
- E-Mail : sclee98@inha.edu

이 문 형(Mun-Hyung Lee)

[학생회원]



- 2022년 2월 : 인하대학교 컴퓨터공학과 (공학사)
- 2022년 3월 ~ 현재 : 인하대학교 전기컴퓨터공학과 석사과정
- 관심분야 : 인공지능
- E-Mail : mun0659@inha.edu

권 장 우(Jang-Woo Kwon)

[정회원]



- 1990년 2월 : 인하대학교 전자공학과 (공학사)
- 1992년 2월 : 인하대학교 전자공학과 (공학석사)
- 1996년 8월 : 인하대학교 전자공학과 (공학박사)
- 1998년 2월 : 특허청 사무관

- 2009년 12월 : 동명대학교 컴퓨터공학과 부교수
- 2012년 2월 : 정보통신산업진흥원 인재양성단장
- 2012년 3월 ~ 현재 : 인하대학교 컴퓨터공학과 교수
- 관심분야 : 인공지능, 인간과 컴퓨터 상호작용, 딥러닝
- E-Mail : jwkwon@inha.ac.kr

김 종 구(Jong-Gu Kim)

[학생회원]



- 2022년 2월 : 인하대학교 컴퓨터공학과 (공학사)
- 2022년 3월 ~ 현재 : 인하대학교 전기컴퓨터공학과 석사과정
- 관심분야 : 머신러닝, 딥러닝, 데이터 분석
- E-Mail : kim27y@inha.edu

박 지 훈(Ji-Hoon Park)

[정회원]



- 2006년 2월 : 인하대학교 환경토목공학부 (공학사)
- 2008년 2월 : 인하대학교 환경공학과 (공학석사)
- 2008년 2월 ~ 현재 : 국립환경과학원 대기환경연구과 연구원

- 관심분야 : 모델링, 대기오염측정망
- E-Mail : pjhdo@korea.kr

최 정 무(Jung-Mu Choi)

[학생회원]



- 2015년 2월 ~ 현재 : 인하대학교 컴퓨터공학과 학사과정
- 관심분야 : 인공지능, 열화상
- E-Mail : cjm4788@inha.edu

정 동 희(Dong-Hee Jung)

[정회원]



- 2012년 2월 : 영남대학교 환경공학과 (공학사)
- 2014년 2월 : 영남대학교 환경공학과 (공학석사)
- 2014년 4월 ~ 현재 : 국립환경과학원 대기환경연구과 전문연구원

- 관심분야 : 미세먼지, VOCs
- E-Mail : ehgml6869@korea.kr

신 혜 정(Hye-Jung Shin)

[정회원]



- 1999년 2월 : 이화여자대학교 환경공학과 (공학사)
- 2001년 2월 : 이화여자대학교 대기환경과 (공학석사)
- 2011년 8월 : 이화여자대학교 대기환경과 (공학박사)
- 2003년 ~ 2006년 : 대구지방환경청 환경연구사

- 2006년 ~ 2007년 : 국립환경과학원 환경연구사
- 2007년 3월 ~ 현재 : 국립환경과학원 환경연구관
- 2009년 9월 ~ 2009년 11월 : 캐나다 환경청 환경연구사
- 관심분야 : 인공지능, 미세먼지, 대기오염측정망 구축
- E-Mail : shjoung@korea.kr