

위성 영상과 관측 센서 데이터를 이용한 PM_{10} 농도 데이터의 시공간 해상도 향상 딥러닝 모델 설계

Spatiotemporal Resolution Enhancement of PM_{10} Concentration Data Using Satellite Image and Sensor Data in Deep Learning

백창선¹⁾ · 염재홍²⁾

Baek, Chang-Sun · Yom, Jae-Hong

Abstract

PM_{10} concentration is a spatiotemporal phenomena and capturing data for such continuous phenomena is a difficult task. This study designed a model that enhances spatiotemporal resolution of PM_{10} concentration levels using satellite imagery, atmospheric and meteorological sensor data, and multiple deep learning models. The designed deep learning model was trained using input data whose factors may affect concentration of PM_{10} , such as meteorological conditions and land-use. Using this model, PM_{10} images having 15 minute temporal resolution and $30m \times 30m$ spatial resolution were produced with only atmospheric and meteorological data.

Keywords : PM_{10} , Deep Learning, Satellite Image, Sensor Data, Spatiotemporal Resolution

초 록

PM_{10} 농도는 시간 및 공간 의존성을 동시에 가지는 시공간 데이터이지만 현실적으로 연속적인 시공간 데이터를 획득하는 것은 쉬운 일이 아니다. 본 연구에서는 위성영상과 대기질 및 기상 관측 센서 데이터를 복합적인 딥러닝 모델에 적용하여 시공간 해상도를 향상시키는 모델을 설계하였다. 설계된 딥러닝 모델은 기상, 토지 이용 등 PM_{10} 농도에 영향을 줄 수 있는 인자를 이용하여 학습하였으며, 대기질 및 기상 관측 데이터만을 이용하여 15분 단위의 시간해상도와 $30m \times 30m$ 의 공간해상도를 갖는 PM_{10} 영상을 생성하였다.

핵심어 : PM_{10} , 딥러닝, 위성 영상, 센서 데이터, 시공간 해상도

Received 2019. 11. 21, Revised 2019. 12. 07, Accepted 2019. 12. 26

1) Dept. of Geoinformation Engineering, Sejong University (E-mail: kwsxfk8332@gmail.com)

2) Corresponding Author, Member, Dept. of Environment, Energy & Geoinformatics, Sejong University (E-mail: jhyom@sejong.ac.kr)

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

1. 서론

도시에서 발생하는 PM_{10} (미세먼지) 농도는 시간 및 공간적 특성을 가지고 있으며 시간과 장소에 구애받지 않고 정확한 농도를 파악하는 것은 모니터링 및 의사결정에 있어 중요한 역할을 한다(Le, 2019). 그러나 운영 및 관리, 비용 등과 같은 한계로 인해 시공간적으로 연속적인 데이터를 확보하는 것은 쉬운 일이 아니다. 따라서 현재 PM_{10} 농도 데이터들은 하나의 관측소가 특정 지역을 대표하는 포인트 샘플링 방법을 통해 수집되고 있다. 그러나 포인트 샘플링 방법을 통해 확보한 PM_{10} 농도 데이터는 특정 지역의 관측 센서 분포 밀도에 의존하는 한계가 있으며, 시공간적으로 중요한 패턴이나 특성을 찾아내는데 제한이 있다. 이러한 한계를 극복하기 위해 제한된 포인트 샘플 데이터로부터 샘플링 되지 않은 위치에 대한 특성을 정확하게 유추하여 시공간적 연속성을 예측하는 것이 필수적이라고 할 수 있다.

2. 연구 배경 및 목적

대기질 예측과 관련한 기존의 연구는 크게 결정론적 방법과 통계학적 방법으로 이루어졌다. 결정론적 방법은 대기역학 이론과 물리 화학적 특성에 기반하여 수학적 방법을 통해 PM_{10} 농도의 집중과 확산에 대한 정보를 예측하는 방법으로 많은 양의 데이터를 필요로 하지는 않지만 오염원 정보와 방출량, 화학적 조성과 같은 지식을 필요로 한다(Shahraiyini and Sodoudi, 2016).

반면, 통계학적 방법은 대기 오염 물질과 다른 변수들 사이의 복잡한 관계를 설명하는데 적합하다는 장점이 있다. 딥러닝은 통계학적 방법의 한 예로써 PM_{10} 농도와 이에 영향을 주는 인자들 사이의 비선형적 관계를 비교적 쉽게 찾을 수 있는 방법이다.

많은 연구들이 딥러닝을 이용하여 PM_{10} 농도 예측을 시도하였다. 초기 딥러닝 모델은 주로 퍼셉트론(perceptron)에 기반한 모델인 ANN (Artificial Neural Network)과 MLP 모델을 이용하였다. 입력 데이터로 PM_{10} 농도 관측 데이터와 기상 관측 데이터를 사용하거나(Hooyberghs *et al.*, 2005; Hamza *et al.*, 2016), 혹은 MODIS (Moderate Resolution Imaging Spectroradiometer) 위성 영상의 AOT (Aerosol Optical Thickness) 정보를 사용(Yao *et al.*, 2012)하는 등 다양한 입력 데이터를 이용해 PM_{10} 농도를 예측하려는 시도가 있었다. 딥러닝을 이용한 결과가 다중 회귀법에 의한 결과보다 더 우수하였으며(Yao *et al.*, 2012), 뉴런(neuron)의 개수와 은닉층(hidden layer)의 수가 많을수록 더 좋은 결과를 보였다(Hamza *et al.*, 2016).

PM_{10} 농도 데이터는 시간에 의존적인 시계열적 특성을 지니고 있다. LSTM은 이러한 시계열 데이터 처리에 강점을 지니는 딥러닝 모델로, PM_{10} 농도 예측을 위해 LSTM을 기반으로 한 다양한 모델들이 제시되었다. PM_{10} 농도의 시계열적 특성만 고려하여 다중 LSTM 레이어로 구성된 모델로 농도를 예측하거나(Zhou *et al.*, 2019), 3D-CNN (Convolution Neural Network)과 LSTM을 결합하여 관측소 사이의 공간적 상관성을 추가적으로 고려한 PM_{10} 농도 예측 모델이 제시되었다(Wen *et al.*, 2019).

위에 제시된 모델들은 지상 관측소가 위치한 지점의 PM_{10} 농도를 예측할 수는 있지만, 그렇지 않은 곳의 정보는 예측하지 못한다는 한계가 있다. 최근에는 여러 딥러닝 모델을 결합하여 PM_{10} 농도의 시간적 특성뿐만 아니라, 공간적인 특성을 예측하는 모델들이 제시되었다. 입력 데이터의 특성 분석, PM_{10} 농도 데이터의 보간 및 예측을 통합한 하나의 모델을 제시하거나(Qi *et al.*, 2018), 컨볼루션(convolution) 레이어를 LSTM 유닛(unit)처럼 구성하여 시공간 분석을 가능하도록 한 모델(Le, 2019)은 관측소가 없는 곳의 PM_{10} 농도를 예측할 수 있게 한다.

본 연구에서는 딥러닝을 이용하여 PM_{10} 농도의 분포를 예측할 수 있게 하는 PM_{10} 농도 시공간 해상도 향상 모델을 제시하였으며, 전체적인 흐름은 Fig. 1과 같다. 자료준비 단계(Step 1)에서는 데이터를 수집 및 전처리하며, 시간 해상도 향상 모델링 단계(Step 2)에서는 딥러닝 모델인 MLP (Multi-Layer Perceptron)와 LSTM (Long Short-Term Memory)을 이용하여 대기질 및 기상 관측 데이터의 시간 해상도 향상과 예측을 위한 학습을 수행한다. 마지막으로 시공간 통합 모델링 단계(Step 3)에서는 시공간 해상도가 향상된 PM_{10} 영상을 생성하기 위하여 Step 1, 2에서 산출된 결과물을 이용하여 딥러닝 기반의 생성 모델(generative model) 중 하나인 Conditional VAE (Variational Auto-Encoder) 모델을 학습한다. 학습된 Conditional VAE 모델은 대기질 및 기상 관측 값을 이용하여 기존 연구들보다 향상된 15분 단위의 시간 해상도와 $30m \times 30m$ 의 공간 해상도를 갖는 영상을 생성할 수 있다.

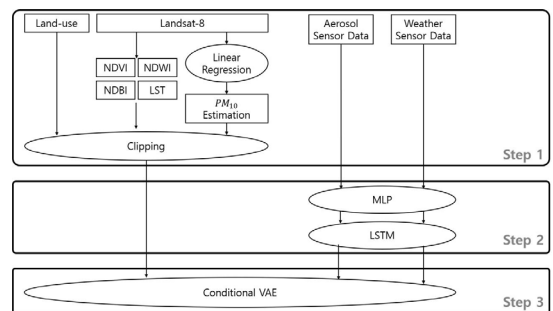


Fig. 1. Work flow

3. PM_{10} 시공간 해상도 향상을 위한 모델 설계

3.1 Step 1: 데이터 확보 및 전처리

본 연구에서는 2015년 1월 1일 0시부터 2018년 12월 31일 23시까지 서울시를 대상으로 하는 대기질 관측 데이터, 기상 관측 데이터, 토지 이용도, Landsat-8 위성 영상 데이터 이용해 모델을 훈련하였다. 대기질 관측 데이터는 에어코리아(Airkorea)에서 제공하는 도시 대기 관측 자료를 이용하였으며(Fig. 2), 25개 관측소(서울시 기준)에서 측정한 1시간 단위의 이산화질소(NO_2), 아황산가스(SO_2), 오존(O_3), 일산화탄소(CO), PM_{10} , $PM_{2.5}$ 농도 관측값을 수집하였다. PM_{10} 농도와 밀접한 관계가 있는 기상 데이터(Hooyerghs *et al.*, 2005)의 경우, PM_{10} 농도 데이터와 동일하게 25개 관측소(Fig. 2)에서 측정한 1시간 단위의 풍향, 풍속, 기온, 습도 자료를 확보하였으며, 환경부에서 제공하는 중분류 토지 이용도와 미지질조사국(USGS: United States Geological Survey)에서 제공하는 Landsat-8 위성영상 데이터를 수집하였다. 중분류 토지이용도는 주거지, 상업지, 내륙수, 혼효림 등 22개의 범례로 구성된 2018년 자료를 사용하였고, 위성 영상의 경우 구름에 가려 정확한 분석이 어

려운 날짜의 영상은 제외하여 25개 날짜에 대한 영상을 확보하였다. 확보한 Landsat-8 위성 영상으로부터 정규 식생 지수(NDVI: Normalized Difference Vegetation Index), 정규 습윤 지수(NDWI: Normalized Difference Wetness Index), 정규 시가지 지수(NDBI: Normalized Difference Built-up Index), 지표면 온도(LST: Land Surface Temperature)를 산출하였다.

Sarawat *et al.*(2017)은 Landsat-8 위성 영상의 가시밴드 대기 반사도 정보와 지상 관측소의 PM_{10} 농도 정보 사이의 상관관계를 이용하여 PM_{10} 농도 데이터의 공간 해상도를 향상시켰다. 본 연구에서는 위 연구에서 제안한 방법을 통해 PM_{10} 농도 예측 레이어를 생성하고, 이를 딥러닝 모델의 입력 데이터로 활용하였다(Fig. 3(a)). 첫 번째로 Landsat-8 위성 영상 1,2,3,4번 밴드의 DN (Digital Number) 값을 태양각을 보정한 TOA (Top of Atmosphere) 반사도로 환산한다. 그 후, 선형 회귀법을 이용하여 각 TOA 반사도의 영상에서 지상 관측소가 위치한 지점의 픽셀값과 PM_{10} 농도 관측값 사이의 회귀 계수를 산출하며, 산출된 회귀 계수를 바탕으로 Eq. (1)을 통해 $30m \times 30m$ 의 공간 해상도를 갖는 PM_{10} 농도 영상을 추정한다.

$$PM_{10} = a_1 R_{\lambda_1} + a_2 R_{\lambda_2} + a_3 R_{\lambda_3} + a_4 R_{\lambda_4} \quad (1)$$

where a_i represents the linear regression coefficient of i th band, and R_{λ_j} means TOA reflectivity of j th band.

컨볼루션 기반의 딥러닝 모델 학습에 용이하도록 PM_{10} 예측 영상, 정규 식생 지수, 정규 습윤 지수, 정규 시가지 지수, 지표면 온도, 토지 이용도를 clipping 하였으며, 여의도 및 종로구, 중구 일대의 데이터를 산출하였다. Clipping 된 영상은 512×512 개의 픽셀들로 구성되어 있으며, 각 픽셀은 $30m$ 의 공간 해상도를 갖는다. 전처리가 끝난 데이터들의 정보는 Table 1과 같으며 clipping 된 결과의 예시는 Fig. 3과 같다.

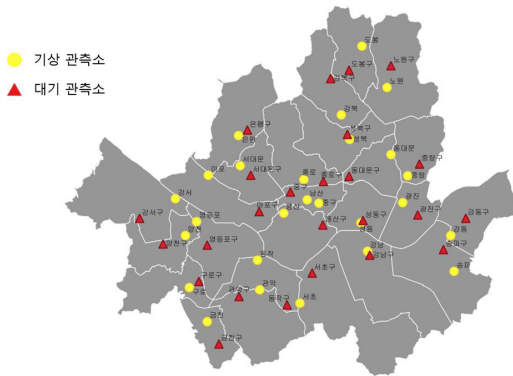


Fig. 2. Location of air quality and meteorological sensors in Seoul

Table 1. Preprocessed data summary

		Attribute	Data quantity	Period
Sensor Data	Air quality	SO_2 , CO_2 , NO_2 , O_3 , PM_{10} , $PM_{2.5}$	$4(year) \times 365(day) \times 24(hour)$	2015-01-01 00:00 ~ 2018-12-31 23:59
	Meteorological	Wind speed, Wind direction, Temperature, Humidity		
Landsat-8	PM_{10} estimation	PM_{10} concentration ($\mu g/m^3$)	25	
	NDVI	Normalized vegetation index		
	NDWI	Normalized wetness index		
	NDBI	Normalized built-up index		
	LST	Land surface temperature ($^{\circ}C$)		
Land-use			1	2018

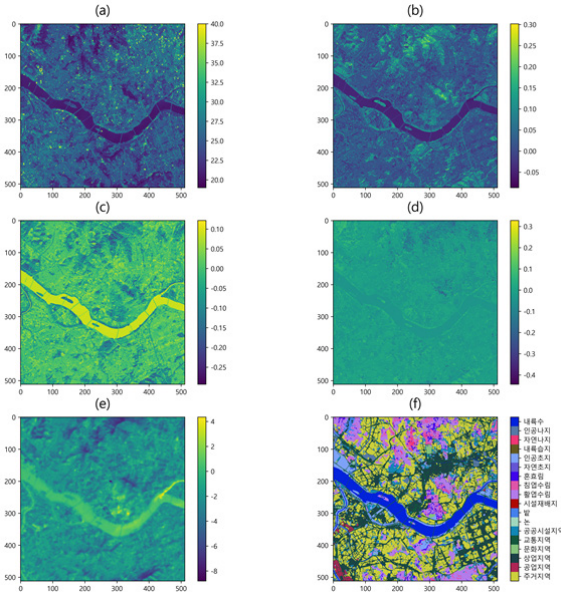


Fig. 3. Example of input data: (a) PM_{10} estimation image, (b) NDVI, (c) NDWI, (d) NDBI, (e) LST, (f) land-use

3.2 Step 2: MLP를 이용한 시간 해상도 향상 모델 및 LSTM을 이용한 예측 모델 설계

서울시 PM_{10} 농도는 관측소에 따라 한 시간 사이에 급변하는 양상을 보이기 때문에 단기적인 모니터링이 필요하며(Itai *et al.*, 2011) 재난에 대한 대비 및 빠른 의사결정을 위해 정확한 예측 또한 중요하다(Zhou *et al.*, 2019). 따라서 현재 한 시간 단위로 제공되는 대기질 관측 데이터의 시간 해상도를 향상시킬 필요성이 있으며, 이를 바탕으로 정확한 PM_{10} 농도를 예측해야 한다.

본 연구에서는 MLP 모델을 이용하여 1시간 단위의 대기질 및 기상 관측 데이터를 15분 단위로 보간(interpolation)하였으며, 해당 데이터를 LSTM 모델에 적용하여 15분 단위로 대기질 및 기상 관측 데이터를 예측할 수 있는 모델을 구축하였다. MLP는 퍼셉트론(perceptron)으로 이루어진 층 여러 개를 순차적으로 붙여놓은 형태의 인공신경망으로 입력층(input layer), 은닉층(hidden layer), 출력층(output layer)로 구성되어 있으며 순전파(forward propagation)와 역전파(back propagation)를 통해 학습을 수행한다. 시간 해상도 향상을 위한 MLP 모델은 Fig. 4와 같다. 3시간 단위 ($t, t+1, t+2$)로 분류한 일련의 대기질 및 기상 관측 데이터 중 t 시간의 데이터는 Fig. 4의 왼쪽 상단에 위치한 MLP 모델에 입력 데이터로 사용하였으며, $t+2$ 시간의 데이터는 Fig. 4 왼쪽 하단에 위치한 MLP 모델의 입력 데이터로 사용하여 해당 시간의 관측값 특성을 학습하였다. 각

MLP 모델은 200개와 300개의 노드를 갖는 2개의 은닉층으로 구성되어 있다. 앞선 2개의 MLP 모델로부터 나온 결과를 합쳐 300, 200개의 노드로 구성된 2개의 은닉층을 갖는 MLP 모델에 전달하면 정답(label) 데이터와 동일한 크기의 결과를 산출하게 되고 이 결과는 MSE (Mean Squared Error) 손실 함수를 통해 $t+1$ 시간의 관측값과 손실 정도를 계산하여 Adam optimizer를 통해 학습을 수행한다. 학습된 MLP는 PM_{10} 농도의 변화와 경향을 고려한 모델로 데이터의 세부적인 변화 특성을 반영한 내삽 결과를 도출할 수 있다. 2015년 1월 1일부터 2017년 12월 31일까지의 데이터를 사용하여 모델을 학습시켰고, 2018년 1월 1일부터 12월 31일까지의 데이터를 이용하여 모델을 검증하였으며, 검증 결과 $2.8604\mu g/m^3$ 의 평균 제곱근 오차(RMSE: Root Mean Square Error) 결과를 보였다. Cubic 보간 방법을 사용하여 시간해상도를 향상시킨 결과와 MLP 보간 결과를 비교한 예시는 Fig. 5과 같다. MLP는 다수의 입력 데이터들이 수많은 매개 변수들을 통해 학습하기 때문에 Cubic 보간법에 비해 실제 데이터의 복잡한 특성을 잘 나타내며 신뢰성이 있는 결과를 보인다.

대기질 및 기상 관측 데이터 예측을 위해 LSTM 모델을 이용하였다. LSTM은 세 개의 개폐장치(input gate, forget gate, output gate)로 구성된 인공신경망의 한 종류로 RNN(Recurrent Neural Networks)의 기울기 소실(vanishing gradient) 문제를 해결하기 위해 고안되었다(Hochreiter and Jürgen, 1997). MLP를 이용하여 시간 해상도가 향상된 15분 단위의 대기질 및 기상 관측 데이터를 96개의 시퀀스(sequence)로 분류한 후, 200개의 유닛, 4개의 은닉층으로 된 LSTM과, 한 개의 완전 연결 레이어(fully-connected layer)로 구성된 모델에 입력 데이터로 활용하였다(Fig. 6). MLP 모델과 동일하게 2015년 1월 1일부터 2017년 12월 31일까지의 데이터를 사용하여 학습을 수행하였고, 2018년 1월 1일부터 12월 31일까지의 데이터를 이용해 검증하였으며, $5.5215\mu g/m^3$ 의 평균 제곱근 오차 결과를 보였다.

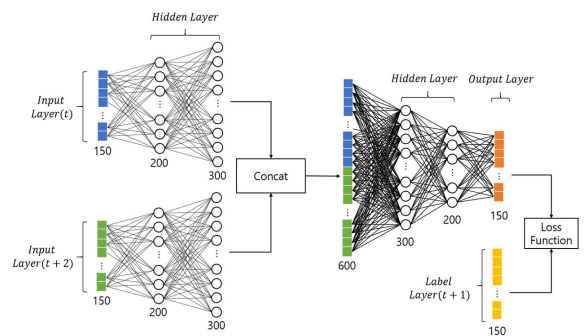


Fig. 4. Interpolation MLP model architecture

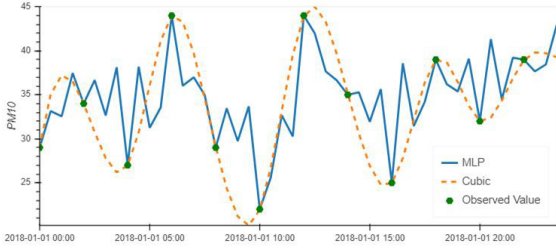


Fig. 5. Example of difference between MLP and cubic interpolation

3.3 Step 3: Conditional VAE를 이용한 15분 단위의 PM_{10} 영상 생성

VAE는 데이터의 확률분포를 학습하는 생성 모델의 한 종류로 인코더(encoder), 디코더(decoder), 잠재 공간(latent space)으로 구성된 자기 감독 학습(self-supervised) 모델이다(Kingma and Welling, 2013). 인코더는 입력 데이터로부터 잠재 변수(latent variable)를 만드는 역할을 하고, 디코더는 인코더가 만든 잠재 변수를 이용하여 다시 입력 데이터로 복원하는 역할을 한다. 잠재 공간은 입력 데이터의 실제 의미를 표현하고 입력 데이터 생성에 관여하는 저차원의 잠재 변수가 군집하는 가상의 공간이며, 여기서 잠재변수는 인코더가 만들어 낸 결과의 평균과 분산을 모수로 하는 정규 분포이다.

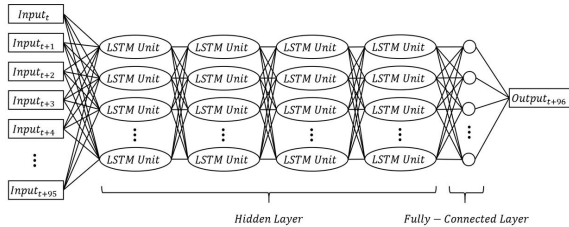


Fig. 6. LSTM architecture for prediction

VAE의 목적 함수는 Eq. (2)와 같다. 우변의 첫 번째 항은 인코더가 입력 데이터 x 를 받아 q 라는 함수를 통해 잠재변수 z 를 만들고 디코더는 인코더가 만든 잠재변수 z 를 입력 데이터로 받아 p 라는 함수를 통해 원 데이터인 x 로 복원하는 것 사이의 크로스 엔트로피(cross-entropy)를 의미한다. 우변의 두 번째 항은 KL (Kullback-Leibler) Divergence를 이용한 인코더와 디코더의 확률분포 차이를 의미하며, VAE는 이 목적 함수 L 이 최소화되는 방향으로 반복 학습을 수행한다.

$$L = -E_{z \sim q(z|x)} [\log p(x|z)] + D_{KL}(q(z|x)||p(z)) \quad (2)$$

Conditional VAE는 인코더와 디코더에 조건을 부여하여 특

정한 결과가 도출되도록 유도하는 VAE 모델로, VAE의 제어가 가능하다는 장점이 있다(Kingma *et al.*, 2014). 본 연구에서는 Conditional VAE 모델에 착안해 잠재 변수에 조건을 부여하는 모델을 설계하였다. 본 연구의 Conditional VAE 모델 아키텍처는 Fig. 7과 같다. Conditional VAE 모델은 위성 영상과 선형 회귀를 통해 추정한 서울시 PM_{10} 농도 영상, 정규 식생 지수, 정규 습윤 지수, 정규 시가지 지수, 지표면 온도, 토지 이용도를 입력 데이터로 사용한다. 입력 데이터들은 6개의 컨볼루션 레이어, ReLU (Rectified Linear Unit) 활성화 함수, MaxPooling 레이어로 구성된 인코더로 전달된다. 컨볼루션 레이어의 경우 여러 번의 실험을 거쳐 가장 좋은 결과를 보인 3×3 kernel size, 1 stride, 1 zero-padding을 하이퍼파라미터(hyperparameter)로 설정하였다. 인코더를 통과한 입력 데이터들은 평균(μ)과 분산(σ) 형태로 표현이 되고, 임의의 가우시안(gaussian) 노이즈와 결합되어 잠재변수를 생성한다. 생성된 잠재변수 z 는 입력 데이터로 사용된 PM_{10} 농도 영상과 동일한 시간의 대기질 및 기상 관측 데이터라는 조건이 합쳐져 조건부 잠재변수(conditioned latent variable)가 된다. 디코더는 6개의 전치된 컨볼루션(transposed convolution) 레이어와 ReLU 활성화 함수로 구성되어 있으며, 전달받은 조건부 잠재변수를 입력 데이터로 복원한다. 디코더를 거쳐 나온 결과는 MSE 손실 함수를 통해 입력 데이터 중 PM_{10} 농도 영상과 손실 정도를 계산하고 이 손실 정도를 바탕으로 Adam optimizer를 이용하여 Conditional VAE 모델을 학습하였다. 2015년 1월 1일부터 2017년 12월 31일 사이의 데이터를 이용하여 모델을 학습하고, 2018년 1월 1일부터 12월 31일 까지의 데이터로 모델을 검증하여 그 결과를 Table 2와 Fig. 8에 나타내었다.

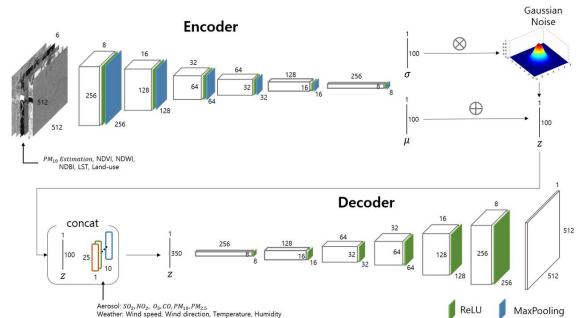


Fig. 7. Conditional VAE model architecture

4. 실험 결과 및 분석

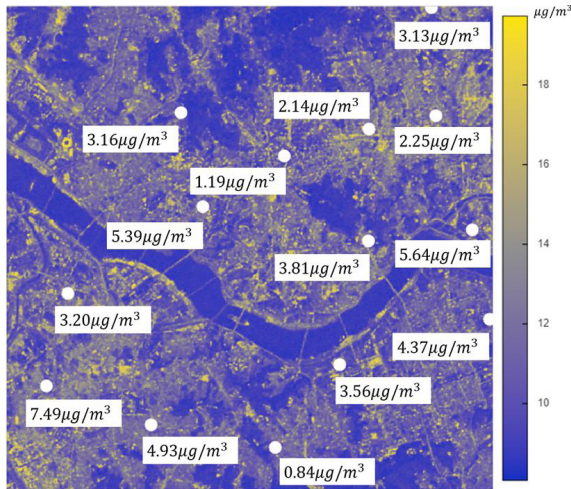
4.1 Conditional VAE 결과 및 분석

특정 시간의 대기질 및 기상 관측 데이터를 임의의 가우시안

Table 2. RMSE result and occurrences by PM_{10} concentration

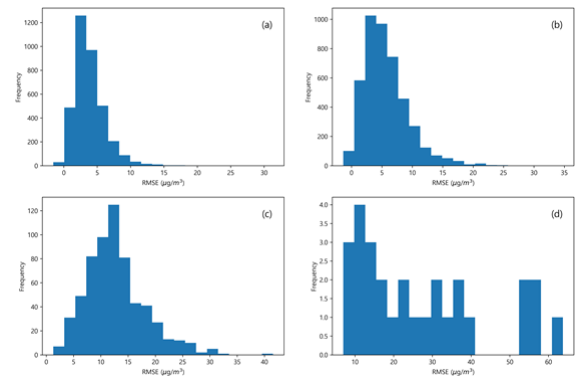
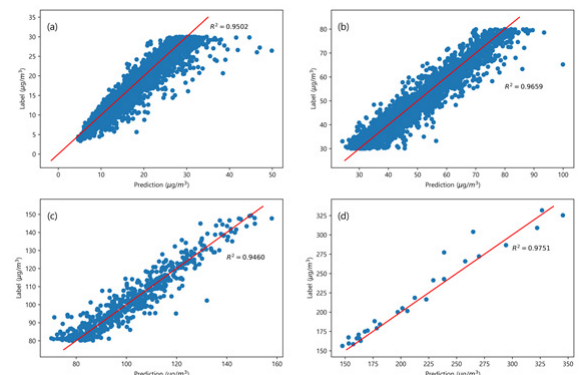
PM_{10} Level	Low ($\sim 30\mu\text{g}/\text{m}^3$)	Medium ($31 \sim 80\mu\text{g}/\text{m}^3$)	High ($81 \sim 150\mu\text{g}/\text{m}^3$)	Very high ($151\mu\text{g}/\text{m}^3 \sim$)
RMSE($\mu\text{g}/\text{m}^3$)	3.6507	5.6034	12.5801	26.7289
Occurrence	3,617	4,463	628	28

노이즈와 결합하여 학습된 디코더에 전달해주면 시공간 해상도가 향상된 PM_{10} 영상을 생성할 수 있다. LSTM 모델에 2018년 1월 1일부터 12월 31일까지의 대기질 및 기상 관측 데이터 적용 결과를 학습된 Conditional VAE 모델에 넣어 시공간 해상도가 향상된 PM_{10} 데이터를 생성하고 이를 검증하였다. 생성된 PM_{10} 영상에서 도시 대기 관측소가 설치된 위치의 픽셀(농도)만 선별하여 실제 지상 관측소에서 관측한 PM_{10} 관측값과 차이를 구한 예시는 Fig. 8과 같으며 이 차이를 바탕으로 평균 제공근 오차를 계산하여 검증하였다. 에어코리아에서 제공하는 PM_{10} 농도 단계에 따라 구분하여 농도별 평균 제공근 오차의 평균값과 발생 횟수를 산출한 결과는 Table 2와 같으며 각 농도별 평균 제공근 오차를 히스토그램으로 나타낸 것은 Fig. 9와 같다.

**Fig. 8. Example of generated PM_{10} image from trained Conditional VAE**

검증 결과, 각 농도 별 상관계수는 모두 강한 양의 상관관계를 보였으며(Fig. 10), 농도별 발생 횟수가 작을수록 평균 제공근 오차가 증가하는 것을 Table 2와 Fig. 9를 통해 확인 할 수 있었다. 미세먼지 농도가 좋음(Fig. 9(a)), 보통(Fig. 9(b)) 일 때는 히스토그램의 평균 제공근 오차 수치가 비교적 낮은 쪽에 집중되어 있는 것을 확인 할 수 있었으나 나쁨(Fig. 9(c)), 매우

나쁨(Fig. 9(d))의 경우 평균 제공근 오차 수치가 큰 쪽에서 빈도수가 증가하였다. 이는 Conditional VAE 모델의 잠재변수가 입력 데이터의 평균과 분산 값으로 이루어진 특성 때문으로, 모델 훈련에 사용한 PM_{10} 데이터의 평균이 $44.2495\mu\text{g}/\text{m}^3$, 분산이 33.4464 였던 것과 입력 데이터 중 고농도 PM_{10} 정보가 상대적으로 부족했던 부분이 복합적으로 작용하여 ‘ 좋음 ($\sim 30\mu\text{g}/\text{m}^3$)’과 ‘보통($31 \sim 80\mu\text{g}/\text{m}^3$)’에서 비교적 좋은 평균 제공근 오차 결과를 보인 것으로 판단된다.

**Fig. 9. RMSE histogram by PM_{10} concentration: (a) Low, (b) Medium, (c) High, (d) Very high****Fig. 10. Scatter plot of prediction and label values: (a) Low, (b) Medium, (c) High, (d) Very high**

4.2 결론

본 연구에서는 위성영상과 지상 관측소 정보를 딥러닝 모

델에 적용하여 15분 단위의 공간 해상도를 갖는 PM_{10} 영상을 생성하였다. 위성영상과 지상 관측값 사이의 회귀 계수를 통해 PM_{10} 관측값의 공간해상도를 향상시켰으며, MLP를 이용하여 관측값의 시간해상도를 향상시켰다. 그 후 LSTM을 이용하여 시간 해상도가 향상된 관측값을 예측하고, Conditional VAE를 이용하여 향상된 시공간해상도 예측 영상을 생성하였다. 제시된 딥러닝 모델을 이용하여 예측한 결과는 $3.6507\mu g/m^3$ (좋은), $5.6034\mu g/m^3$ (보통), $12.5801\mu g/m^3$ (나쁨), $26.7289\mu g/m^3$ (매우 나쁨)의 평균 제공근 오차 결과를 보였다. 제시된 Conditional VAE 딥러닝 모델은 관측값만을 이용하여 영상을 생성할 수 있다는 특징이 있다. 그러나 데이터 불균형 문제로 인해 고농도 PM_{10} 예측에 대해서는 부정확한 결과를 보였다. 언더 샘플링(under-sampling) 및 오버 샘플링(over-sampling)을 통해 학습 데이터를 가공하거나 고농도 PM_{10} 데이터에 가중치를 부과하는 모델을 설계(Krawczyk, 2016)하여 데이터 불균형 문제를 해결한다면 고농도 PM_{10} 에 대한 예측 정확도가 향상될 것으로 판단된다.

감사의 글

이 연구는 한국연구재단 이공분야기초연구사업(NRF-2018R1D1A1B07043821)의 지원으로 수행되었습니다.

References

- Abderrahim, H., Chellali, M.R., and Hamou, A. (2016), Forecasting PM 10 in Algiers: Efficacy of multilayer perceptron networks, *Environmental Science and Pollution Research*, Vol. 23, No. 2, pp. 1634-1641.
- Hochreiter, S. and Jürgen S. (1997), Long short-term memory, *Neural Computation*, Vol. 9, No. 8, pp. 1735-1780.
- Hooyberghs, J., Mensink, C., Dumont, G., Fierens, F., and Brasseur, O. (2005), A neural network forecast for daily average PM10 concentrations in Belgium, *Atmospheric Environment*, Vol. 39, No. 18, pp. 3279-3289.
- Kingma, D.P. and Welling, M. (2013), Auto-encoding variational bayes, *2nd Proceedings of the International Conference on Learning Representations*, ICLR, 2-4 May, Scottsdale, Arizona, USA
- Kingma, D.P., Mohamed, S., Rezende, D.J., and Welling, M. (2014), Semi-supervised learning with deep generative models, *In Advances in Neural Information Processing Systems*, pp. 3581-3589.
- Kloog, I., Koutrakis, P., Coull, B.A., Lee, H.J., and Schwartz, J. (2011), Assessing temporally and spatially resolved PM2.5 exposures for epidemiological studies using satellite aerosol optical depth measurements, *Atmospheric Environment*, Vol. 45, No. 35, pp. 6267-6275.
- Krawczyk, B. (2016), Learning from imbalanced data: open challenges and future directions, *Progress in Artificial Intelligence*, Vol. 5, No. 4, pp. 221-232.
- Le, V.D. (2019), *Spatiotemporal Deep Learning Model for Citywide Air Pollution Interpolation and Prediction*, Master's thesis, Seoul National University, Seoul, Korea, 52p.
- Saleh, S.A.H. and Hasan, G. (2014), Estimation of PM10 concentration using ground measurements and Landsat 8 OLI satellite image, *Journal of Geophysics and Remote Sensing*, Vol. 3 No. 2, pp. 2169-0049.
- Saraswat, I., Mishra, R.K., and Kumar, A. (2017), Estimation of PM10 concentration from Landsat 8 OLI satellite imagery over Delhi, India, *Remote Sensing Applications: Society and Environment*, Vol. 8, pp. 251-257.
- Shahraiyini, H.T. and Sodoudi, S. (2016), Statistical modeling approaches for PM10 prediction in urban areas: A review of 21st-century studies, *Atmosphere*, Vol. 7, No. 2, 15p.
- Wen, C., Liu, S., Yao, X., Peng, L., Li, X., Hu, Y., and Chi, T. (2019), A novel spatiotemporal convolutional long short-term neural network for air pollution prediction, *Science of The Total Environment*, Vol. 654, pp. 1091-1099.
- Yao, L., Lu, N., and Jiang, S. (2012), Artificial neural network (ANN) for multi-source PM2.5 estimation using surface, MODIS, and Fig. rological data, *International Conference on Biomedical Engineering and Biotechnology*, IEEE, 28-30 May, Macau, China, pp. 1228-1231.
- Zhou, Y., Chang, F.J., Chang, L.C., Kao, I.F., and Wang, Y.S. (2019), Explore a deep learning multi-output neural network for regional multi-step-ahead air quality forecasts, *Journal of Cleaner Production*, Vol. 209, pp. 134-145.