



# Development of PM<sub>2.5</sub> short-term forecast model using Artificial Intelligence – Focused on Seoul.

Jeong-Beom Lee<sup>1)</sup>, Geon-Woo Yun<sup>1)</sup>, Youn-Seo Koo<sup>1)</sup>, Hui-Young Yun<sup>1)</sup>, Dae-Ryun Choi<sup>1)</sup>, Ji-Seok Koo<sup>2)</sup>

<sup>1)</sup>Dept. of Environmental & Energy Eng., Anyang University, Anyang, Korea, <sup>2)</sup>Enitech Co.,Ltd

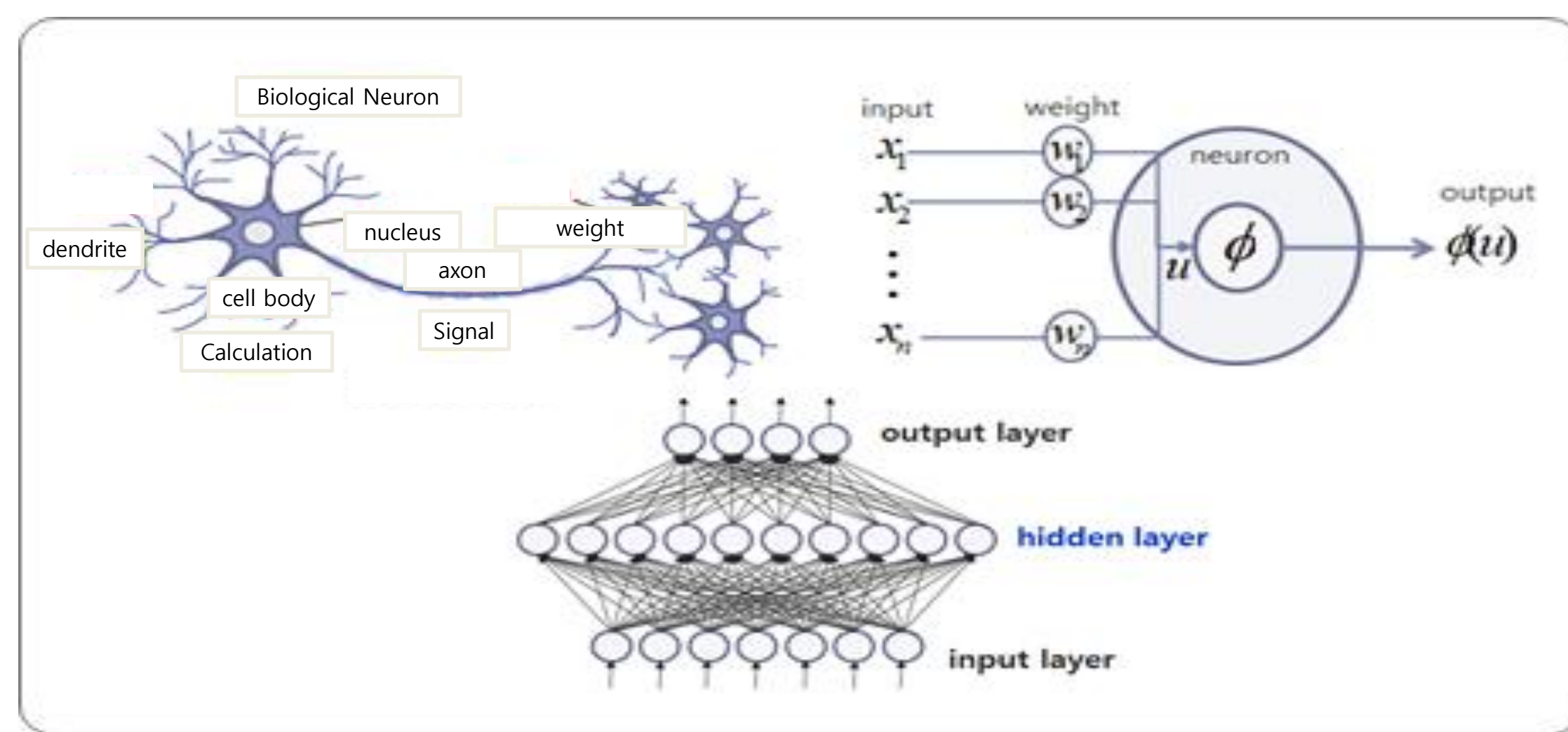
## Introduction

- Forecasting model performance for PM<sub>2.5</sub> using chemical transport model is often overestimated compared to the measurements(Koo et al, 2008;2012;2015, Choi et al, 2018;2019) in Korea.
- In order to improve model performance for PM<sub>2.5</sub> forecasting, we developed PM forecasting system using artificial intelligence with big data such as air quality and weather observations as well as forecasting model data.
- It is important to number of high concentration of PM<sub>2.5</sub> data to accurately predict the episode. However, the number of high concentration of PM<sub>2.5</sub> data is insufficient. Therefore we created the data to improve model performance using AI for accuracy of high concentrations events of PM<sub>2.5</sub>.
- We analyzed developed AI forecasting PM<sub>2.5</sub> model performance for 3-days in Seoul.

## Methodology

### Deep Neural Network

- The artificial intelligence technique used DNN(Deep Neural Network).
- DNN is an extended model that includes multiple hidden layers between the input and output layers to enable deep learning in existing ANN.
- The calculation of weights and biases between layers is key.



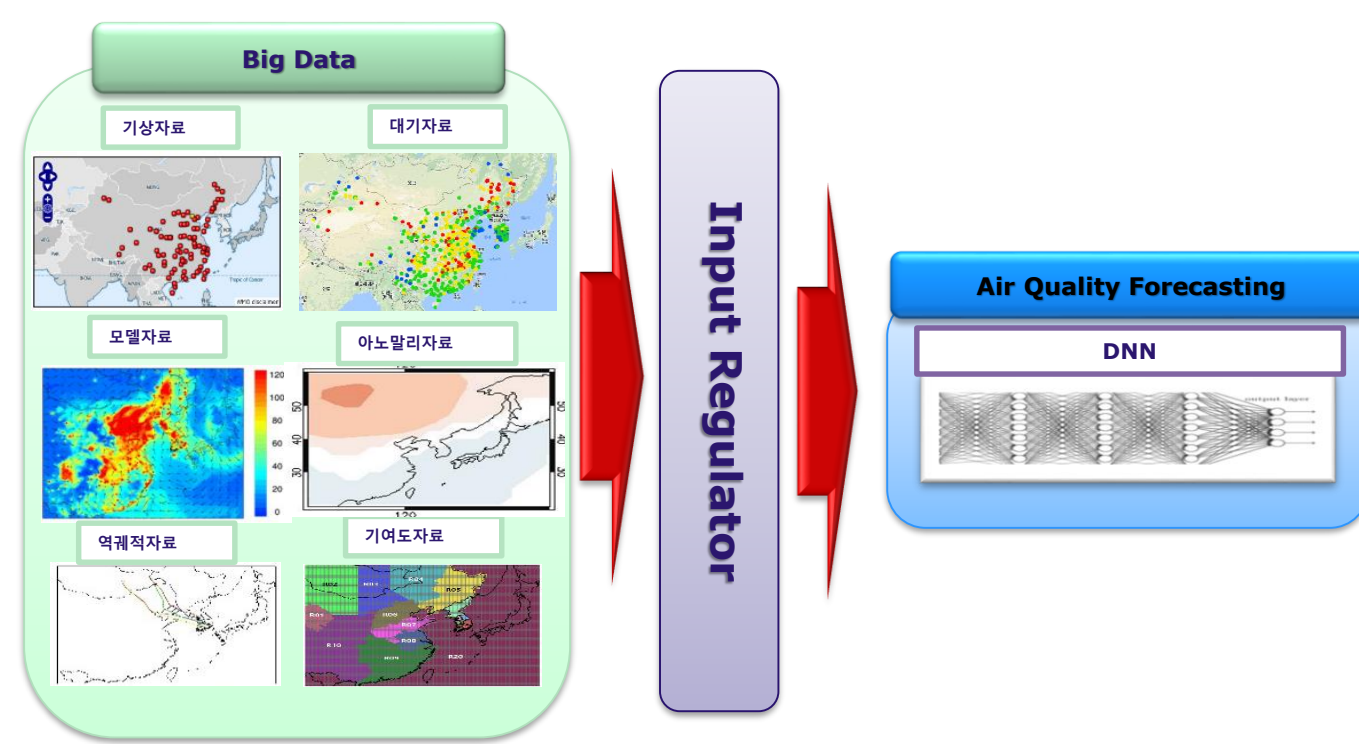
### DNN input data

- Observed data is meteorological and six air pollutants(PM<sub>10</sub>, PM<sub>2.5</sub>, O<sub>3</sub>, NO<sub>2</sub>, SO<sub>2</sub>, CO).
- Numerical model value, WRF weather forecasts, Anomaly, Cosine similarity, Back trajectory, Contribution, Julian day were used.
- The input data used for learning is normalized.

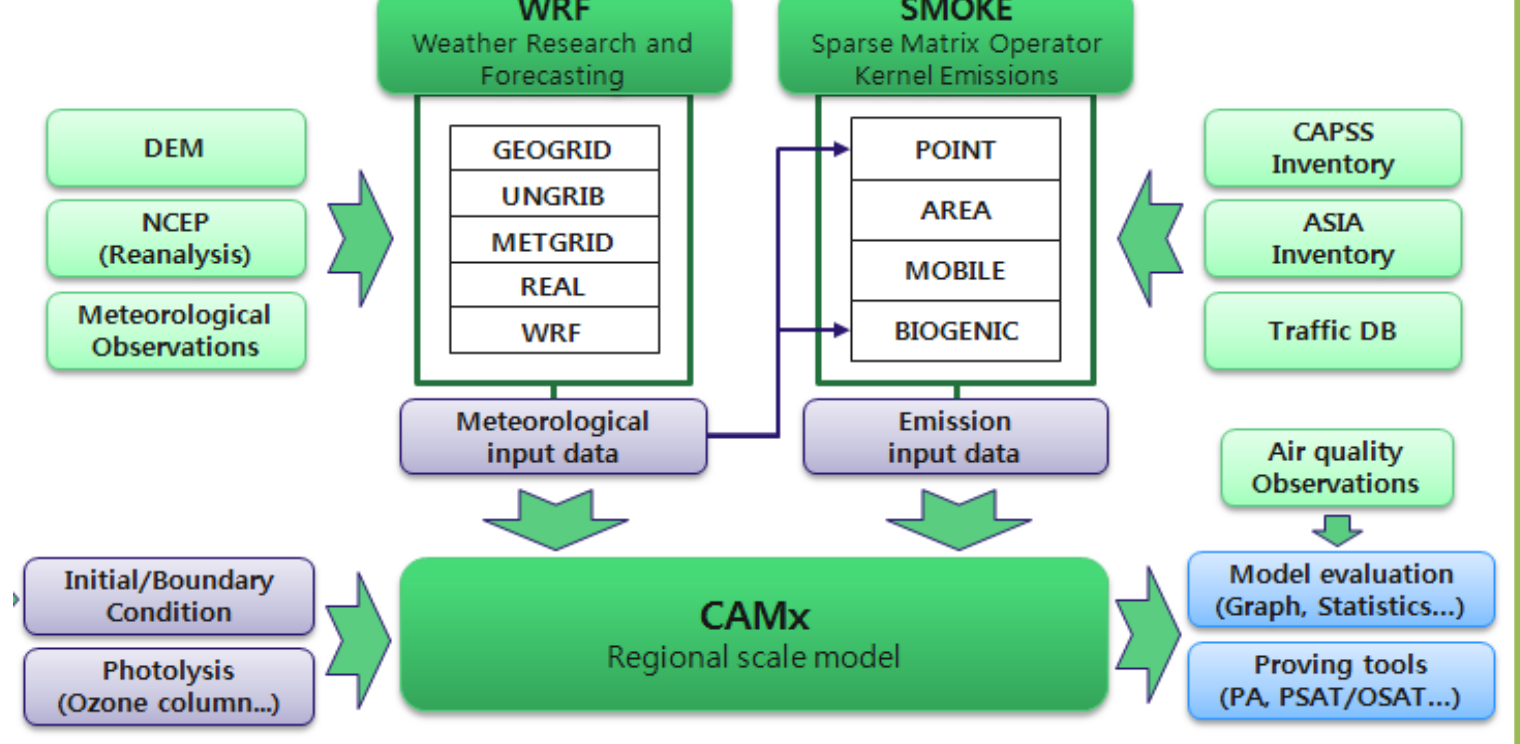
#### Normalization equation

$$\text{If } \min(X_{\text{raw}}) < 0 \\ X_{\text{norm}} = (X_{\text{raw}} - \min(X_{\text{raw}})) / (\max(X_{\text{raw}}) - \min(X_{\text{raw}})) - 0.5$$
$$\text{If } \min(X_{\text{raw}}) \geq 0 \\ X_{\text{norm}} = (X_{\text{raw}} - \min(X_{\text{raw}})) / (\max(X_{\text{raw}}) - \min(X_{\text{raw}}))$$

#### DNN input data flowchart

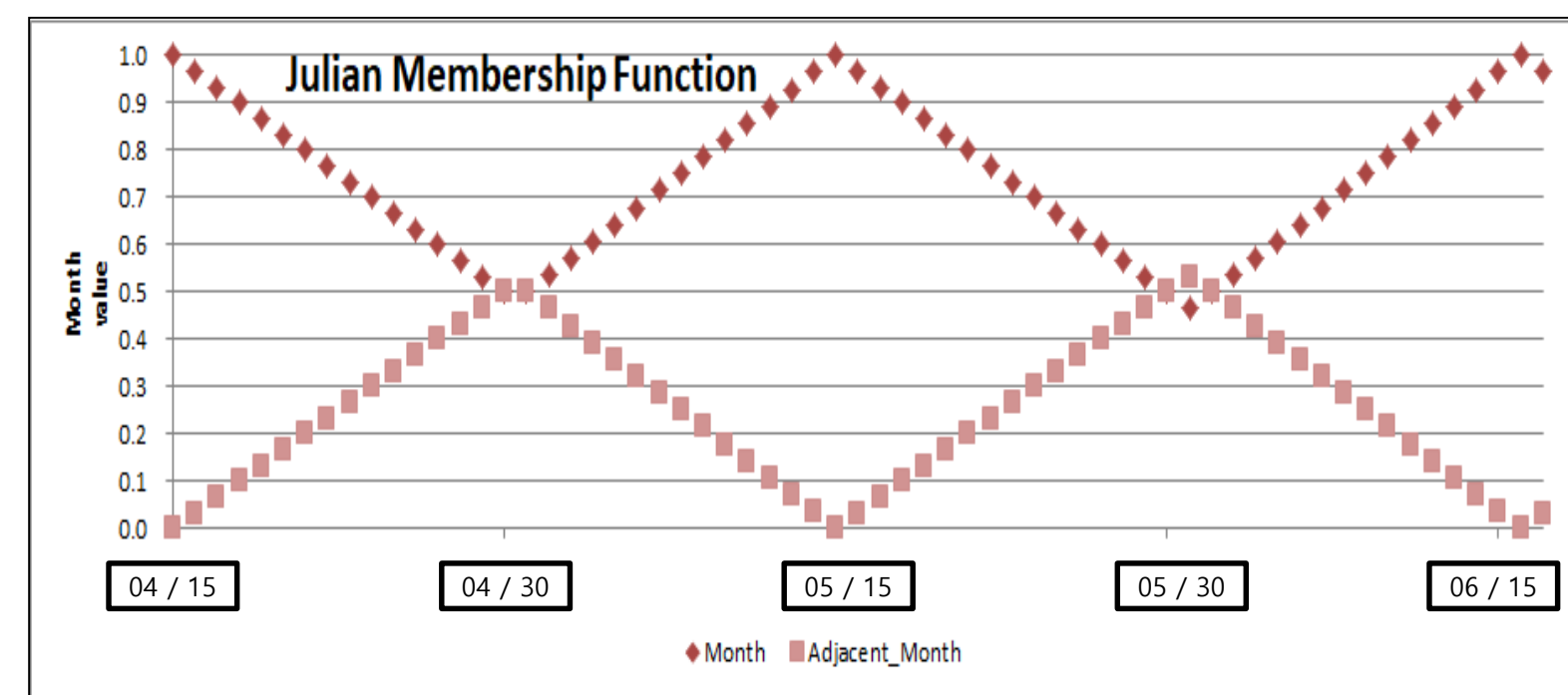


#### Regional modeling system configuration



### Fuzzy theory to Julian day data (Jeon et al., 2018)

1. Adjacent Month definition  
① If(Day = 15) then Adjacent\_Month is not exist  
② else-If(Month = 1 and Day < 15) then Adjacent\_Month = 12  
③ else-If(Month = 12 and Day > 15) then Adjacent\_Month = 1  
④ else-If(Day < 15) then Adjacent\_Month = Month - 1  
⑤ else-If(Day > 15) then Adjacent\_Month = Month + 1  
2. Month value definition  
① If(Day < 15) then Month\_Value =  $\frac{1}{30} \times 240 = 8$   
② else-If(Day > 15) then Month\_Value =  $\frac{1}{30} \times 240 = 9$   
③ Day = 15, Month\_Value = 1  
3. Adjacent\_Month\_value definition  
Adjacent\_Month\_value = 1 - Month\_value



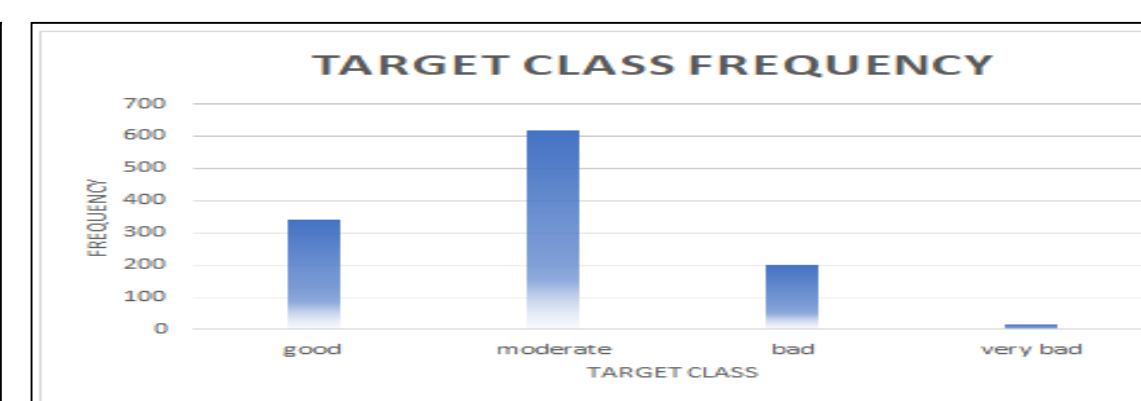
### Data Creation method

- Most observed PM<sub>2.5</sub> data are Good and Moderate.
- Learning is not done properly because the data is not balance.
- Therefore, data creation method is implemented to create sufficient high concentration of PM<sub>2.5</sub> data.

$$X_{\text{new}} = X_{\text{old}} + (X_{\text{old}} * 0.1 * \text{uniform}(x,y))$$

$X_{\text{new}}$  : Created Data.  $X_{\text{old}}$  : Original Data.

$\text{uniform}(x,y)$  : random number from x to y



### Model performance evaluation method

- Change the concentration value of PM<sub>2.5</sub> to the index value.
- Reference : <http://www.airkorea.or.kr/>

#### Index assessment

Column	Forecasted			
	Good	Moderate	Bad	Very Bad
OBS	Good	s1	b1	c1
	Moderate	s2	b2	c2
	Bad	s3	b3	c3
	Very Bad	s4	b4	c4
I : <span style="background-color: #e0f0ff;"> </span> II : <span style="background-color: #fff0e0;"> </span> III : <span style="background-color: #ffe0e0;"> </span> IV : <span style="background-color: #ffcccc;"> </span>				
Method				
Accuracy (A)				
$= \frac{1}{N} \sum_{i=1}^N \frac{ s_i - b_i }{ s_i  +  b_i } \times 100(\%)$				
HIT rate (HIT)				
$= \frac{1}{N} \sum_{i=1}^N \frac{ c_i + d_i }{ c_i  +  d_i } \times 100(\%)$				
Probability of Detection (POD)				
$= \frac{1}{N} \sum_{i=1}^N \frac{ c_i }{ c_i  +  d_i } \times 100(\%)$				
False Alarm Rate (FAR)				
$= \frac{1}{N} \sum_{i=1}^N \frac{ a_i }{ a_i  +  b_i } \times 100(\%)$				

#### Statistic assessment

$$MBIAS = \frac{1}{N} \sum_{i=1}^N (Model - Obs)$$
$$NMB = \frac{\sum_{i=1}^N (Model - Obs)}{\sum_{i=1}^N Obs} \times 100$$
$$IOA = 1 - \frac{\sum_{i=1}^N (|Model - Obs| + |Obs - Obs|)^2}{\sum_{i=1}^N (|Model - Model|^2 + |Obs - Obs|^2)}$$
$$R = \frac{\sum_{i=1}^N (Model - Model) \times (Obs - Obs)}{\sqrt{\sum_{i=1}^N (Model - Model)^2 \times \sum_{i=1}^N (Obs - Obs)^2}}$$

## Results and Discussion

### Results of model performance evaluation (not create data (Standard model))

		Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Index Assessment				Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Statistic Assessment			
Model name	Day	ACC	HIT	POD	FAR	MB	NMB	IOA	R
Standard	D+0	72.22	63.47	73.17	21.05	-0.09	-2.71	0.93	0.91
	D+1	71.43	52.27	70.45	16.22	-2.96	-8.57	0.86	0.90
	D+2	65.87	51.16	69.77	28.57	-2.66	-7.73	0.82	0.85

### Results of model performance evaluation (Create data( $X_{\text{old}}$ = Julian day))

		Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Index Assessment				Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Statistic Assessment			
Model name	Day	ACC	HIT	POD	FAR	MB	NMB	IOA	R
Create Data (Standard value Julian day)	D+0	70.63	65.85	78.05	20.00	0.48	1.43	0.94	0.90
	D+1	64.29	63.64	81.82	33.33	0.81	2.37	0.83	0.85
	D+2	57.14	60.47	79.07	43.33	1.21	3.52	0.79	0.77

### Results of model performance evaluation (Create data( $X_{\text{old}}$ =numericalPM<sub>2.5</sub>))

		Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Index Assessment				Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Statistic Assessment			
Model name	Day	ACC	HIT	POD	FAR	MB	NMB	IOA	R
Create Data (Standard value 'Numerical PM <sub>2.5</sub> ' )	D+0	73.02	68.29	82.93	22.73	0.29	0.89	0.91	0.89
	D+1	74.60	50.00	65.91	6.45	-4.07	-11.82	0.87	0.90
	D+2	66.67	46.51	65.12	24.32	-3.09	-8.99	0.82	0.83

### Results of model performance evaluation (Create data)

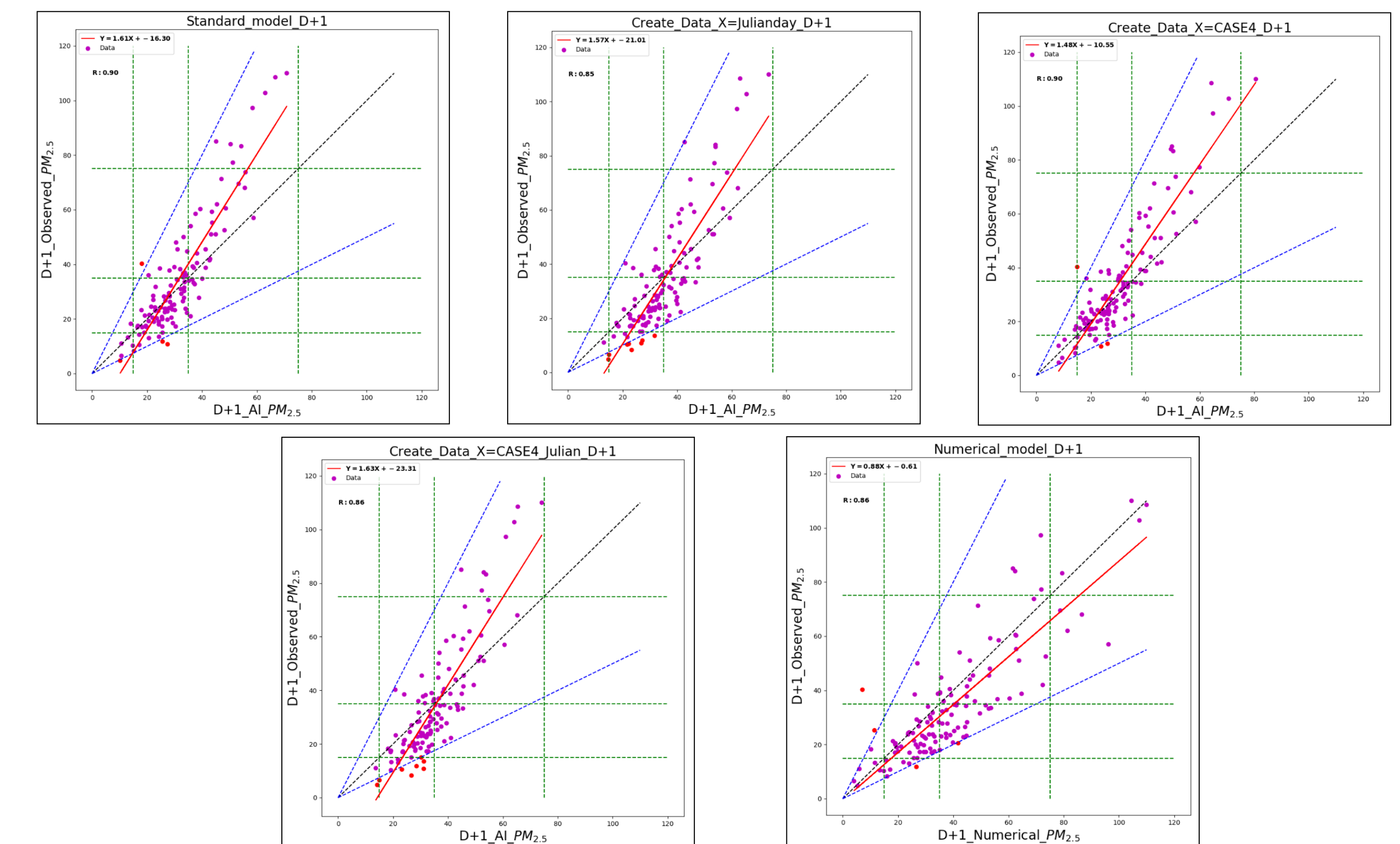
( $X_{1\text{old}}$  = Julian day //  $X_{2\text{old}}$  = Numerical PM<sub>2.5</sub>)

		Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Index Assessment				Observed PM <sub>2.5</sub> vs AI PM <sub>2.5</sub> Statistic Assessment			
Model name	Day	ACC	HIT	POD	FAR	MB	NMB	IOA	R
Create Data (Standard value 'Numerical PM <sub>2.5</sub> &Julian day)	D+0	70.63	75.61	82.93	24.44	2.60	7.84	0.93	0.88
	D+1	65.08	61.36	79.55	31.37	0.89	2.58	0.83	0.86
	D+2	59.52	60.47	79.07	39.29	0.68	1.98	0.81	0.80

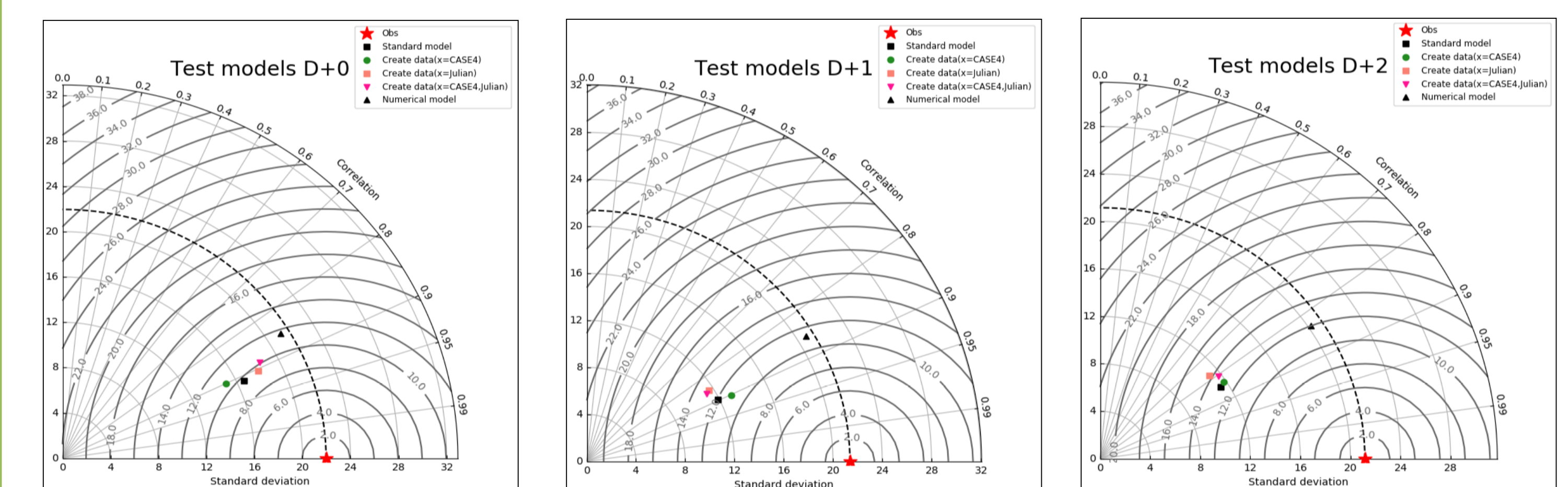
### Results of model performance evaluation compared to numerical model

		Observed PM <sub>2.5</sub> vs Numerical PM <sub>2.5</sub> Index Assessment				Observed PM <sub>2.5</sub> vs Numerical PM <sub>2.5</sub> Statistic Assessment			
Model name	Day	ACC	HIT	POD	FAR	MB	NMB	IOA	R
Numerical Model	D+0	64.29	68.29	85.37	33.96	3.84	11.58	0.92	0.86
	D+1	61.11	68.18	86.36	39.68	5.23	15.18	0.91	0.86
	D+2	60.32	65.12	79.07	43.33	4.60	13.37	0.90	0.83

### Comparison of correlation graphs in test models – Focused D+1.



### Comparison of Taylor-diagram



## Conclusions

- AI results showed ACC increased and FAR decreased compared with numerical mode because AI tend to reduce overestimation of PM<sub>2.5</sub> of a numerical model.
- In D+0, the POD index of AI models with created high concentration of PM<sub>2.5</sub> events data is increased and ACC and FAR are similar compared with standard model.
- In D+1, the POD and FAR index of AI models with created high concentration of PM<sub>2.5</sub> events data using Julian-day or Julian-day&Numerical-PM<sub>2.5</sub> are increased, but ACC is decreased compared with standard model.
- The results of comparisons in various aspects in this study suggest that developed AI forecast model is able to replace numerical model for air quality PM<sub>2.5</sub> forecasting in Seoul.
- We believe further studies with development of data created method are necessary to improve performance of AI model.

## Acknowledgements

This subject is supported by the National Institute of Environmental Research.