

# **Imputation of missing PM<sub>2.5</sub> observations in an air-quality monitoring network of air quality monitoring stations by a new KNN method**

**Formatted:** Font: Italic, Complex Script Font: Italic

## **Supplementary Material**

Idit Belachsen and David M. Broday\*

Faculty of Civil and Environmental Engineering, Technion, Israel Institute of Technology, Haifa,  
Israel

### **Corresponding author:**

David Broday, Faculty of Civil and Environmental Engineering, Technion, Israel  
E-mail: [dbroday@technion.ac.il](mailto:dbroday@technion.ac.il)

**Table S1.** The ID numbers and names of the 59 AQM stations in the imputed dataset. In bold are the 36 AQM stations that have accumulated missing observations of  $\leq 4$  years.

AQM station ID and name			
1: <b>Afula</b>	2: <b>Antokolsky</b>	3: Ehad-Haam	4: <b>Holon</b>
5: <b>Ironi D</b>	6: <b>Kvish 4</b>	7: <b>Petah-Tikva rd.</b>	8: Rakevet Hagana
9: Rakevet Hashalom	10: Rakevet Komemiyut	11: Rakevet Yoseftal	12: <b>Remez</b>
13: Rishon-Lezion	14: Tahana Merkazit	15: <b>Yad Lebanim</b>	16: <b>Yefet Yafo</b>
17: <b>Ahuza G</b>	18: <b>Atzmaut B</b>	19: Begin	20: Igud
21: Kakal	22: <b>Kiryat Ata</b>	23: <b>Kiryat Bialik</b>	24: <b>Kiryat Binyamin</b>
25: <b>Kiryat Tivon</b>	26: <b>Nave-Shaanan</b>	27: <b>Nesher</b>	28: Park Hacarmel
29: <b>Bar-Ilan</b>	30: <b>Efrata</b>	31: Nave-Ilan	32: <b>Ashdod Igud</b>
33: <b>Ashkelon South</b>	34: <b>Dalya</b>	35: <b>Gedera</b>	36: <b>Gvaraam</b>
37: Hayovel	38: <b>Kiryat Malahi</b>	39: <b>Nir-Israel</b>	40: Ofek
41: <b>Ort</b>	42: <b>Rova TV</b>	43: <b>Sderot</b>	44: <b>Sde-Yoav</b>
45: <b>Yahalom</b>	46: <b>Beer-Sheva</b>	47: <b>East Negev</b>	48: <b>Kfar-Masarik</b>
49: Barkai	50: Hefziba	51: Kfar-Saba	52: <b>Pardes-Hana</b>
53: <b>Raanana</b>	54: Shfeya	55: Bnei-Atarot	56: Modein
57: Ashalim	58: Gush-Ezion	59: Neot-Hakikar	

**Table S2.** The performance metrics used in this work, with  $y_i$  the true measured value and  $\hat{y}_i$  the imputed value.

Abbreviation	Term	Formula
NRMSE	Normalized root mean squared error	$\sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} = \frac{RMSE}{mean(Y)}$
$R^2$	Coefficient of determination	$1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$
NMAE	Normalized mean absolute error	$\frac{\sum_{i=1}^N  y_i - \hat{y}_i }{\frac{1}{N} \sum_{i=1}^N y_i} = \frac{MAE}{mean(Y)}$
NMB	Normalized mean bias	$\frac{\sum_{i=1}^N (y_i - \hat{y}_i)}{\frac{1}{N} \sum_{i=1}^N y_i} = \frac{MB}{mean(Y)}$
Normalized CRMSE	Normalized centered root mean squared error	$\sqrt{\frac{1}{N} \sum_{i=1}^N [(y_i - \bar{y}) - (\hat{y}_i - \bar{\hat{y}})]^2} = \frac{CRMSE}{SD(Y)}$

**Table S3.** The tuned hyperparameters of the wkNNr model (see section 2.4.1).  $k$ - the number of weights accounted for,  $q$ - the power in Eq. 4.

Hyperparameter	$k$	$q$	Correlation type
Search space	5-100	1-2.5	Spearman/Pearson
Optimal values:			
very short	10	2.5	Spearman
short	25	2	Spearman
medium-length	30	1.5	Spearman
long	35	1.5	Spearman

**Table S4.** The tuned hyperparameters of the iiET model (see section 2.4.2).  $n\_estimators$ - the number of trees in the ensemble,  $min\_samples\_leaf$ - the minimum number of samples in each leaf node,  $min\_samples\_split$ - the minimum number of samples required to split a node, and  $max\_iter$ - the number of iterations.

Hyperparameter	$n\_estimators$	$min\_samples\_split$	$min\_samples\_leaf$	$max\_iter$
Search space	40-100	2-10	1-10	1-10
Optimal value:				
very short	70	4	1	2
short	80	5	3	2
medium-length	80	4	2	2
long	100	9	2	2