

# **Report: MIC-based Correlation Analysis and LSTM Model for Air Pollution-Related Mortality/Morbidity Assessment**

Author: Yuanming Zhang

Author's Declaration: All Rights Reserved.

# Table of Contents

<b>List of Figures</b>	iv
<b>List of Tables</b>	vi
<b>List of Abbreviations</b>	vii
<b>1 Introduction</b>	1
1.1 Background . . . . .	1
1.2 Literature Review . . . . .	3
1.3 Challenges of the Problem . . . . .	7
1.4 Research Framework . . . . .	9
<b>2 MIC-based Correlation Analysis</b>	11
2.1 Datasets Description . . . . .	11
2.2 Mutual Information Based Association . . . . .	14
2.3 Maximal Information Coefficient for Air Pollutants Selection . . . . .	18
2.4 Experiments and Discussions . . . . .	20
2.4.1 Tools for Experiments . . . . .	20

2.4.2	Numerical Results . . . . .	20
<b>3</b>	<b>LSTM NetWork</b>	<b>27</b>
3.1	Introduction of LSTM Network . . . . .	27
3.2	Data Reorganization . . . . .	30
3.3	LSTM Model for Health Outcome Assessment . . . . .	31
3.3.1	LSTM Network for Feature Information Extraction . . . . .	31
3.3.2	Health Outcome Assessment . . . . .	34
3.4	Results and Discussions . . . . .	36
3.4.1	Performance Metrics . . . . .	37
3.4.2	Parameter Settings . . . . .	38
3.4.3	Numerical Results . . . . .	38
<b>4</b>	<b>Conclusions</b>	<b>52</b>
<b>Bibliography</b>		<b>56</b>
<b>APPENDICES</b>		<b>69</b>

# List of Figures

2.1	Daily mean values of temperature and air pollutants concentration for Toronto, ON, Canada from 2000 to 2021. . . . .	13
2.2	Daily mean values of temperature and air pollutants concentration and daily mortality for Chicago, IL, the U.S. from 1987 to 2000. . . . .	15
2.3	Relations of different air pollutants for Toronto, ON, Canada. . . . .	16
2.4	Comparison of PCCs and MICs with data of all seasons for Toronto, ON, Canada. . . . .	22
2.5	Comparison of PCCs and MICs with data of cold seasons (from October to March) for Toronto, ON, Canada. . . . .	23
2.6	Comparison of PCCs and MICs with data of warm seasons (from April to September) for Toronto, ON, Canada. . . . .	23
3.1	Elementary structure of RNN . . . . .	28
3.2	LSTM network for sequential information extraction . . . . .	32
3.3	Health outcome assessment with weighted evaluation of lags . . . . .	35
3.4	Training losses for different lengths of air pollution exposure sequence ( $m=1-12$ ). . . . .	40
3.5	Comparison between prediction and actual daily mortality ( $m=1-12$ ). .	41
3.6	Training losses with air pollutants $PM_{10} + O_3 + CO$ . . . . .	43
3.7	Comparison between prediction and actual daily mortality with air pollutants $PM_{10} + O_3 + CO$ . . . . .	44
3.8	Training losses with air pollutants $PM_{10} + O_3 + CO + NO_2$ . . . . .	45

3.9 Comparison between prediction and actual daily mortality with air pollutants $\text{PM}_{10} + \text{O}_3 + \text{CO} + \text{NO}_2$ . . . . .	46
3.10 Training losses (left side) and comparison between prediction and actual mortality (right side) for single and multiple air pollutant (s) ( $m=7$ ). . . . .	48
3.11 Comparison between the LSTM model and the GAM for daily mortality assessment. . . . .	51

# List of Tables

2.1	MICs between air pollutants and daily mortality from different causes for Chicago, IL, the U.S. All-cause Mortality is the all-cause non-accidental daily mortality. CVD, COPD, PNEINF, PNEU are RESP are short for cardiovascular deaths, chronic obstructive pulmonary disease, pneumonia and influenza, pneumonia and respiratory deaths. They represent the daily mortality from corresponding causes. Age-1 is for under 65 years of age, Age-2 is for 65 to 74, and Age-3 is for over 75.	26
3.1	Main parameters setting for the LSTM model.	38
3.2	Performance comparison for different $m$ on the training and test sets (TrS: training set, TeS: test set).	39
3.3	Performance comparison with different air pollutants (TrS: training set, TeS: test set).	42
3.4	Comparison of performance metrics with different air pollutant (s) (TrS: training set, TeS: test set, $m=7$ ).	47
3.5	Comparison performance metrics between the LSTM model and the GAM (TrS: training set, TeS: test set).	50

# List of Abbreviations

CO carbon monoxide

NO nitric oxide

NO<sub>2</sub> nitrogen oxide

NO<sub>x</sub> nitrogen oxides

O<sub>3</sub> ground ozone

PM<sub>10</sub> particulate matter less than or equal to 10 micrometers

PM<sub>2.5</sub> particulate matter less than or equal to 2.5 micrometers

SO<sub>2</sub> sulphur dioxide

T daily mean temperature

**ADAM** Adaptive Moment Estimation

**APHEA** Air Pollution and Health: a European Approach

**GAM** Generalized Additive Model

**GLM** Generalized Linear Model

**LSTM** Long Short-Term Memory

**MAE** Mean Absolute Error

**MI** Mutual Information

**MIC** Maximal Information Coefficient

**NAPS** National Air Pollution Surveillance

**NMMAPS** National Morbidity and Mortality Air Pollution Study

**PCC** Pearson Correlation Coefficient

**RMSE** Root Mean Square Error

**RNN** Recurrent Neural Network

# Chapter 1

## Introduction

### 1.1 Background

Short- and long-term exposure to air pollutants has been associated with various adverse public health effects in numerous epidemiology studies. Exposure to major regulated air pollutants like ground-level ozone, particulate matters, oxides of sulfur and oxides of nitrogen has become one of the high environmental risks and has been shown to have a relationship with increasing cardiovascular disease [1, 2, 3, 4], respiratory disease [5, 6, 7], and public mortality and morbidity [8, 9, 10]. In order to monitor and improve air quality, learn about the health impacts of air pollutants, and reduce the burden of social healthcare, a series of funds have been granted in several countries or regions to research on ambient air quality monitoring and statistical relationship between air pollution exposure and specific health consequence.

Among these air pollution-related air quality and public health studies, three influential programs are worth mentioning. The first one is [National Morbidity and](#)

Mortality Air Pollution Study (NMMAPS) [11, 12], which was initiated in 1996 with the main purpose of studying the association between particulate matters and daily mortality, as well as the adverse effects of several regulated air pollutants on public health in United States. This research work was conducted by investigators from Johns Hopkins University and Harvard University, and was funded and supported by the Health Effects Institute. Air pollution data and daily mortality of 108 large cities from 1987 to 2000 in the United States were collected. Time-series analysis of several air pollutants and association between air pollutants and daily mortality were also explored. The second one is the Air Pollution and Health: a European Approach (APHEA) project, including APHEA and APHEA2 [13, 14, 15]. Theses projects began in 1993, which also focused on investigating the adverse effects of short-term exposure to particulate matter and other air pollutants on public health in European cities. Quantitative evaluations of their short-term effects using data from 15 cities in APHEA and 8 in APHEA2 were investigated, where association between exposure to ambient particles and mortality or hospital admissions was identified. The third one is the National Air Pollution Surveillance (NAPS) Program [16, 17] started in 1969, which is managed by Environment and Climate Change Canada and has multiple purposes including monitoring the variation of Air Quality Health Index and Canadian Environmental Sustainability Indicators. Ambient air pollution data of NAPS are collected from more than 250 monitoring stations distributed across Canada and are publicly accessible through the Canada-Wide Air Quality Database. The above-mentioned programs have been promoting the development of quan-

titative and statistical methods for adverse impacts assessment of air pollution on public health, and have achieved significant results [18, 19]. Collected historical data of air pollution and medical records and implementation of innovative statistical tools in these research are still having a far-reaching influence on current air pollution-related epidemiology studies.

## 1.2 Literature Review

One of the major topics in air pollution epidemiology research is to investigate the impacts of short- or long-term exposure to air pollutants of interest on specific public health outcome. Various statistical methods and epidemiological designs have been explored for this problem since mid-twentieth century. These research methods generally fall into two categories: cross-sectional study and longitudinal study, with the former one investigating air pollutants exposure and health outcome at a single point of time and the latter one over a period of time [20, 21].

In earlier studies, as systematic monitoring and data collection were not available, miscellaneous methods under the above-mentioned two categories were explored. Most of the research were conducted for air pollution-caused health incidents or based on observation of specifically designed experiment, where population groups, types of pollutants, exposures, etc., were regulated. In 1961, Ciocco *et al.* [22] investigated the impacts of air pollution on residents in Donora and Pittsburgh using data from community survey ten years after the Donora smog disaster in 1948 and showed higher morbidity and mortality were related to exposure to higher air pol-

lutants concentration. In 1970, multiple regression analysis method was introduced to quantitatively estimate the long-term effects of air pollutants using data from particular locations of England and showed a strong relationship between specific air pollutants and several diseases [23]. Similar methods were also used later by Lipfert [24], Gibbons and McDonald [25], Ostro [26], etc.

Later, with the need to investigate specific health outcome caused by exposure to air pollutants of interest over a period of time, cohort study was used with prospective cohort investigating future health consequences and retrospective cohort for past relationship. In 1973, a prospective cohort study was conducted by Waller *et al.* [27] to investigate the lasting impacts of exposure to London fog in the 1950s, and higher prevalence of respiratory symptoms in the exposed residents was reported. From 1978 to 1981, Kerigan *et al.* [28] conducted a three-year cohort study to investigate the influence of ambient air pollution on children's respiratory health in Hamilton, Ontario. This method was also applied to a series of air pollution epidemiology studies such as cancer [29, 30, 31], diabetes [32, 33], cardiovascular disease [34, 35], mortality [36, 37, 38], etc. Panel study is another longitudinal research design similar to cohort study used in air pollution epidemiology, with the main difference that the former one focuses on the same group of population while the latter on population with shared characteristics. Air pollution panel study was also used in a variety of epidemiological problems, such as asthma [39, 40], cardiac disease [41], pulmonary disease [42, 43], neurobehavioral disorder [44], inflammation [45, 46], etc.

In addition, case-crossover analysis is also an extensively used longitudinal method

in air pollution epidemiology when estimating their acute impacts on same or similar population group through contrast before and after some health consequence [47]. Since the development of case-crossover method in 1991 [48], it has been applied to a variety of experimental designs for air pollutants-related epidemiological studies. In 1999, case-crossover designs were made for assessing the effects of air pollution on mortality in Philadelphia of Pennsylvania [49] and Seoul of Korea [50]. The effects of acute exposure to air pollution on sudden cardiac arrest and asthma attack were investigated by Checkoway *et al.* [51] and Im *et al.* [52] in 2000. In the following decades, case-crossover analysis was used in impacts assessment of air pollution on myocardial infarction [53] (2003), coronary mortality [54] (2005), otitis media [55] (2010), respiratory disease [56] (2014), incident pneumonia [57] (2018), etc.

With development of air pollutants monitoring systems and availability of continuously collected data, longitudinal analysis using time series methods were applied to investigate the cumulative impacts of air pollution exposure. Two of extensively discussed models are [Generalized Linear Model \(GLM\)](#) and its extension [Generalized Additive Model \(GAM\)](#), where the health consequence is modeled using Poisson process as a function of air pollution exposure. In 1979, a logistic regression model was presented by Korn and Whittemore [39] to analyze the impacts of exposure to varying concentrations of air pollution on asthma, with the motivation to overcome the problems of independence requirement and missing response data in linear regression. In the following years of time series study for air pollution epidemiology, Poisson regression was proposed to analyze a series of air pollutants-related epidemiological issues

like morbidity [58], respiratory disease [59, 60], lung cancer [61], mortality [62, 63], etc. When it came to around 2000, **GLM** and its extension **GAM** were extensively applied to air pollution epidemiology by a number of researchers, where quite a few significant work were motivated by the aforementioned three programs. In 2000, Dominici *et al.* [18] presented a general semiparametric log-linear regression model to estimate the PM<sub>10</sub>-associated mortality rate in 20 large US cities and a hierarchical regression model was also developed to investigate the spatial correlation of responding coefficients among different cities. In the same year, Sheppard and Damian [64] proposed a Poisson regression model for health effects assessment of short-term air pollutant exposure, which combined exposure distribution and measurement error with the disease model. This model was applied to particulate matter-associated asthma hospital admission in Seattle and showed that the measurement error did not affect the results of the time series regression model. Similar methods were also used in **APHEA** projects. In 1998, Spix *et al.* [65] used Poisson regression models and standard confounder models to investigate the air pollution-related hospital admissions of different age groups in five West Europe cities and showed that ozone had a significant impact on respiratory diseases. In 2001, researchers from **APHEA2** project [66] applied Poisson regression to particle matter-related mortality analysis in 29 European cities, taking advantage of its non-parametric smoothing of seasonal patterns, and confirmed the impacts of ambient particles on mortality while finding the characteristics-related heterogeneity of effect parameters. Besides these single-pollutant time series models, multi-pollutant models under **GAM** framework were

also investigated. In 2004, Zeka and Schwartz [67] presented a two-pollutant Poisson regression model by linearly correlating two kinds of air pollutants to investigate the effects of particulate matters and several other air pollutants on daily mortality based on the [NMMAPS](#) datasets. They confirmed the effects of PM<sub>10</sub> on daily mortality and reported a carbon monoxide-related daily mortality increase. In order to estimate the cumulative effects of exposure, Zanobetti *et al.* [68] developed [GAMs](#) with distributed lags that related some health outcome on a given day with air pollution exposure in prior days. Timescale effects of exposure was also investigated by Dominici *et al.* [69] by using predictors from Fourier decomposition of air pollution time series in [GAM](#) and later by Burr *et al.* [70] through combining the lagged predictors in [GAM](#). Until today, [GAM](#) is still one of the most popular time series methods used in air pollution epidemiology, with application in multi-site time series [71], bias correction [72], multi-pollutant effects [73, 74], temporal trends [75], etc.

In summary, longitudinal study has become the mainstream method in air pollution epidemiology at present, of which time series methods, especially the [GAMs](#) are widely used. Moreover, complex epidemiological designs with combination of time series and other longitudinal methods have been gradually developed for meticulous investigation of the health effects from cumulative air pollution exposure.

### 1.3 Challenges of the Problem

As the most widely used time series model in air pollution epidemiology, [GAM](#) provides a perspicuous air pollution exposure-health response assessment framework

taking advantage of Poisson regression for various health events, non-parametric splines for exposure effects estimation and confounding elements smoothing, while the following issues and challenges still exist in its epidemiological applications.

1. Different air pollutants are correlated with each other. Air pollutants are usually mixed with each other physically and chemically from their generating sources to spreading, especially for some strongly related ones like different oxides of nitrogen or particulate matters with different diameters. As people are generally exposed to multiple air pollutants simultaneously, finding the primary risk factors for some particular health consequence is usually not easy in the **GAM**-based modelling and analysing process. Capturing and evaluating the mixed health impacts of correlated air pollutants by smoothing the confounding effects may not be enough and keep an open question for further investigation. In addition, collinearity (also known as multicollinearity) brought by the associations reduces the interpretability of response coefficients for air pollutant of interest and weakens the power of **GAMs**.

2. Performance of the **GAMs** is limited. **GAM** presents a straightforward exposure-response relationship by fitting the model using time-varying air pollution concentrations and health events data, which facilitates analysis of the degree to which health outcomes respond to air pollution variations. However, in most air pollution epidemiological applications, to what extent the fitted model can reflect the true relationship between air pollutants of interest and particular health consequence is generally not evaluated. Either underfitting or overfitting may degrade its performance of generalization, and thus reliability of the exposure-response coefficients

obtained from the [GAMs](#) might need further examination.

3. Public health outcomes are affected by cumulative effects of exposure. In general air pollution epidemiology application of [GAMs](#), it is usually assumed that adverse health outcome at one time period is associated with exposure at a previous period. This kind of formulation is referred to as single lag model, while its alternative is distributed lag model, which involves cumulative effects of exposures at multiple time periods before the health event occurs and conforms better to realistic situation. Parametric designs are usually used to model the relationship between previous distributed exposures and lagged responses. Most of these formulations focus on the cumulative effects of a single air pollutant for parameters fitting and coefficients interpretation issues. Further studies are still needed to investigate the effects of exposure to multiple air pollutants in distributed lag models.

## 1.4 Research Framework

Aiming at the above-mentioned issues, the following main research work related to quantitative health risk analysis in air pollution epidemiology is made as follows:

- 1) [Chapter 2: Mutual Information \(MI\)-based Maximal Information Coefficient \(MIC\)](#) is used to evaluate the association between different air pollutants and find the most related factors for health consequence of interest, taking advantage of information entropy to capture both the linear and nonlinear relation.
- 2) [Chapter 3: An Long Short-Term Memory \(LSTM\)](#) model is developed to assess the impacts of exposure to multiple air pollutants on health outcome of inter-

est with weighted evaluation of distributed lags. An [LSTM](#) neural network is first designed to extract health outcome-related feature information from multi-pollutant exposure sequence with temporal dependence. Then estimation layers with weighted evaluation of the extracted features from the exposure that has distributed lags are constructed to assess the health consequence of interest.

# Chapter 2

## MIC-based Correlation Analysis

In this chapter, datasets used for analysis throughout this report are first introduced. Then [MI](#) and [MIC](#) are applied to correlation analysis between different air pollutants and air pollution variables selection.

### 2.1 Datasets Description

The ambient air pollutants data used throughout this work are mainly from the publicly available Canada-Wide Air Quality Database with pollution data from the [NAPS](#) program and the database for [NMMAPS](#). Information of these two databases are briefly introduced as follows.

The [NAPS](#) database includes continuous and hourly measurements of [nitric oxide \(NO\)](#), [nitrogen oxide \(NO<sub>2</sub>\)](#), [nitrogen oxides \(NO<sub>x</sub>\)](#), [carbon monoxide \(CO\)](#), [sulphur dioxide \(SO<sub>2</sub>\)](#), [ground ozone \(O<sub>3</sub>\)](#), [particulate matter less than or equal to 2.5 micrometers \(PM<sub>2.5</sub>\)](#) and [particulate matter less than or equal to 10 micrometers](#)

( $\text{PM}_{10}$ ). These data are available from the Government of Canada Open Data Portal [76] and detailed description of the program and datasets is available at Open Government [77]. As the data availability is reduced for the impact of the COVID-19 pandemic, 22-year air pollution data from 2000 to 2021 are collected. In addition, daily mean temperature ( $T$ ) of the same periods are collected from the Canadian Centre for Climate Services [78]. All categories of data are preprocessed as follows to facilitate the numerical experimentation. The raw air pollution time series data in the NAPS datasets are measured on an hourly basis and their 24-hour means are first calculated since all the analyses throughout this work are conducted on a daily basis. Then a second-order polynomial interpolation is performed for the missing values. At last, the outliers including negative values of the interpolated time series are replaced using the mean of values before and after the interpolated day.

Daily mean temperature and seven air pollutants data from 2000-2021 for Toronto, ON, Canada are shown as an example in Figure 2.1. It can be roughly seen that  $T$  and  $\text{O}_3$  have obvious yearly periodicity. People are prone to be exposed to higher  $\text{O}_3$  concentrations in higher temperature periods, as  $\text{O}_3$  is closely related to solar irradiance in its formation process. Another obvious point is that daily concentrations of  $\text{CO}$ ,  $\text{NO}$ ,  $\text{NO}_2$ ,  $\text{NO}_x$  and  $\text{SO}_2$  all have downward trends in the given periods, possibly due to various stricter air pollution emission standards. It can also be seen that the  $\text{PM}_{2.5}$  series keeps relatively stable over this time span. The  $\text{PM}_{10}$  data are excluded for the reason that the vast majority of them are missing.

The NMMAFS database includes daily measurements of  $\text{NO}_2$ ,  $\text{CO}$ ,  $\text{SO}_2$ ,  $\text{O}_3$ ,

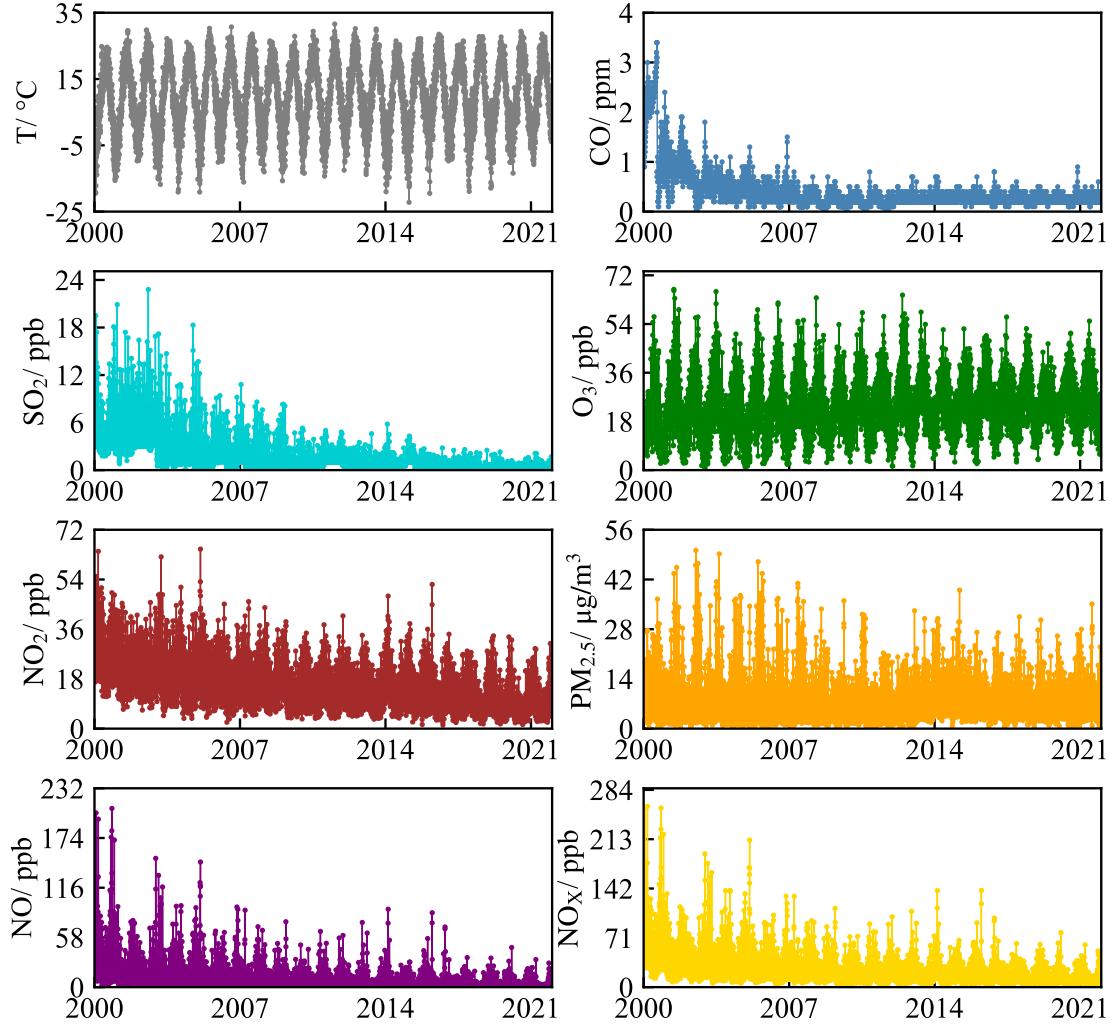


Figure 2.1: Daily mean values of temperature and air pollutants concentration for Toronto, ON, Canada from 2000 to 2021.

$\text{PM}_{2.5}$  and  $\text{PM}_{10}$ . The air pollution data, weather condition and health outcomes are assembled and organized by city in an R package NMMAPSLite [79, 80]. As the updated data after 2000 are never publicly accessible, time periods used for analysis throughout this work are from 1987 to 2000. Data of Chicago are selected for our investigation, which includes the daily mean temperature ( $T$ ),  $\text{CO}$ ,  $\text{SO}_2$ ,  $\text{O}_3$ ,  $\text{NO}_2$ ,  $\text{PM}_{10}$  and daily non-accidental mortality from several diseases. These data used for investigation are visualized and shown in Figure 2.2 after being processed as follows:

- 1) as the original temperature data from Jan. 1, 1998 to Dec. 31, 2000 are missing, online weather data of Chicago [81] for the three years are collected and supplemented (using the Fahrenheit scale).
- 2) the original air pollutants data series are detrended by subtracting a 365 day moving average. They are restored to the true monitor values by adding the trimmed mean and the corresponding trend item [79] to facilitate our following analysis.
- 3) daily non-accidental mortality data for three age groups (under 65, 65 to 74, and above 75) are aggregated as the daily morality observation.
- 4) other small portions of missing values are also interpolated and outliers are dealt with as the above **NAPS** datasets.

It can also be seen in Figure 2.2 that **T** and **O<sub>3</sub>** data series have obvious yearly periodicity as in the **NAPS** datasets. The daily all-cause non-accidental mortality data series show roughly regular fluctuations on a yearly basis as well. Compared to the **NAPS** datasets, other air pollutants series present no evident periodicity or downward trends. In addition, **PM<sub>2.5</sub>** data are not included as most of them are missing.

## 2.2 Mutual Information Based Association

In air pollution epidemiology study, as different air pollutants are physically and chemically mixed with each other in diffusion process and people are generally exposed to different air pollutants at the same time, health impacts assessment models

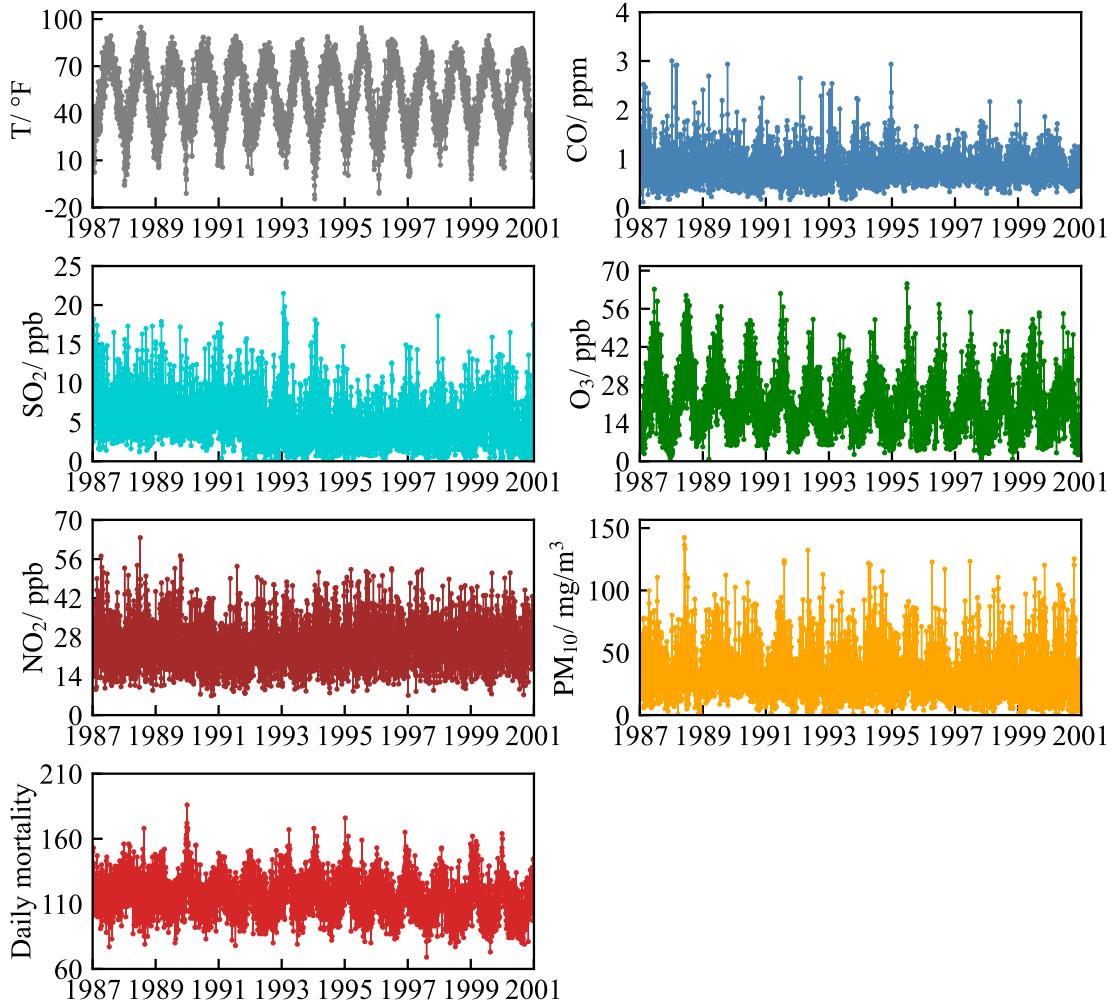


Figure 2.2: Daily mean values of temperature and air pollutants concentration and daily mortality for Chicago, IL, the U.S. from 1987 to 2000.

with multiple air pollutants are more realistic. Thus, correlation analysis for different air pollutants and weather conditions is usually made to identify the most significant risk factor(s) and facilitate the modelling process. However, associations between different risk factors are not always linear in various health assessment problems, and hence the [Pearson Correlation Coefficient \(PCC\)](#) used for relationship analysis and risk factors screening [82] cannot always work properly. The relationship between several air pollutants of Toronto is shown as an example in Figure 2.3. The relation

between NO and  $\text{NO}_x$  is approximately linear, while the linear relation between  $\text{O}_3$  and temperature,  $\text{O}_3$  and  $\text{PM}_{2.5}$  as well as  $\text{NO}_2$  and CO is not obvious. Using PCC to evaluate their relevance may not be able to capture their true association, especially when choosing the air pollutants of interest that contribute to specific health problems.

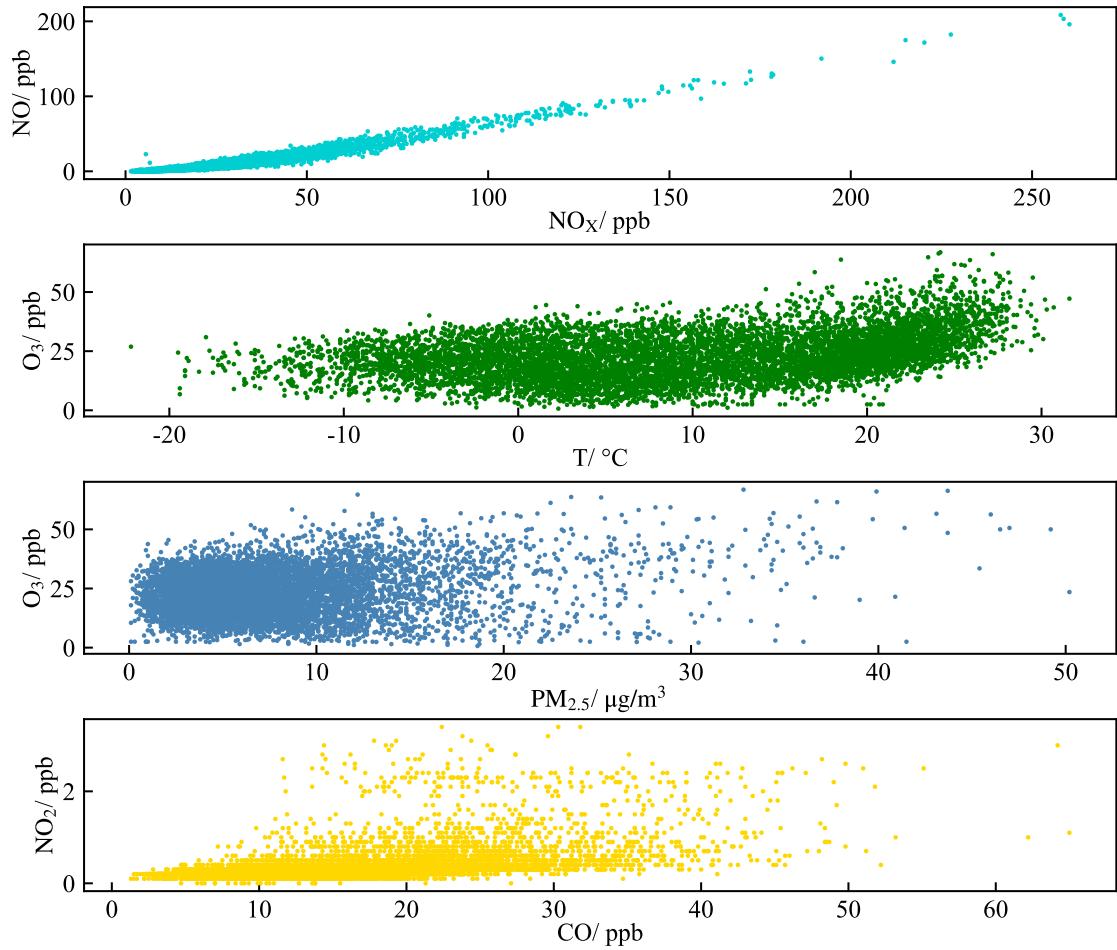


Figure 2.3: Relations of different air pollutants for Toronto, ON, Canada.

In order to evaluate the strength of association between different health risk factors in a general form, MI is used to quantify both the linear and nonlinear relation in this work. MI is a measure of strength in statistics that evaluates mutual

dependence of two random variables, which quantifies how much information of one random variable can be obtained after observing another one using entropy [83]. For two different air pollution variables  $X_1$  and  $X_2$ , their MI is expressed as

$$MI(X_1, X_2) = H(X_1) - H(X_1|X_2), \quad (2.1)$$

where  $H(X_1)$  is the information entropy for  $X_1$ .  $H(X_1|X_2)$  is the conditional information entropy of  $X_1$  after observing  $X_2$ , which can be expressed as:

$$H(X_1|X_2) = H(X_1, X_2) - H(X_2), \quad (2.2)$$

where  $H(X_2)$  is the information entropy for  $X_2$  and  $H(X_1, X_2)$  is the joint entropy of  $X_1$  and  $X_2$ . The information entropy formula [84] for a random variable  $X$  with discrete distribution is

$$H(X) = \sum_x p_X(x) \log_2 p_X(x)^{-1}, \quad (2.3)$$

where  $p_X(x)$  is the probability of taking value  $x$  for variable  $X$ .

With Eqn. (2.2), Eqn. (2.3) and conditional probability formula, MI in Eqn. (2.1) can be further formulated using probabilities as

$$MI(X_1, X_2) = \sum_{x_1} \sum_{x_2} p_{X_1 X_2}(x_1, x_2) \log_2 \left[ \frac{p_{X_1 X_2}(x_1, x_2)}{p_{X_1}(x_1)p_{X_2}(x_2)} \right], \quad (2.4)$$

where  $p_{X_1 X_2}(x_1, x_2)$  is the joint probability distribution of  $X_1$  and  $X_2$ , and  $p_{X_1}(x_1)$  and  $p_{X_2}(x_2)$  are their marginal distributions. A larger MI value indicates a stronger association between  $X_1$  and  $X_2$ , while a lower one indicates that they are weakly associated. Zero value means that they have no statistical correlation.

## 2.3 Maximal Information Coefficient for Air Pollutants Selection

MI in Eqn. (2.4) can capture both linear and nonlinear associations between different air pollutants, while there still exist some limitations when applied to our air pollution datasets. The first issue is that the air pollutants data are measured using different metrics with units on multiple scales. Therefore, MIs between different variables may have various scales, with which comparison of MIs is not available. The second one is that calculation of the joint probability distribution in Eqn. (2.1) is time-consuming and difficult. Therefore, MIC is used for correlation evaluation instead, which is a quantity developed from MI to measure the association between two random variables in a normalized form. Its main idea is to use many possible binning methods to partition the scatter plot of variables into different bins and calculate the MIs with estimation of the joint probability distribution accordingly. Then the gridding method with the maximal normalized MI is chosen as the MIC formulation shown in Eqn. (2.5),

$$MIC(X_1, X_2) = \max_{n_{X_1} \times n_{X_2} \leq N} \frac{MI(X_1, X_2)}{\log_2 \min(n_{X_1}, n_{X_2})}, \quad (2.5)$$

where  $n_{X_1}$  and  $n_{X_2}$  are the numbers of partitions for pairwise air pollutants data of  $X_1$  and  $X_2$ , and  $N$  is a parameter related to the sample size  $S$  and is set to  $S^{0.6}$  in the following experiments. The range of  $MIC(X_1, X_2)$  falls into  $[0, 1]$ , with which the association between pairwise air pollutants measured using different metrics can be evaluated for comparative analysis.

Note that for the formulation in Eqn. (2.5): 1)  $MI(X_1, X_2)$  is the mutual information calculated based on partition of pairwise air pollution samples to  $n_{X_1} \times n_{X_2}$ , with which their joint and marginal probability distributions are estimated accordingly; 2)  $N = S^{0.6}$  is an empirical parameter setting recommended by Reshef *et al* [85]; 3) theoretically, the base of logarithm (also for the above equations with logarithm operation) can be replaced with  $e$ , 10 or other values. Here, base 2 is used to evaluate the amount of information with unit bit. 4) the maximal value of  $MI(X_1, X_2)$  with partitions  $n_{X_1} \times n_{X_2}$  is  $\log_2 \min(n_{X_1}, n_{X_2})$ , and thus  $MI(X_1, X_2)$  is normalized to  $[0, 1]$  through the division operation; 5) in practice, calculating all  $n_{X_1} \times n_{X_2}$  combinations is time-consuming. A sub-optimal  $MIC(X_1, X_2)$  value is usually calculated instead; 6)  $MI(X_1, X_2)$  is a general association metric between  $X_1$  and  $X_2$ , while it does not provide information concerning the direction of correlation as **PCC**; 7) although the performance of **MIC** is arguable in comparison with **MI** [86], it provides a general metric for association measurement and facilitates the comparative analysis.

With **MIC** of Eqn. (2.5), both the linear and nonlinear associations between pairwise air pollution variables can be captured and quantified in a general form based on information entropy. Moreover, this metric can also be used to identify the most significant risk factor(s) for the health consequence of interest through evaluation and sorting of their associations. In the following section, we will show the effect of **MIC** metric in correlation analysis and health risk factors identification through experimental tests on the **NAPS** and **NMMAPS** datasets.

## 2.4 Experiments and Discussions

### 2.4.1 Tools for Experiments

With the above **MIC** formulation, quantitative experiments and analysis are made based on the aforementioned datasets. Several statistical tools are utilized to process the original data and facilitate the experimental process and the main tools are briefly introduced as follows.

- i) The experiments are mainly performed using Python (version 3.7.13) language in a virtual environment.
- ii) SciPy (version 1.7.3) [87] is used for calculation of **PCCs** with Python. SciPy is a scientific computation library designed for efficient numerical computation based on NumPy, including various tools for applications in optimization, linear algebra, signal and image processing, etc. Its **scipy.stats** submodule for statistics is utilized for the following experiments.
- iii) minepy (version 1.2.6) [88] is used for calculation of **MICs** with Python. minepy is a library developed for efficient computation of maximal information-related association. Its **minepy.MINE** submodule is utilized for computation of **MICs** in the following experiments.

### 2.4.2 Numerical Results

First, experiments on the **NAPS** datasets are performed to show the effect of **PCC** and **MIC** in association evaluation. Then **MIC** is used to identify the most significant

risk factor(s) in different health consequences based on the [NMMAPS](#) datasets.

The association of different air pollutants (including temperature) evaluated using [PCCs](#) (absolute values) and [MICs](#) for Toronto, ON, Canada are presented as in the heat map of Figure 2.4. It can be seen from the figure that

- a. The values of correlation measured using [MICs](#) are relatively milder than using [PCCs](#) in general.
- b. Compared to [PCCs](#), the influence of [T](#) on concentrations of air pollutants is smaller evaluated using [MICs](#) except for [NO](#) and [NO<sub>x</sub>](#). Both [PCCs](#) and [MICs](#) show that [CO](#), [NO](#), [NO<sub>2</sub>](#), [NO<sub>x</sub>](#) and [SO<sub>2</sub>](#) are weakly associated with [T](#), while [MICs](#) show a relatively stronger relationship between [NO](#) or [NO<sub>x</sub>](#) and [T](#). This implies that [PCCs](#) are not able to capture all the association and cannot fully reflect their relevance for risk factors screening and variable selection.
- c. [CO](#) and [PM<sub>2.5</sub>](#) are correlated with the other six elements evaluated using both [PCCs](#) and [MICs](#), with [MICs](#) showing weaker relations. [SO<sub>2</sub>](#) and [O<sub>3</sub>](#) are also shown to have some relevance with other five pollutants using both [PCCs](#) and [MICs](#). However, [MIC](#) shows a mild relationship between [SO<sub>2</sub>](#) and [O<sub>3</sub>](#), where their [PCC](#) shows they are almost independent. This is another case that [PCC](#) does not fully capture the association between different pollutants.
- d. Although the [MIC](#) values are relatively smaller, both [PCCs](#) and [MICs](#) show strong correlations of [NO](#), [NO<sub>2</sub>](#), and [NO<sub>x</sub>](#), as these pollutants usually come from human activity and have the same sources like industrial emission, fossil fuel electric utilities and motor vehicles.

e. Compared to **PCC**, as **MIC** can capture multiple types of association instead of only linearity, it provides a more effective approach for dependency analysis and has the potential to facilitate risk factor(s) selection for air pollution-related health assessment models like **GAMs** as well as analysis of their risk contributions therein.

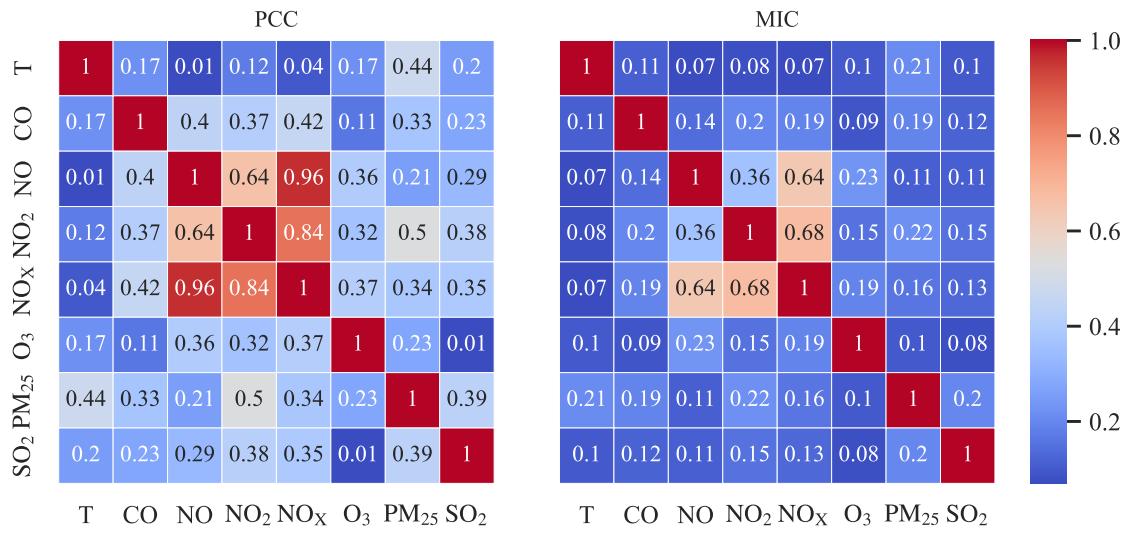


Figure 2.4: Comparison of PCCs and MICs with data of all seasons for Toronto, ON, Canada.

As the health risk in air pollutant-related risk assessment models like **GAMs** is usually evaluated for different seasons, relations of different air pollutants using data of both warm (from April to September) and cold (from October to March) seasons are evaluated and presented in Figure 2.5 and Figure 2.6, respectively. It can be seen from the figures that in both cold and warm seasons association evaluated using **MICs** are milder than using **PCCs** in general, following the patterns in Figure 2.4. However, **MIC** values for **T** and **NO**, **T** and **NO<sub>2</sub>**, and **T** and **NO<sub>x</sub>** in cold seasons, as well as **MIC** values for **T** and **NO<sub>x</sub>**, **CO** and **O<sub>3</sub>**, **NO** and **PM<sub>2.5</sub>**, **NO<sub>2</sub>**

and  $O_3$  in warm seasons are all larger than their **PCC** values. This means that in these different season scenarios **MIC** still shows its capability to capture more association between different air pollutants compared to **PCC**, and that **PCC** may lead to incorrect association conclusions in nonlinear situations.

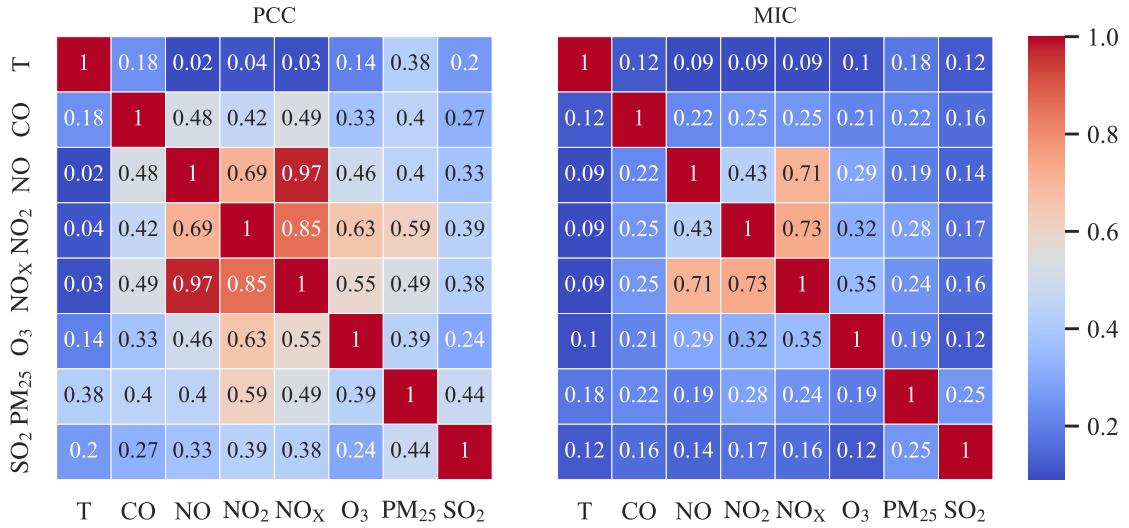


Figure 2.5: Comparison of PCCs and MICs with data of cold seasons (from October to March) for Toronto, ON, Canada.

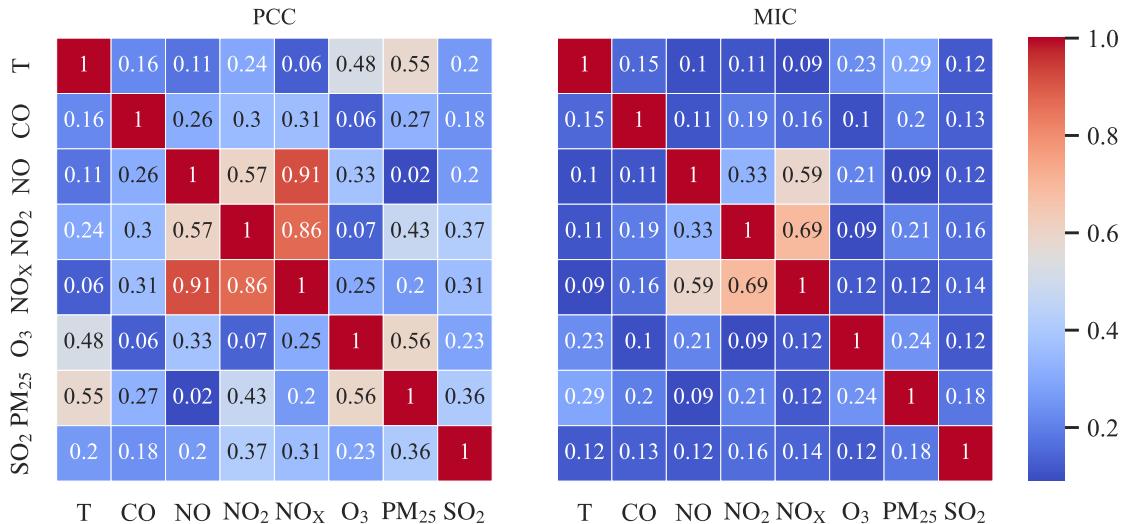


Figure 2.6: Comparison of PCCs and MICs with data of warm seasons (from April to September) for Toronto, ON, Canada.

Note that in the above analysis: 1) values of **PCC** are mainly used for compar-

ison and they might be meaningless in situations where the association of different elements is not linear; 2) values of MIC mean more for comparative analysis than for correlation strength identification; 3) MIC measures the general correlation without reporting correlation type or direction.

Experiments are further performed using datasets of Chicago, IL, the U.S. MIC is used to evaluate the correlation between different risk factors and mortality from various causes, in order to identify the most significant air pollution (and temperature) factors concerning the outcome. Test results are presented in Table 2.1 and it can be seen that

- a. Generally, older people (Age-3 category) are more sensitive to T than the other two categories for different health outcomes. SO<sub>2</sub> has relatively less influence on all health outcomes of three age categories compared to other risk factors. CO has relatively more influence on all health outcomes of three age categories than NO<sub>2</sub>, O<sub>3</sub> and PM<sub>10</sub>.
- b. For different health outcomes of age categories Age-1 and Age-2, the influence of PM<sub>10</sub> is greater than that of both NO<sub>2</sub> and O<sub>3</sub>. For the health outcomes PNEINF, PNEU, and RESP of Age-3, PM<sub>10</sub> is also the most associated factor compared with NO<sub>2</sub> and O<sub>3</sub>, while for all-cause mortality, CVD and COPD of Age-3, O<sub>3</sub> has the maximum relevance of the three air pollutants.
- c. For age category Age-1, CO has the greatest impacts on all health outcomes. With age increasing (from Age-2 to Age-3), T gradually becomes the most associated risk factor for different health outcomes except COPD, where CO keeps

the most associated one, although the influence of  $T$  increases and becomes the second most related factor.

- d. For age category Total, the association between each health outcome and different air pollutants (temperature) is influenced by composition of the three age categories. All the health outcomes of Total are most associated with  $T$  and  $CO$  and least associated with  $SO_2$ . Both  $O_3$  and  $PM_{10}$  have greater impacts than  $NO_2$  for different health outcomes.
- e. Besides the above trend, for a specific health consequence, the most associated risk factor(s) can be selected by sorting the corresponding  $MIC$  values in an descending order.

For practical application in air pollution epidemiology, it is worth mentioning that: 1)  $MIC$  is a normalized form of  $MI$  that can evaluate the association between two different risk factors (air pollutants exposure, temperature, etc.), or between a risk factor and the health outcome for comparative analysis and main risk factors identification; 2)  $MIC$  is a measurement between two random variables, which cannot decide if the impact of one risk factor on health outcome of interest outweighs those of several other factors; 3) from both visual analysis of Figure 2.2 and the  $MIC$  values in Table 2.1, temperature ( $T$ ) plays an important role in daily all-cause non-accidental mortality and mortality from several other diseases, especially for elder people, and it might be worth further investigation of its role and formulation in  $GAMs$ -based health outcome assessment models.

Health Outcomes	Age Category	MIC					
		T	CO	NO <sub>2</sub>	O <sub>3</sub>	PM <sub>10</sub>	SO <sub>2</sub>
All-cause Mortality	Age-1	0.0637	0.0722	0.0389	0.0478	0.054	0.0314
	Age-2	0.0777	0.0736	0.0393	0.0418	0.0478	0.0347
	Age-3	0.1217	0.0729	0.0413	0.0681	0.0553	0.0273
	Total	0.1443	0.0748	0.0485	0.0624	0.0575	0.0321
CVD	Age-1	0.0634	0.0732	0.0431	0.0422	0.0503	0.0302
	Age-2	0.0686	0.0719	0.0397	0.0409	0.0504	0.0372
	Age-3	0.1161	0.0779	0.0416	0.0664	0.0524	0.0277
	Total	0.1186	0.0767	0.0409	0.0597	0.0595	0.0358
COPD	Age-1	0.0503	0.0695	0.0451	0.0421	0.0532	0.0268
	Age-2	0.0468	0.0746	0.039	0.0443	0.0498	0.0243
	Age-3	0.0579	0.0722	0.0415	0.0496	0.047	0.0267
	Total	0.0595	0.0681	0.0389	0.0443	0.0491	0.0271
PNEINF	Age-1	0.0549	0.0761	0.0382	0.0441	0.0551	0.0295
	Age-2	0.0588	0.0739	0.0389	0.0381	0.0495	0.0248
	Age-3	0.0861	0.0746	0.0406	0.0517	0.0565	0.0224
	Total	0.0937	0.0735	0.0412	0.0551	0.0553	0.023
PNEU	Age-1	0.0546	0.0763	0.0381	0.0432	0.055	0.0294
	Age-2	0.0582	0.0732	0.0387	0.038	0.0494	0.0248
	Age-3	0.0857	0.0746	0.0409	0.052	0.0562	0.0225
	Total	0.0923	0.0737	0.0408	0.0553	0.0551	0.0227
RESP	Age-1	0.0581	0.0708	0.0418	0.0477	0.0562	0.0252
	Age-2	0.055	0.076	0.0386	0.0442	0.0491	0.0226
	Age-3	0.0891	0.0767	0.0449	0.0561	0.057	0.027
	Total	0.0954	0.0672	0.0432	0.0581	0.0504	0.0284

Table 2.1: MICs between air pollutants and daily mortality from different causes for Chicago, IL, the U.S. All-cause Mortality is the all-cause non-accidental daily mortality. CVD, COPD, PNEINF, PNEU are RESP are short for cardiovascular deaths, chronic obstructive pulmonary disease, pneumonia and influenza, pneumonia and respiratory deaths. They represent the daily mortality from corresponding causes. Age-1 is for under 65 years of age, Age-2 is for 65 to 74, and Age-3 is for over 75.

# Chapter 3

## LSTM NetWork

In this chapter, an [LSTM](#) model is developed to evaluate the adverse impacts of air pollutants on public health. [LSTM](#) network is used to extract dependencies from air pollution time series with distributed lags and the dependent features are utilized for health outcomes assessment.

### 3.1 Introduction of LSTM Network

[LSTM](#) network is a kind of artificial neural network proposed by Hochreiter and Schmidhuber [89] in 1997. It is an improved version of [Recurrent Neural Network \(RNN\)](#) [90, 91] that is developed to deal with sequential data and originates from the recursive neural network, Hopfield network [92]. The elementary structure of [RNN](#) is shown in Figure 3.1, where  $\mathbf{x}_{t-1}$ ,  $\mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  are sequential vector inputs,  $\mathbf{h}_{t-1}$ ,  $\mathbf{h}_t$  and  $\mathbf{h}_{t+1}$  are sequential outputs, and the repeating module  $\mathbf{A}$  is some neural network. Arrow represents the direction of information flow.

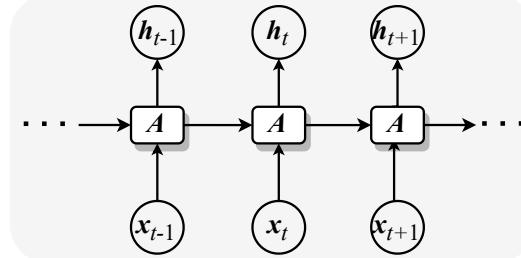


Figure 3.1: Elementary structure of RNN

One unique feature of the [RNN](#) in Figure 3.1 is that information from the previous processing module  $\mathbf{A}$  can be utilized by its successor. The passage of “recurrent” information makes [RNN](#) the natural to deal with sequential data like voice or video sequences, as well as other kinds of time series data. However, the following several main disadvantages [93, 94] limit its application in practice:

- 1) Basic [RNNs](#) suffer from the problem of exploding or vanishing gradients. As the algorithm of backpropagation through time (BPTT) is usually applied for model training, the error gradients are propagated and calculated based on the chain rule. Once the gradient of a processing module  $\mathbf{A}$  becomes smaller or larger, multiplications of these gradients from  $\mathbf{As}$  of different timestamps tend to explode or vanish, which makes [RNNs](#) difficult to train.
- 2) Basic [RNNs](#) can hardly process long-term sequential data. Influenced by the gradient issues, it is difficult for [RNNs](#) to capture the long-term dependencies in longer sequences, where the training process becomes unstable.
- 3) Other important but more general issues met by most artificial neural networks like complexity of training, overfitting and generalization, depth of the networks, etc., also exist for traditional [RNNs](#).

Aiming at the above issues of [RNN](#), [LSTM](#) networks are specifically designed to mitigate the gradient exploding or vanishing problem and the long-term dependency problem. The key ideas of [LSTM](#) networks is to use a “cell state” to store information and use “three gates” to regulate the information flow. More specifically, the “cell state” stores information over time, where for each timestamp it interacts with the “three gates” for state update. The first step is the interaction with the “forget gate”, which decides how much old information from previous timestamp needs to be removed from the “cell state”. Then the “input gate” regulates how much new information from present timestamp is to be added to the “cell state”. At last, through the “output gate”, the “cell state” is optionally passed to the next time step after manipulation and update. A detailed mathematical description of the above process used in the air pollution-related health assessment problem is presented in the following section [LSTM Model for Health Outcomes Assessment](#).

With the gate mechanism regulating information flow and alleviating the gradient issues, [LSTM](#) network can store “short-term memory” in “long” time steps and capture the long-term dependency. Thus, it has been successfully applied to a series of machine learning problems that need to handle the time-dependency issue across long time gaps of sequential data, such as natural language processing [95], speech recognition [96], machine translation [97], and time series forecasting [98], etc. As to the air pollution-related health impacts problem, people are generally exposed to multiple air pollutants across long periods before some health events occur, where the air pollutant concentration levels are typical time series. This inspires the exploration

of applying **LSTM** network to capture the dependent impacts of previous air pollution exposure on future public health consequences, i.e., distributed lags in this research.

## 3.2 Data Reorganization

In order to facilitate the following formulation and description of **LSTM** network-based health consequences assessment model, air pollution time series data and health outcome data of [section 2.1](#) are reorganized as follows.

First, the original air pollutants data and health outcomes data are scaled through standardization as Eqn. (3.1),

$$x_{s,i} = \frac{x_i - \mu_x}{\sigma_x}, \quad (3.1)$$

where  $x_i$  is the  $i$ -th data in the original time series datasets for each category,  $\mu_x$  and  $\sigma_x$  are the mean and standard deviation values, and  $x_{s,i}$  is the corresponding standardized value. After standardization, each category of the original data is scaled to have a mean value of 0 and a standard deviation of 1. The main purposes of the above standardization operation are to

- 1) reduce the influence of scales from different kinds of data, as ambient air pollutants are generally measured and recorded with different metrics and scales;
- 2) accelerate the process of searching for the optimal network parameters during the training of **LSTM** model;
- 3) improve the generalization ability of the network by reducing overfitting.

Second, the datasets are further reorganized with each element being a vector

consisting of a series of the standardized data. Assume the length of vector after reorganization is  $m$ , the vectorized data sample can then be expressed as,

$$\mathbf{x}_{s,i} = [x_{s,i-m+1}, x_{s,i-m+2}, \dots, x_{s,i}], \quad (3.2)$$

where  $\mathbf{x}_{s,i}$  is the  $i$ -th reorganized data sample. The new reorganized datasets are

$$[\mathbf{x}_{s,1}, \mathbf{x}_{s,2}, \dots, \mathbf{x}_{s,n-m+1}], \quad (3.3)$$

where  $n$  is the length of original data series and it becomes  $n-m+1$  after reorganization. The purpose of above operation is to facilitate the inputs during training of the following **LSTM** model. It is worth mentioning that parameter  $m$  has the same meaning as the time lag used in **GAM**-type health assessment models and the expression in Eqn. (3.2) includes distributed lags which means sequential exposures from timestamp  $i - m + 1$  to timestamp  $i$ .

### 3.3 LSTM Model for Health Outcome Assessment

#### 3.3.1 LSTM Network for Feature Information Extraction

For some public health outcome of interest at the time period  $t$ , the **LSTM** network used for sequential information extraction from multiple air pollution time series of length  $m$ , i.e., exposure from previous  $m$  time periods is shown in Figure 3.2. The **LSTM** network has  $r$  stacked layers to process the input sequence and each layer has  $m$  repeating neural network modules  $\mathbf{A}$ . The input to the **LSTM** model is  $[\mathbf{x}_{s,t-m+1}, \mathbf{x}_{s,t-m+2}, \dots, \mathbf{x}_{s,t}]$ , i.e., multiple air pollutant concentration (exposure) sequence of previous  $m$  periods, and its output is the extracted feature vector sequence

with temporal dependency as  $[\mathbf{h}_{r,t-m+1}, \mathbf{h}_{r,t-m+2}, \dots, \mathbf{h}_{r,t}]$ .

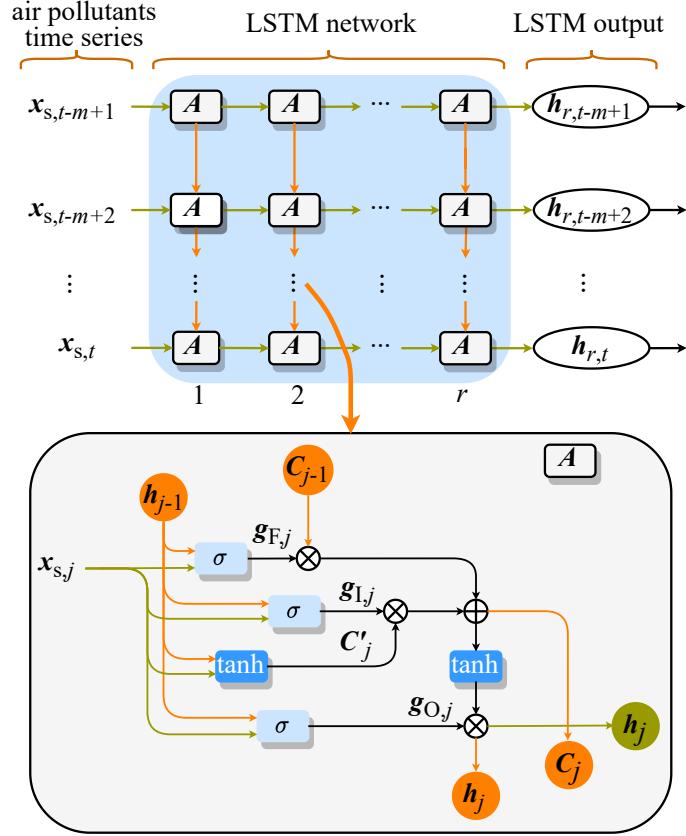


Figure 3.2: LSTM network for sequential information extraction

Note that in Figure 3.2: 1) each element  $\mathbf{x}_{s,j}$  ( $j = t - m + 1, \dots, t.$ ) of the input air pollution sequence is a vector consisting of multiple air pollutants, which is a generalized expression of the scalar element in Eqn. (3.2) for single air pollutant; 2) dimension of the output of module  $\mathbf{A}$  for each layer, i.e.,  $\mathbf{h}_{k,j}$  ( $k = 1, \dots, r, j = t - m + 1, \dots, t.$ ) is  $l$ , which can be different from the dimension of input  $\mathbf{x}_{s,j}$ ; 3)  $m$ ,  $r$  and  $l$  are three parameters that determine the network structure.

The repeating module  $\mathbf{A}$  uses three specifically designed ‘‘gates’’ and a ‘‘cell state’’ to regulate and extract information from the input, where ‘‘gates’’ are essentially sigmoid functions (between 0 and 1) that decide the information flow, with 0 for

“completely blocked” and 1 for “completely passed”, and “cell state” is used to store the dependent information extracted from the input sequence. As operations for different layers are alike, layer number  $k$  is neglected in the following description.

For period  $j$ , how much information from “cell state”  $\mathbf{C}_{j-1}$ , i.e., stored information from last period  $j - 1$  should be blocked is decided by  $\mathbf{g}_{F,j}$  of the “forget gate” as

$$\mathbf{g}_{F,j} = \sigma(\mathbf{w}_F \cdot \mathbf{x}_{s,j} + \mathbf{u}_F \cdot \mathbf{h}_{j-1} + \mathbf{b}_F), \quad (3.4)$$

where  $\mathbf{w}_F$ ,  $\mathbf{u}_F$  and  $\mathbf{b}_F$  are network weights and biases.  $\sigma(\cdot)$  is sigmoid function. Value of  $\mathbf{g}_{F,j}$  is decided by both the air pollution concentration  $\mathbf{x}_{s,j}$  in the  $j$ -th time period and extracted information  $\mathbf{h}_{j-1}$  from the  $j - 1$ -th time period.

Then the “input gate” as Eqn. (3.5) is used to decide how much new information from the candidate “cell state” as Eqn. (3.6) that could be added to  $\mathbf{C}_{j-1}$ ,

$$\mathbf{g}_{I,j} = \sigma(\mathbf{w}_I \cdot \mathbf{x}_{s,j} + \mathbf{u}_I \cdot \mathbf{h}_{j-1} + \mathbf{b}_I) \quad (3.5)$$

$$\mathbf{C}'_j = \tanh(\mathbf{w}_C \cdot \mathbf{x}_{s,j} + \mathbf{u}_C \cdot \mathbf{h}_{j-1} + \mathbf{b}_C), \quad (3.6)$$

where  $\mathbf{w}_I$ ,  $\mathbf{u}_I$ ,  $\mathbf{w}_C$ , and  $\mathbf{u}_C$  are network weights.  $\mathbf{b}_I$  and  $\mathbf{b}_C$  are network biases.  $\tanh$  is the hyperbolic tangent activation function. Together with the “forget gate”, the “cell state” output from the  $j$ -th module is updated as  $\mathbf{C}_j$  in Eqn. (3.7),

$$\mathbf{C}_j = \mathbf{g}_{F,j} \odot \mathbf{C}_{j-1} + \mathbf{g}_{I,j} \odot \mathbf{C}'_j, \quad (3.7)$$

where  $\odot$  is the Hadamard product. The new “cell state”  $\mathbf{C}_j$  blocks part of old information  $\mathbf{C}_{j-1}$  from last time period  $j - 1$  and adds part of new information  $\mathbf{C}'_j$  from current time period. Finally,  $\mathbf{C}_j$  is passed to the next time period  $j + 1$  as well as to the next network layer after being processed by the “output gate” as Eqn.

(3.8) and Eqn. (3.9),

$$\mathbf{g}_{\text{O},j} = \sigma(\mathbf{w}_{\text{O}} \mathbf{x}_{\text{s},j} + \mathbf{u}_{\text{O}} \mathbf{h}_{j-1} + \mathbf{b}_{\text{O}}) \quad (3.8)$$

$$\mathbf{h}_j = \mathbf{g}_{\text{O},j} \odot \tanh(\mathbf{C}_j), \quad (3.9)$$

where  $\mathbf{w}_{\text{O}}$ ,  $\mathbf{u}_{\text{O}}$  and  $\mathbf{b}_{\text{O}}$  are weights and biases. The output  $\mathbf{h}_j$  of a module  $\mathbf{A}$  is decided by both the updated “cell state”  $\mathbf{C}_j$  and the “output gate”  $\mathbf{g}_{\text{O},j}$ .

With the above regulation through  $m$  time steps of  $r$  layers, dependent information  $\mathbf{h}_{r,j}$  ( $j = t - m + 1, t - m + 2, \dots, t$ ) is extracted from the input air pollution sequence, which can be further used to assess the air pollution-related health outcome of interest. Note that in the above description: 1) the inputs to each layer of LSTM network are outputs of module  $\mathbf{As}$  from previous layer except for the first one where its inputs are the air pollution sequences; 2) for the first time period of each layer ( $j=1$ ), its input  $\mathbf{h}_{j-1}$  and  $\mathbf{C}_{j-1}$  from last time step are zeros; 3) the meaning of extracted feature vector  $\mathbf{h}_{r,j}$  ( $j = t - m + 1, t - m + 2, \dots, t$ ) is determined by the health outcome of interest and may be different for various health events.

### 3.3.2 Health Outcome Assessment

With the sequential information  $[\mathbf{h}_{r,t-m+1}, \mathbf{h}_{r,t-m+2}, \dots, \mathbf{h}_{r,t}]$  extracted from the input air pollution sequence through the LSTM network, health outcome of interest are assessed using combinations of these features. As these extracted features from previous exposure of  $m$  time periods may have different contributions to the health outcome at time period  $t$ , weighed evaluation of the lagged effects is used for health outcome assessment as shown in Figure (3.3)

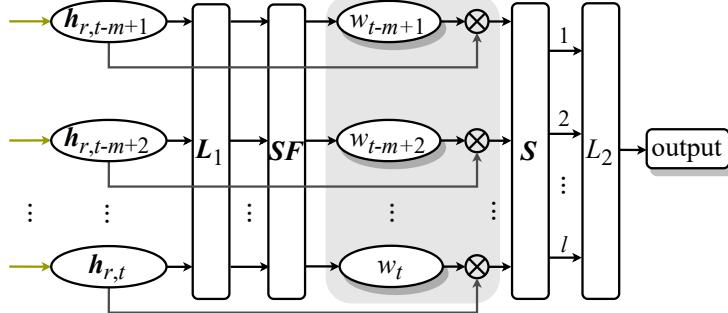


Figure 3.3: Health outcome assessment with weighted evaluation of lags

First, a linear transformation layer  $\mathbf{L}_1$  with  $m$  processing units is used to combine the  $l$  features in each  $\mathbf{h}_{r,j}, j = t - m + 1, t - m + 2, \dots, t$  as Eqn. (3.10),

$$\mathbf{L}_1(j) = \sum_{i=1}^l w_{i,j} \mathbf{h}_{r,j}(i) + b_j, j = t - m + 1, t - m + 2, \dots, t. \quad (3.10)$$

where  $w_{i,j}$  and  $b_j$  are weight and bias.

Then, a softmax layer  $\mathbf{SF}$  is used to generate weights for the output sequence of  $\mathbf{L}_1$  by mapping their values to  $(0,1)$  as Eqn. (3.11),

$$w_j = \mathbf{SF}(j) = \frac{e^{\mathbf{L}_1(j)}}{\sum_{j=t-m+1}^t e^{\mathbf{L}_1(j)}}, j = t - m + 1, t - m + 2, \dots, t. \quad (3.11)$$

where  $w_j$  is the weight evaluating the influence of air pollution exposure in time period  $j$  on some health outcome.

After that, a sum layer  $\mathbf{S}$  is used to combine each feature of  $\mathbf{h}_{r,j}$  from  $m$  time periods with weighted evaluation of their impacts as Eqn. (3.12),

$$\mathbf{S}(i) = \sum_{j=t-m+1}^t w_j \mathbf{h}_{r,j}(i), i = 1, 2, \dots, l. \quad (3.12)$$

At last, the impacts of  $l$  features are combined to assess the health outcome of interest through a linear transformation layer  $L_2$  as Eqn. (3.13).

$$output = L_2[\mathbf{S}(i), i = 1, 2, \dots, l] = \sum_{i=1}^l w'_i \mathbf{S}(i) + b, \quad (3.13)$$

where  $w'_i$  and  $b$  are weight and bias.

Note that in the above formulation: 1) exposure to air pollution from  $t - m + 1$  to  $t$  is used to assess the health outcome at  $t$ , where other periods of interest can be used instead; 2)  $w_j$  in Eqns. (3.11) and (3.12) is the weight for the  $j$ -th feature vector  $\mathbf{h}_{r,j}$  and is shared by all  $l$  features therein to evaluate the impacts of distributed previous exposures; 3)  $w'_i$  in Eqn. (3.13) is the weight evaluating combined impacts of the  $i$ -th element in all the extracted feature vectors; 4) the LSTM network is used to process and extract chronically dependent information from the input multiple air pollutants series through each layer, as well as to capture and fit the relation between input sequence and health outcome through multiple layers together with the subsequent evaluation layers; 5) output from  $L_2$  is the health outcome of interest at time period  $t$ , which can be mortality count, morbidity count, cardiovascular disease count, etc.

### 3.4 Results and Discussions

Performance of the proposed model was tested through numerical experiments on datasets of Chicago. Daily all-cause non-accidental mortality for all three age groups was used for assessment. Besides T, air pollutants of interest were selected from NO<sub>2</sub>, CO, SO<sub>2</sub>, O<sub>3</sub> and PM<sub>10</sub>. All categories of data were preprocessed, reorganized and split into two parts, with 70% for model training and the rest 30% for performance evaluation. The following three experiments were performed.

- i) For different lengths of input air pollutants sequence  $m$ , experiments were performed to demonstrate the performance in capturing the impacts of distributed

lags (previous exposure).

- ii) Single-pollutant and multipollutant models are constructed and tested to show the performance in adapting to both single and multiple air pollutant (s).
- iii) Comparison with a standardized **GAM** model widely used in air pollution epidemiology as Eqn. (3.14) is made to illustrate the advantage of proposed model,

$$\log(\mu_t) = \beta_0 + \beta_1 \chi_t + \beta_2 DOW_t + \text{NS}(T, df = 3)_t + \text{NS}(Time, df = 6/\text{yr})_t \quad (3.14)$$

where  $\mu_t$  is the daily mortality,  $\chi_t$  is the daily mean concentration level of air pollutant of interest,  $\beta_0$ ,  $\beta_1$  and  $\beta_2$  are linear coefficients,  $DOW_t$  is a binary indicator function for day of week with 0 for weekdays and 1 for weekends, and NS is the natural cubic spline regression function with degree of freedom  $df = 3$  for daily mean temperature  $T$  and  $df = 6$  per year for  $Time$ .

### 3.4.1 Performance Metrics

Performance metrics used in all numerical experiments are **Root Mean Square Error (RMSE)** and **Mean Absolute Error (MAE)** as Eqns. (3.15) and (3.16), with RMSE evaluating the deviation between predicted daily mortality and its actual value and MAE evaluating the absolute prediction error,

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_{P,i} - y_{T,i})^2} \quad (3.15)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_{P,i} - y_{T,i}| \quad (3.16)$$

where  $y_{P,i}$  and  $y_{T,i}$  are the  $i$ -th predicted daily mortality and its actual value, and  $n$  is the total number of samples used. Both metrics are calculated using the standardized

air pollutants and health outcome data in all the following experiments.

### 3.4.2 Parameter Settings

PyTorch (version 1.8.1+cpu) is used to construct the [LSTM](#) model. Main hyper parameters of of the [LSTM](#) network during training and test are listed in Table 3.1. The optimizer used is the [Adaptive Moment Estimation \(ADAM\)](#) algorithm, where a variable learning rate with a decay rate of 0.95 per 100 epochs is used to accelerate the convergence of model training process.

Parameter	Value	Description
$r$	3	number of hidden layers
$m$	1–12	length of input sequence
$l$	13	dimension of extracted feature
$epochs$	5000	training time periods
$lr$	5e-2	initial learning rate
$decay$	0.95	decay of $lr$ every 100 epochs
$seed$	65	random seed

Table 3.1: Main parameters setting for the LSTM model.

### 3.4.3 Numerical Results

#### Tests with different $m$

With  $\text{T}$ ,  $\text{PM}_{10}$  and  $\text{O}_3$ , losses during the training process, performance metrics on both training and test sets, and comparison between predicted mortality and actual count for different lengths of input air pollution sequence  $m$  are shown in Fig. 3.4, Table 3.2 and Fig. 3.5. Variation of losses in Fig. 3.4 demonstrates

that for different lengths of input air pollution sequence the proposed **LSTM** model has good convergence properties, and lagged impacts of distributed exposure can be accommodated for daily mortality assessment. Performance metrics in Table 3.2 shows that both the **RMSE** and **MAE** values have general downward trends on the training set with the length of input exposure sequence, i.e.,  $m$  increasing. With a longer term of exposure sequence, more useful features and information can be extracted through the **LSTM** network for mortality assessment. For issues of data noises, distribution differences and possible overfitting, downward trends of the performance metrics are not obvious on the test set.

	RMSE		MAE			RMSE		MAE	
	TrS	TeS	TrS	TeS		TrS	TeS	TrS	TeS
$m=1$	0.614	1.357	0.441	1.040	$m=7$	0.080	1.459	0.047	1.156
$m=2$	0.192	1.909	0.098	1.455	$m=8$	0.048	1.387	0.030	1.114
$m=3$	0.095	1.679	0.042	1.323	$m=9$	0.055	1.460	0.032	1.157
$m=4$	0.094	1.574	0.047	1.219	$m=10$	0.066	1.394	0.034	1.108
$m=5$	0.092	1.561	0.050	1.223	$m=11$	0.059	1.398	0.031	1.127
$m=6$	0.042	1.472	0.026	1.145	$m=12$	0.049	1.372	0.028	1.101

Table 3.2: Performance comparison for different  $m$  on the training and test sets (TrS: training set, TeS: test set).

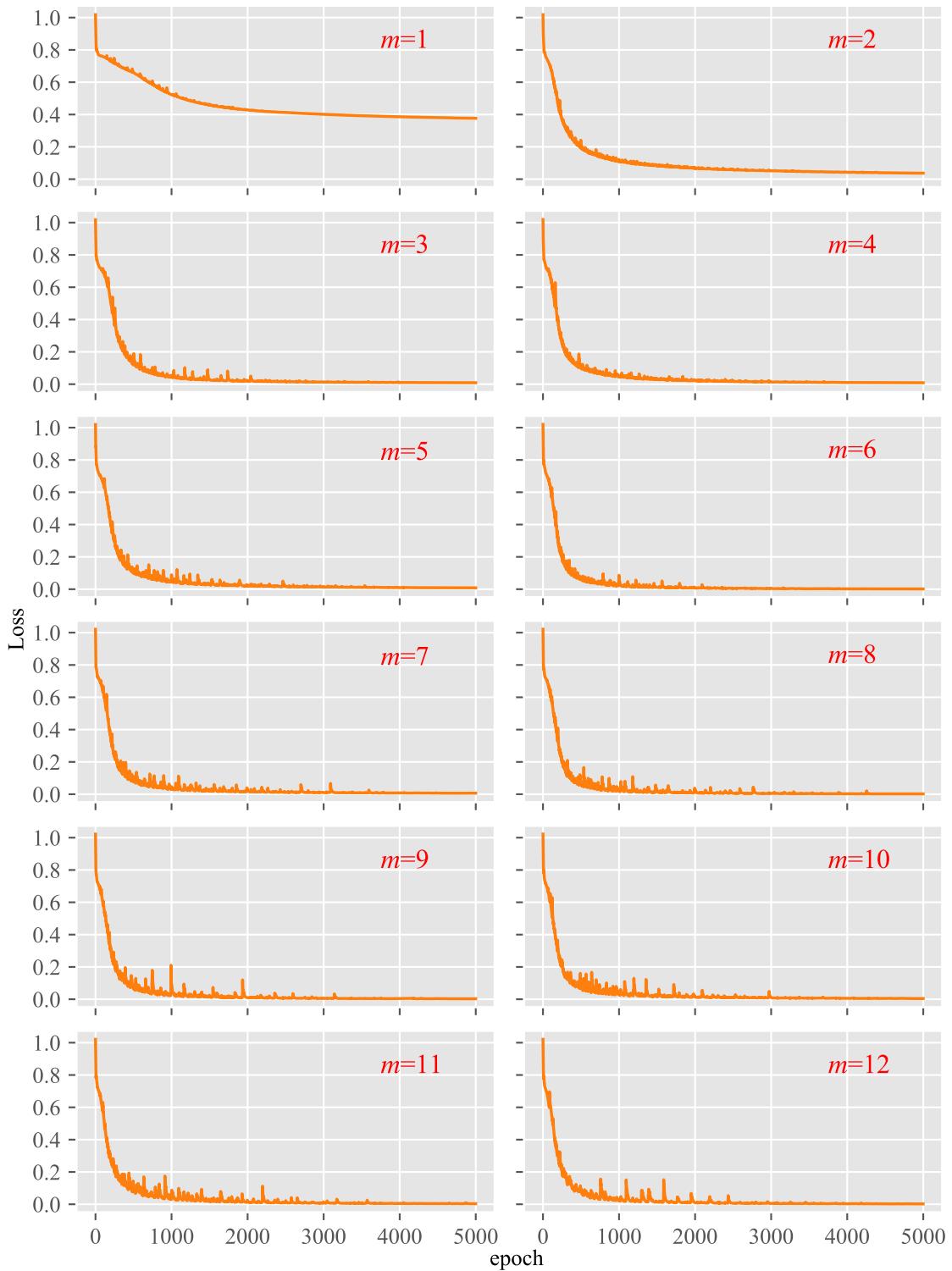


Figure 3.4: Training losses for different lengths of air pollution exposure sequence ( $m=1\text{--}12$ ).

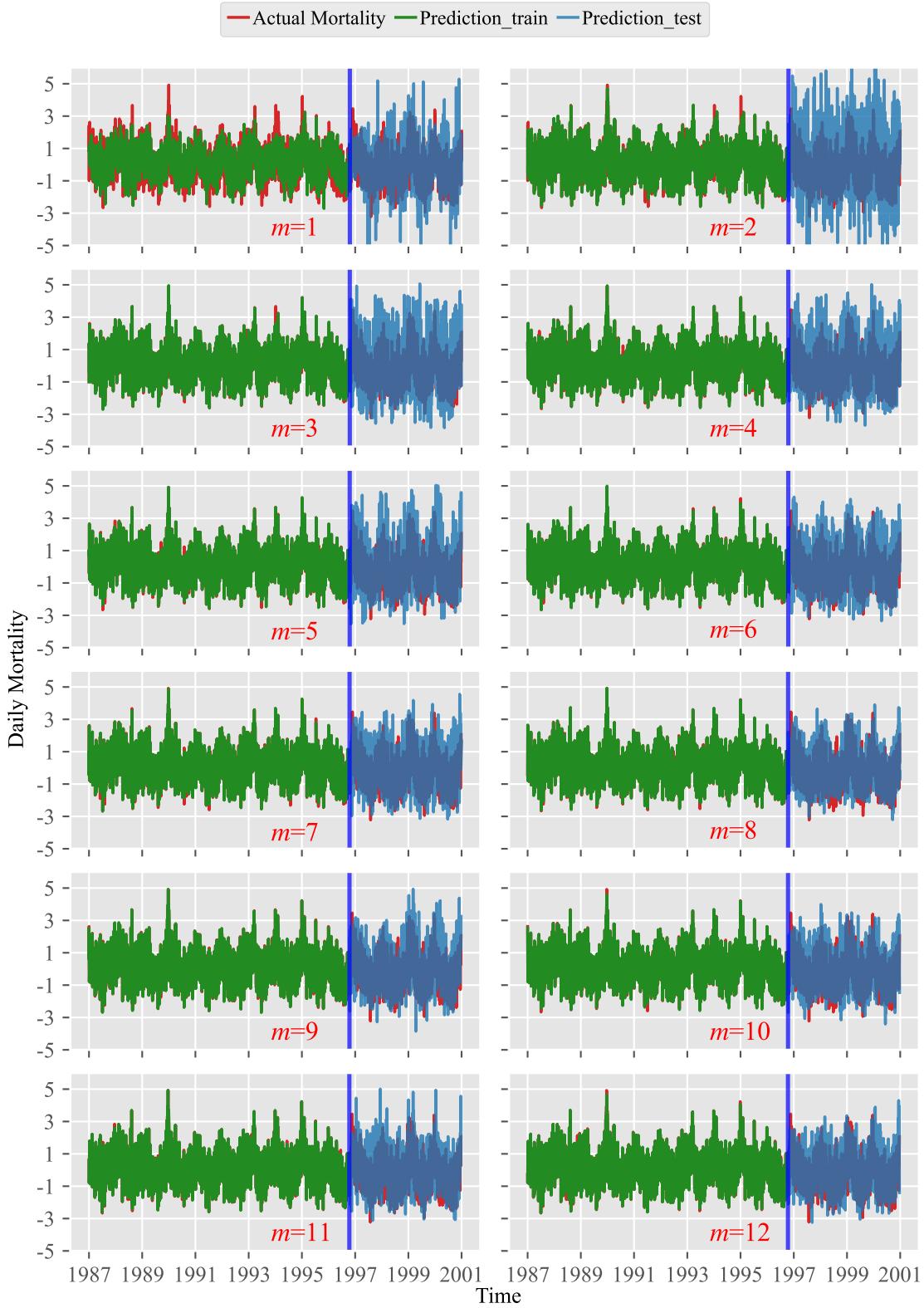


Figure 3.5: Comparison between prediction and actual daily mortality ( $m=1\text{--}12$ ).

Another two experiments were performed to further demonstrate the performance of the proposed model with different lengths of air pollution sequence  $m$ . Besides  $\text{T}$ , the air pollutants were selected based on the **MIC** between each pollutant and the mortality, and the most associated three and four air pollutants, i.e.,  $\text{PM}_{10} + \text{O}_3 + \text{CO}$  and  $\text{PM}_{10} + \text{O}_3 + \text{CO} + \text{NO}_2$  were used for tests. Experimental results are shown in Table 3.3 and Fig. 3.6 to Fig. 3.9. Similar to Fig. (3.4), the training losses in Figs. 3.6 and 3.8 further demonstrate the good convergence properties of the model, and comparisons in Figs. 3.7 and 3.9 roughly show that the **LSTM** model can capture short-term variations and long-term trends in both tests. Downward trends are not obvious for both performance metrics in Table 3.3 with  $m$  increasing.

	PM <sub>10</sub> + O <sub>3</sub> + CO				PM <sub>10</sub> + O <sub>3</sub> + CO + NO <sub>2</sub>				
	RMSE		MAE		RMSE		MAE		
	TrS	TeS	TrS	TeS	TrS	TeS	TrS	TeS	
$m=1$	0.493	1.480	0.360	1.146	$m=1$	0.505	1.475	0.347	1.123
$m=2$	0.067	1.887	0.038	1.439	$m=2$	0.066	1.665	0.038	1.318
$m=3$	0.084	1.551	0.039	1.230	$m=3$	0.051	1.725	0.028	1.358
$m=4$	0.060	1.539	0.026	1.240	$m=4$	0.046	1.559	0.028	1.225
$m=5$	0.081	1.493	0.041	1.190	$m=5$	0.030	1.490	0.019	1.192
$m=6$	0.073	1.428	0.029	1.125	$m=6$	0.029	1.455	0.021	1.156
$m=7$	0.035	1.396	0.022	1.112	$m=7$	0.023	1.487	0.014	1.202
$m=8$	0.049	1.404	0.029	1.107	$m=8$	0.034	1.368	0.022	1.099
$m=9$	0.082	1.433	0.035	1.144	$m=9$	0.048	1.386	0.026	1.114
$m=10$	0.138	1.500	0.053	1.179	$m=10$	0.056	1.470	0.028	1.169
$m=11$	0.074	1.471	0.037	1.154	$m=11$	0.037	1.320	0.020	1.053
$m=12$	0.067	1.389	0.039	1.120	$m=12$	0.027	1.397	0.016	1.117

Table 3.3: Performance comparison with different air pollutants (TrS: training set, TeS: test set).

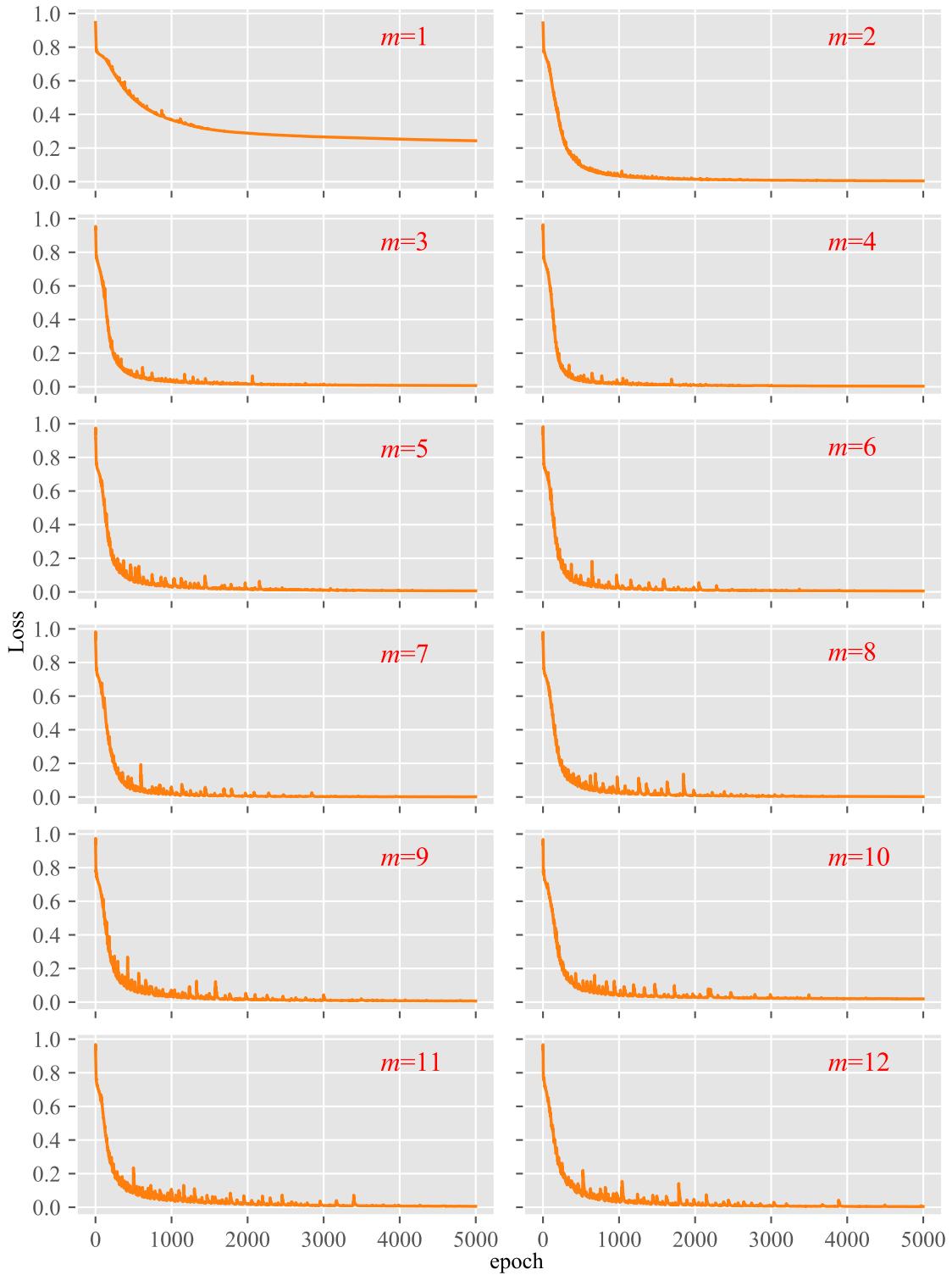


Figure 3.6: Training losses with air pollutants  $\text{PM}_{10} + \text{O}_3 + \text{CO}$ .

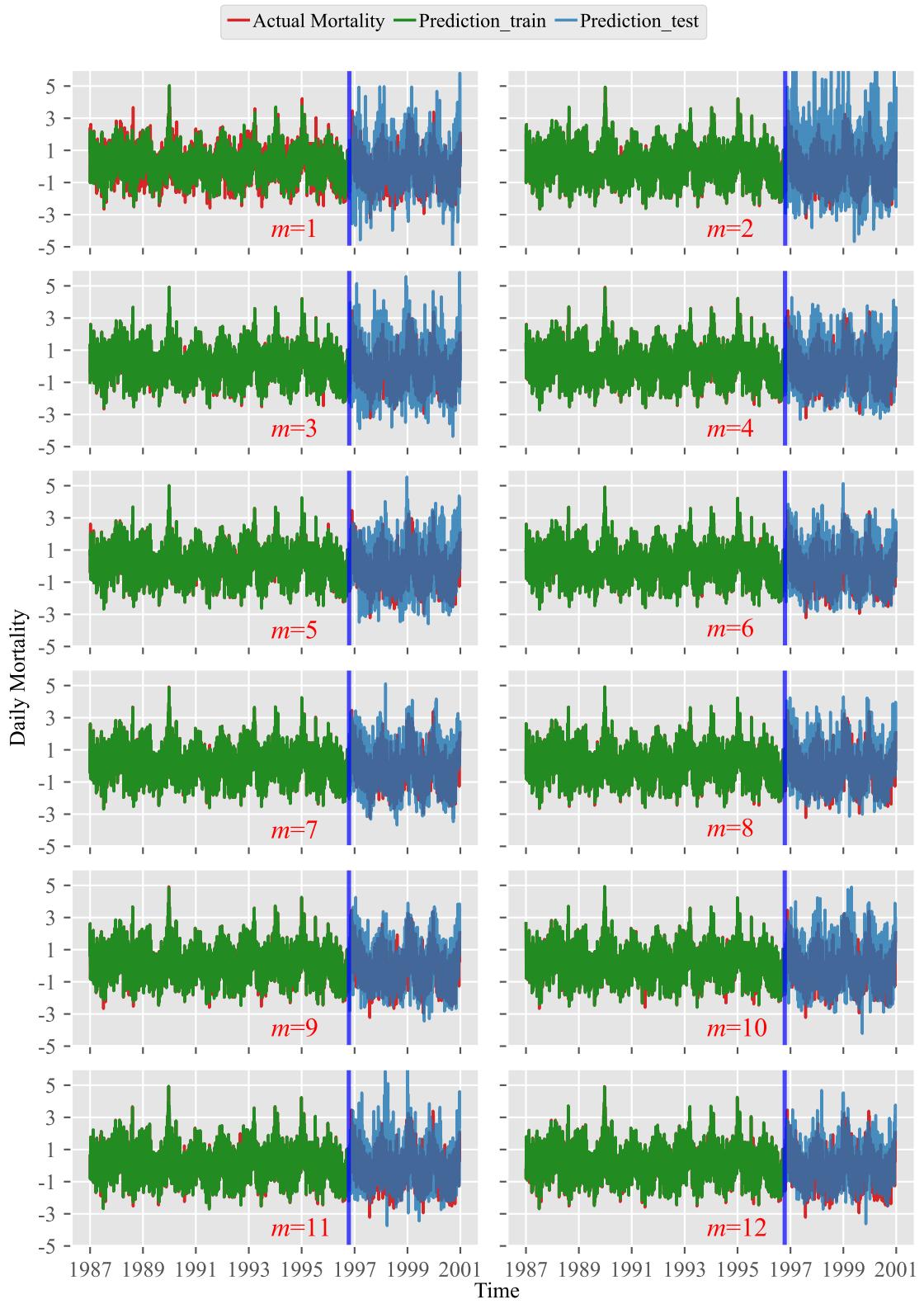


Figure 3.7: Comparison between prediction and actual daily mortality with air pollutants  $\text{PM}_{10} + \text{O}_3 + \text{CO}$ .

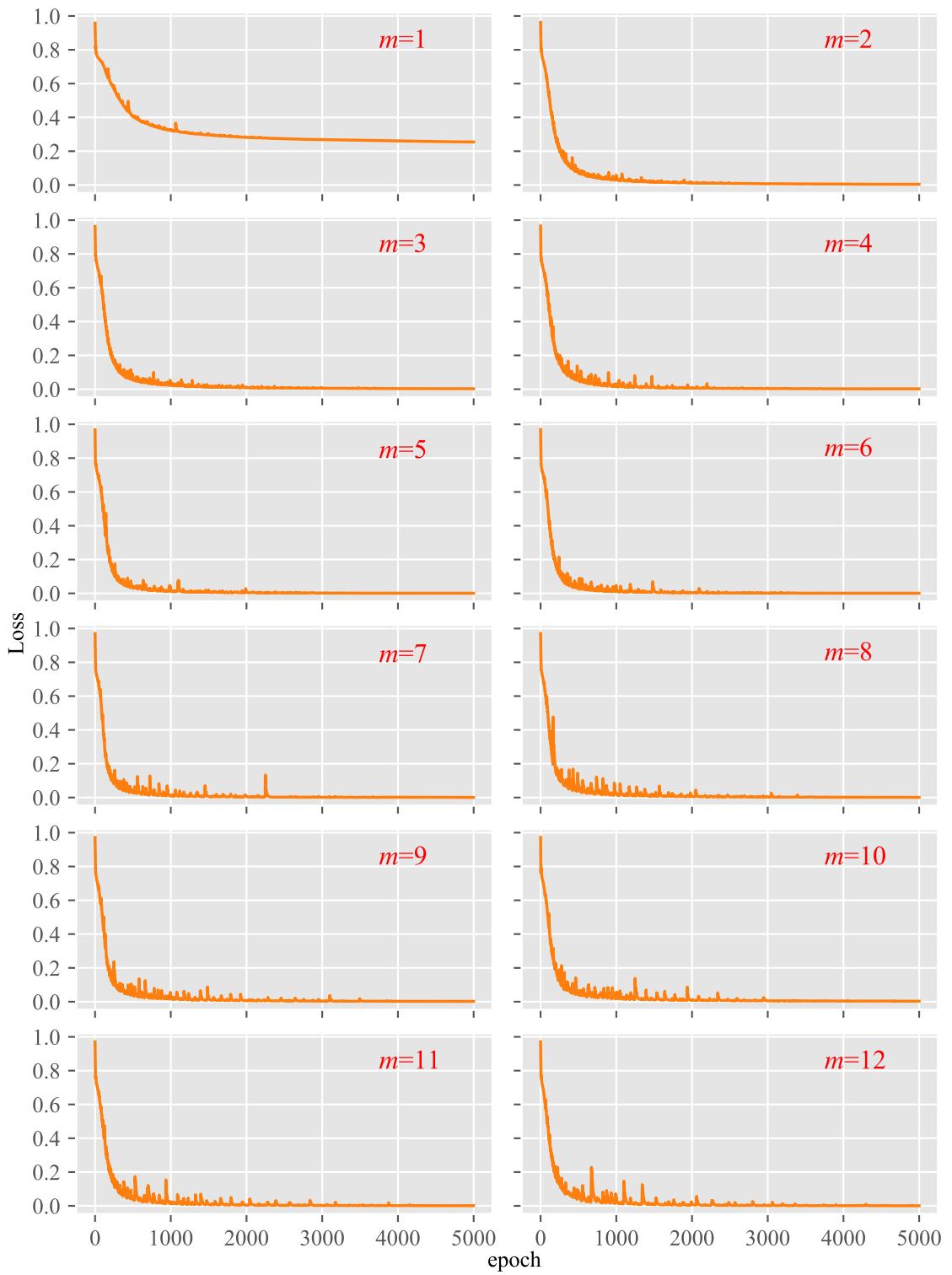


Figure 3.8: Training losses with air pollutants  $\text{PM}_{10} + \text{O}_3 + \text{CO} + \text{NO}_2$ .

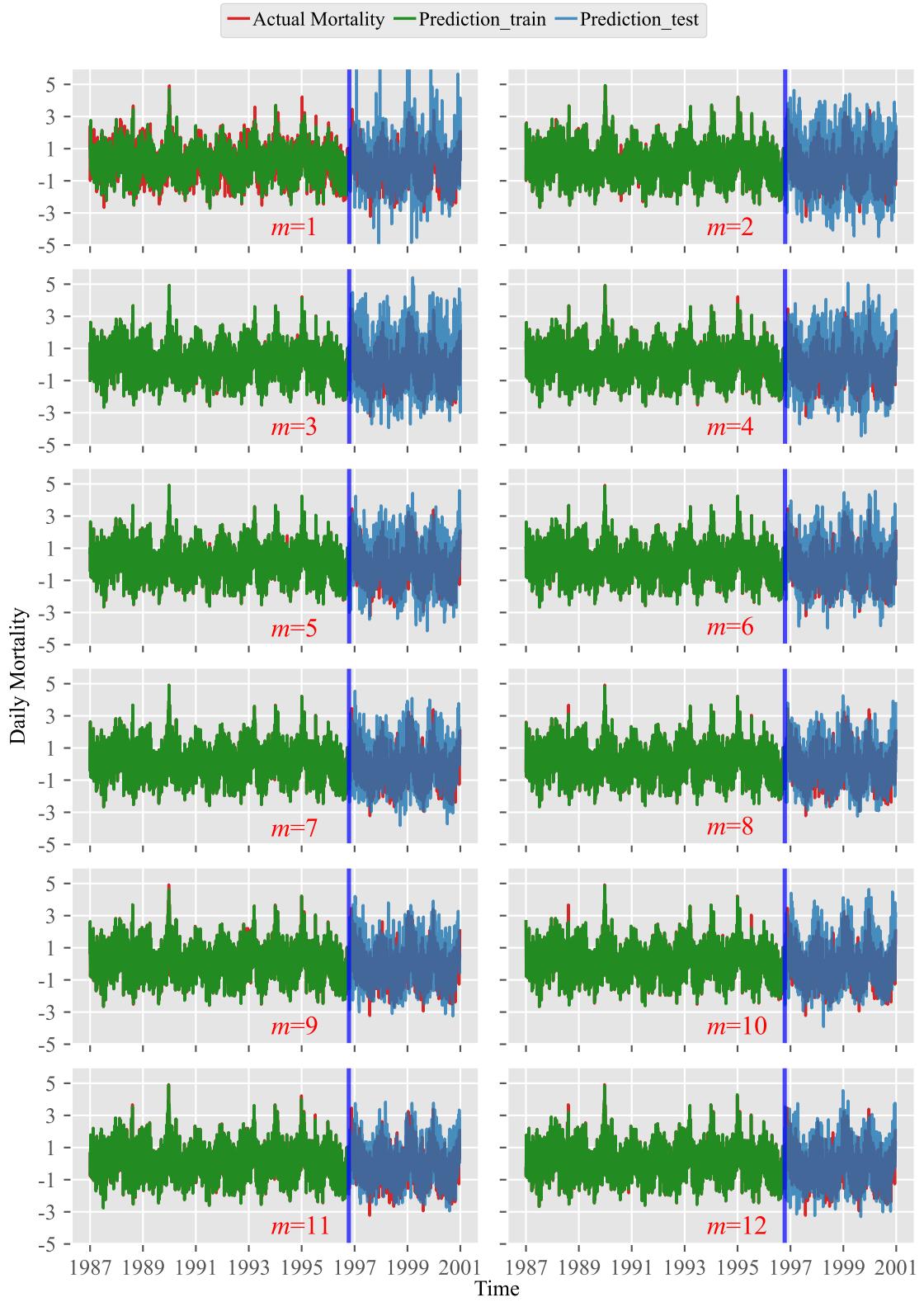


Figure 3.9: Comparison between prediction and actual daily mortality with air pollutants  $\text{PM}_{10} + \text{O}_3 + \text{CO} + \text{NO}_2$ .

## Tests for multiple air pollutants

Test results for different combinations of air pollutants and  $\mathbf{T}$  are shown in Table 3.4 and Fig. 3.10. As shown in Fig. 3.10, the proposed model can accommodate multiple air pollutants and have good convergence properties for different combinations. As shown in Table 3.4, both RMSE and MAE have roughly downward trends on the training set with the kinds of air pollutants increasing, where more types of air pollutants can provide more mortality-related information to the LSTM network for feature extraction and outcome assessment. Both metrics on the test set are not as perfect as their counterparts on the training set. The first reason is that data distribution of test set may not be exactly the same as that of the training set, which lowers the model's generalization ability. The second one is that noise from air pollutants series of the training set may lead to overfitting that lowers the models' performance on the training set. Nonetheless, the proposed model shows great potential in multipollutant-related health outcome assessment as shown by comparison between prediction and actual mortality on the right-side subfigures in Fig. 3.10.

single/multiple air pollutant (s)	RMSE		MAE	
	TrS	TeS	TrS	TeS
PM <sub>10</sub>	0.079	1.494	0.044	1.87
PM <sub>10</sub> + O <sub>3</sub>	0.081	1.460	0.047	1.157
PM <sub>10</sub> + O <sub>3</sub> + CO	0.063	1.383	0.027	1.096
PM <sub>10</sub> + O <sub>3</sub> + CO + NO <sub>2</sub>	0.042	1.393	0.021	1.113
PM <sub>10</sub> + O <sub>3</sub> + CO + NO <sub>2</sub> + SO <sub>2</sub>	0.028	1.482	0.011	1.174

Table 3.4: Comparison of performance metrics with different air pollutant (s) (TrS: training set, TeS: test set,  $m=7$ ).

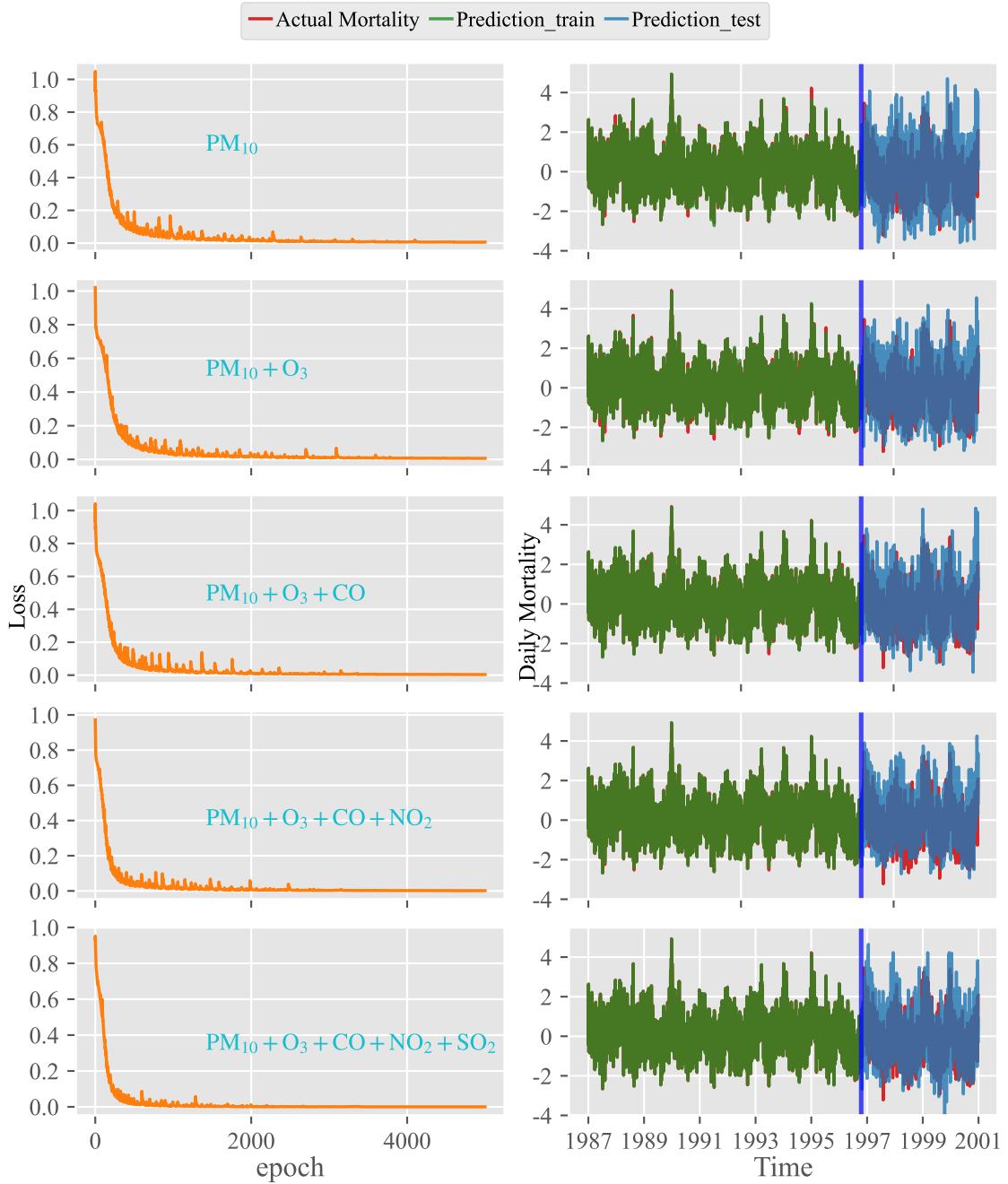


Figure 3.10: Training losses (left side) and comparison between prediction and actual mortality (right side) for single and multiple air pollutant (s) ( $m=7$ ).

### Comparison with GAM

Performance of the [LSTM](#)-based model is compared with a standard [GAM](#) used in air pollution epidemiology as Eqn. (3.14) and the results are shown in Table 3.5

and Fig. 3.11. R (version 4.3.1) is used for the experiments of the **GAM** with the package mgcv (version 1.8.42) and splines (version 4.3.1). The selected air pollutant of interest is **PM<sub>10</sub>**, which has attracted lots of attention from researchers in this field. Distributed exposure sequence from period  $t - m + 1$  to  $t$  is input to the **LSTM** model for daily mortality assessment at period  $t$  and single exposure at  $t - lag$  is used by the **GAM** for assessment. These two formulations are a distributed-lag model and a single-lag model. In addition, as **GAMs** used in air pollution-related health risk assessment is generally for retrospective study, all data are used to calculate its parameters. As shown in Table 3.5, for all four experiments both the **RMSE** and **MAE** of the **LSTM** model on the training set outperform their counterparts of the **GAM**. On the one hand, this demonstrates that the proposed **LSTM** model can better fit the health outcome, i.e., daily mortality using distributed lags than the **GAM** using only a single lag. On the other hand, it also shows that in retrospective study, the **LSTM** model can better assess the health outcome taking advantage of the impacts from distributed lags. Due to problems of data noise, data distribution differences and model overfitting as mentioned in previous experiments, performance of the **LSTM** model on the test set are not as good as on the training set and are exceeded by their counterparts of the **GAM** as in Table 3.5. However, both **RMSE** and **MAE** of the **GAM** are calculated based on its training set, i.e., the whole dataset and their counterparts of **LSTM** are calculated on the test set.

In addition, there are several points worth mentioning for this comparative experiment: 1) the **GAM** applied in air pollution epidemiology is generally used to fit

historical data and lacks the capability of predicting future health outcome with new air pollution exposure data, i.e., it is used for retrospective study instead of prospective study; 2) limited by its generalized linear formulation, the **GAM** is more suitable to evaluate the impacts of single lag and single pollutant instead of distributed lags and multiple pollutants, which is more realistic. 3) the **LSTM** model has better adaptability in accommodating distributed lags and multiple air pollutants and can synthetically evaluate the impacts from historical air pollution exposure. 4) the **LSTM** model can be applied to both retrospective and prospective study while the **GAM** is mainly used in the former one.

LSTM						GAM	
	RMSE		MAE			RMSE	MAE
	TrS	TeS	TrS	TeS		TrS	TrS
$m=2$	0.2360	1.9299	0.1354	1.4055	$lag=1$	0.8818	0.6992
$m=4$	0.0961	1.6528	0.0519	1.2745	$lag=3$	0.8823	0.6994
$m=6$	0.1495	1.4961	0.0587	1.1822	$lag=5$	0.8819	0.6992
$m=8$	0.0744	1.4096	0.0363	1.1193	$lag=7$	0.8824	0.6994

Table 3.5: Comparison performance metrics between the LSTM model and the GAM (TrS: training set, TeS: test set).

Note that in the above numerical experimentation: 1) parameters of the **LSTM** model can be further tuned for performance improvement; 2) the downward trends for both metrics may not remain with  $m$  or types of air pollutants continue increasing due to redundant information brought by associated elements; 3) either the exposure periods of interest or combinations of air pollutants of interest can be adjusted for specific applications; 4) other health outcomes like morbidity or mortality from air

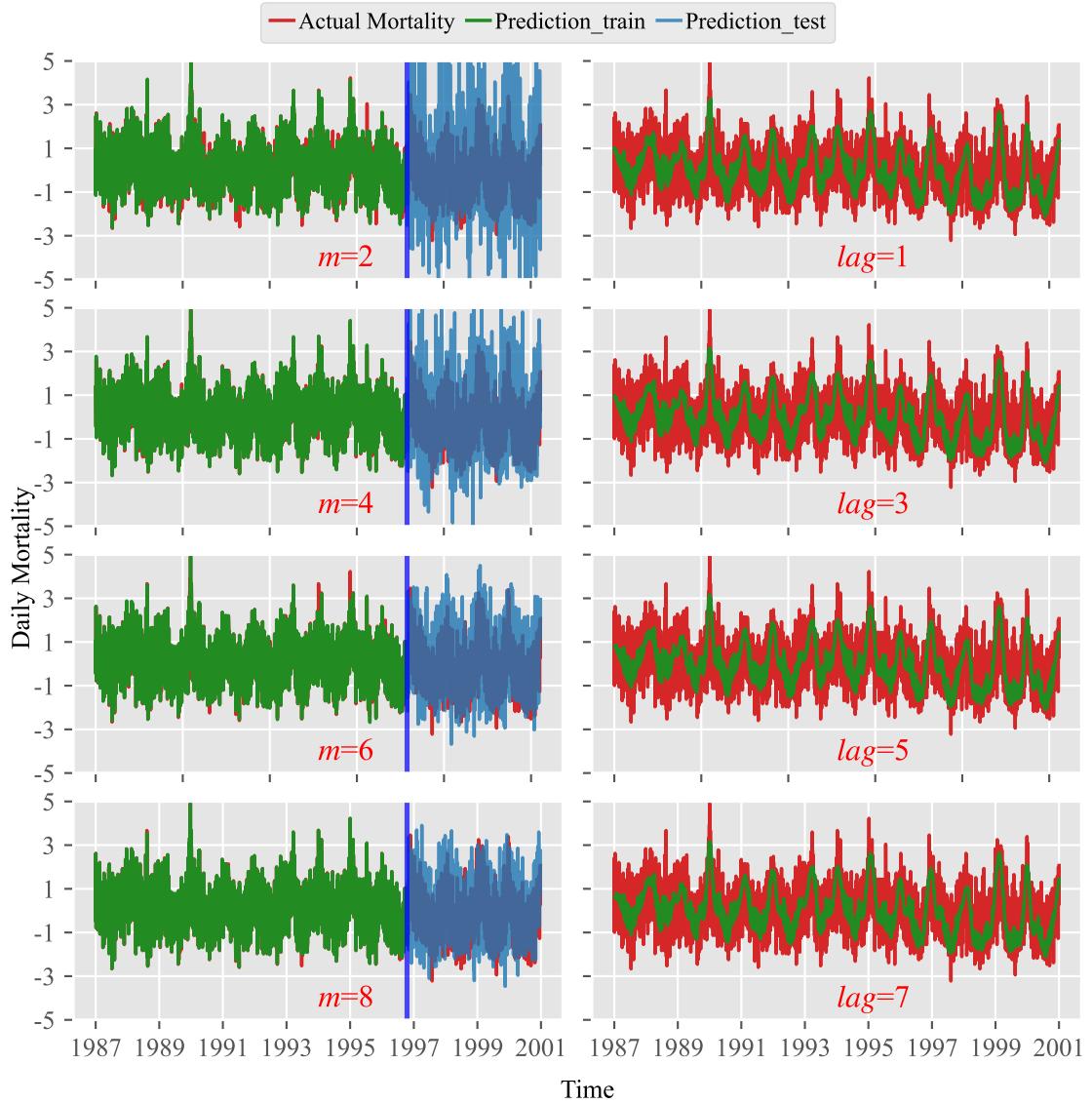


Figure 3.11: Comparison between the LSTM model and the GAM for daily mortality assessment.

pollutants of interest can also be evaluated instead.

# Chapter 4

## Conclusions

In this report, an **LSTM** neural network model is developed to assess the public health consequences from exposure to air pollution, which aims to accommodate the impacts of both distributed lags of exposure and multiple air pollutants. **MIC** is first used to evaluate the association between different air pollutants and between a specific air pollutant and health outcome of interest by means of standardizing information entropy-based **MI**. It is shown through experiments on the Toronto and Chicago datasets that the **MIC** can better capture both the linear and nonlinear relation between different elements than the **PCC**, and can be further used to select the most related air pollutant (s) for a specific health outcome. Then an **LSTM** neural network model for air pollution impacts assessment is formulated. In this formulation, an **LSTM** network is constructed to extract health outcome-related feature information from distributed exposure to a single or multiple air pollutant (s), and neural network layers with adaptive weights for the extracted features are developed to assess the

health outcome. Experiments on the Chicago datasets demonstrate that 1) the **LSTM** model is adaptable to single and multiple air pollutant (s); 2) the **LSTM** model performs well in accommodating the impacts of exposure with distributed lags for both fitting and predicting the health outcome; 3) performance of the model roughly improves with increase of the length of exposure series where more outcome-related temporal information is utilized; 4) compared to single pollutant, performance of the model with multiple air pollutants is better; 5) the proposed model outperforms **GAM** in health outcome fitting and prediction. At last, preliminary exploration of the exposure-response factor based on the **LSTM** model is made, and experiments on the **NMMAPS** datasets show that the proposed model brings a new viewpoint to and exhibits great potential in air pollution-related epidemiology.

Although the association between different elements can be evaluated comprehensively with the **MIC** and the proposed **LSTM** model shows good performance in air pollution-related health consequence assessment, there also exist limitations in their applications. These problems and several promising research opportunities that deserve attention are further discussed as follows.

- 1) **Redundant information.** For pairwise elements, the **MIC** can effectively capture their dependencies. However, for multiple air pollutants selected for some health outcome using the **MIC**, they may actually not be the most related factors as a whole due to the redundant health outcome-related information among them. This redundancy may weaken the strength of the selected air pollutants for outcome assessing and further raises the question of how to select the most

related multi-pollutant. A promising research direction would be to select the most related candidate air pollutants for the health outcome first, and then select the least associated candidates therein.

- 2) **Model refinement.** Main structure parameters including number of hidden layers and output features of the **LSTM** network as well as the batch size and type of the optimizer keep unchanged in previous experimentation. These hyper parameters can be further tuned for each experiment to improve the model performance. Possible parameter tuning methods include but not limited to grid search [99], randomized search [100], Bayesian optimization [101], etc.
- 3) **Distributed lags.** Historical air pollution exposure has important lagged impacts on the health outcome. Longer-term exposure sequence can provide more outcome-related features. However, with the input sequence length keep increasing, more data noise is also brought into the **LSTM** model, with which the model performance may be weakened. In addition, incorporating more distributed lags increases the computational complexity and challenges the model's capability in capturing long-term dependencies. Thus, optimizing the length of input air pollution exposure sequence is a promising research direction to facilitate the application of the model to air pollution epidemiology and the methods mentioned for hyper parameter tuning may be referenced.
- 4) **Pooled and (or) meta analysis.** Based on the application of the **LSTM** model in individual studies, data or results from multiple areas can be further combined for collaborative and synthetic analyses. Either retrospective or prospective anal-

yses can be made to provide overall assessments and more reliable results.

# Bibliography

- [1] Yun-Chul Hong, Jong-Tae Lee, Ho Kim, and Ho-Jang Kwon. Air pollution: a new risk factor in ischemic stroke mortality. *Stroke*, 33(9):2165–2169, 2002.
- [2] Ken Donaldson, Nicholas Mills, William MacNee, Simon Robinson, and David Newby. Role of inflammation in cardiopulmonary health effects of PM. *Toxicology and Applied Pharmacology*, 207(2):483–488, 2005.
- [3] Gerard Hoek, Ranjini M Krishnan, Rob Beelen, Annette Peters, Bart Ostro, Bert Brunekreef, and Joel D Kaufman. Long-term air pollution exposure and cardio-respiratory mortality: a review. *Environmental Health*, 12:1–16, 2013.
- [4] Jamie I Verhoeven, Youssra Allach, Ilonca CH Vaartjes, Catharina JM Klijn, and Frank-Erik de Leeuw. Ambient air pollution and the risk of ischaemic and haemorrhagic stroke. *The Lancet Planetary Health*, 5(8):542–552, 2021.
- [5] Joel Schwartz. Particulate air pollution and chronic respiratory disease. *Environmental Research*, 62(1):7–13, 1993.
- [6] C Arden Pope Iii, Richard T Burnett, Michael J Thun, Eugenia E Calle, Daniel Krewski, Kazuhiko Ito, and George D Thurston. Lung cancer, cardiopulmonary mortality, and long-term exposure to fine particulate air pollution. *The Journal of the American Medical Association*, 287(9):1132–1141, 2002.
- [7] Frank J Kelly. Influence of air pollution on respiratory disease. *European Medical Journal Respiratory*, 2:96–103, 2014.

- [8] Douglas W Dockery, C Arden Pope, Xiping Xu, John D Spengler, James H Ware, Martha E Fay, Benjamin G Ferris Jr, and Frank E Speizer. An association between air pollution and mortality in six US cities. *New England Journal of Medicine*, 329(24):1753–1759, 1993.
- [9] Jonathan M Samet, Francesca Dominici, Frank C Curriero, Ivan Coursac, and Scott L Zeger. Fine particulate air pollution and mortality in 20 US cities, 1987–1994. *New England Journal of Medicine*, 343(24):1742–1749, 2000.
- [10] Paul H Fischer, Marten Marra, Caroline B Ameling, Gerard Hoek, Rob Beele, Kees de Hoogh, Oscar Breugelmans, Hanneke Kruize, Nicole AH Janssen, and Danny Houthuijs. Air pollution and mortality in seven million adults: the dutch environmental longitudinal study (duels). *Environmental Health Perspectives*, 123(7):697–704, 2015.
- [11] Jonathan M Samet, Francesca Dominici, Scott L Zeger, Joel Schwartz, and Douglas W Dockery. The National Morbidity, Mortality, and Air Pollution Study. Part I: Methods and Methodologic issues. *Research Report (Health Effects Institute)*, 94(pt 1):5–14, 2000.
- [12] Jonathan M Samet, Scott L Zeger, Francesca Dominici, Frank Curriero, Ivan Coursac, Douglas W Dockery, Joel Schwartz, and Antonella Zanobetti. The National Morbidity, Mortality, and Air Pollution Study. Part II: Morbidity and Mortality From Air Pollution in the United States. *Research Report (Health Effects Institute)*, 94(pt 2):5–79, 2000.
- [13] Klea Katsouyanni, Joel Schwartz, Claudia Spix, Giota Touloumi, Denis Zmirou, Antonella Zanobetti, Bogdan Wojtyniak, JM Vonk, A Tobias, A Pönkä, S Medina, L Bachárová, and H R Anderson. Short term effects of air pollution on health: a European approach using epidemiologic time series data: the APHEA protocol. *Journal of Epidemiology and Community Health*, 50(1):12–18, 1996.

- [14] Alexandros Gryparis, Bertil Forsberg, Klea Katsouyanni, Antonis Analitis, Giota Touloumi, Joel Schwartz, Evangelia Samoli, Sylvia Medina, HR Anderson, Emilia Maria Niciu, HE Wichmann, Bohumir Kriz, Mitja Kosnik, Jiri Skorkovsky, JM Vonk, and Zeynep Dörbudak. Acute effects of ozone on mortality from the “Air Pollution and Health: A European Approach” project. *American Journal of Respiratory and Critical Care Medicine*, 170(10):1080–1087, 2004.
- [15] Klea Katsouyanni and APHEA Group. APHEA project: Air pollution and health: A European approach. *Epidemiology*, 17(6):S19, 2006.
- [16] Danford G Kelley. The National Air Pollution Surveillance Network in Canada. *Journal of the Air Pollution Control Association*, 29(8):794–795, 1979.
- [17] Kenneth L Demerjian. A review of national monitoring networks in North America. *Atmospheric Environment*, 34(12-14):1861–1884, 2000.
- [18] Francesca Dominici, Jonathan M Samet, and Scott L Zeger. Combining evidence on air pollution and daily mortality from the 20 largest US cities: a hierarchical modelling strategy. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 163(3):263–302, 2000.
- [19] Francesca Dominici, Aidan McDermott, Scott L Zeger, and Jonathan M Samet. On the use of generalized additive models in time-series studies of air pollution and health. *American Journal of Epidemiology*, 156(3):193–203, 2002.
- [20] Richard B Davies. From cross-sectional to longitudinal analysis. *Analyzing Social and Political Change: A Casebook of Methods*, pages 20–40, 1994.
- [21] Thomas Götschi, Joachim Heinrich, Jordi Sunyer, and Nino Künzli. Long-term effects of ambient air pollution on lung function: a review. *Epidemiology*, 19(5):690–701, 2008.

- [22] Antonio Ciocco and Donovan J Thompson. A follow-up of Donora ten years after: methodology and findings. *American Journal of Public Health and the Nations Health*, 51(2):155–164, 1961.
- [23] Lester B Lave and Eugene P Seskin. Air pollution and human health: The quantitative effect, with an estimate of the dollar benefit of pollution abatement, is considered. *Science*, 169(3947):723–733, 1970.
- [24] Frederick W Lipfert. Statistical studies of mortality and air pollution: Multiple regression analyses stratified by age group. *Science of the Total Environment*, 15(2):103–122, 1980.
- [25] Diane I Gibbons and Gary C McDonald. Illustrating regression diagnostics with an air pollution and mortality model. *Computational Statistics and Data Analysis*, 1:201–220, 1983.
- [26] Bart D Ostro. The effects of air pollution on work loss and morbidity. *Journal of Environmental Economics and Management*, 10(4):371–382, 1983.
- [27] RE Waller, AG Brooks, and MW Adler. The 1952 fog cohort study. *British Journal of Preventive and Social Medicine*, 27(1):68, 1973.
- [28] Anthony T Kerigan, Charles H Goldsmith, and L David Pengelly. A three-Year cohort study of the role of environmental factors in the respiratory health of children in Hamilton, Ontario: epidemiologic survey design, methods, and description of cohort. *American Review of Respiratory Disease*, 133(6):987–993, 1986.
- [29] Frank E Speizer. Assessment of the epidemiological data relating lung cancer to air pollution. *Environmental Health Perspectives*, 47:33–42, 1983.
- [30] Göran Pershagen and Lorenzo Simonato. Epidemiological evidence on air pollution and cancer. In *Air Pollution and Human Cancer*, pages 63–74. Springer, 1990.

- [31] Kari Hemminki and Göran Pershagen. Cancer risk of air pollution: epidemiological evidence. *Environmental Health Perspectives*, 102(suppl 4):187–192, 1994.
- [32] Beate Ritz, Fei Yu, Guadalupe Chapa, and Scott Fruin. Effect of air pollution on preterm birth among children born in Southern California between 1989 and 1993, 2000.
- [33] Sung Kyun Park, Marie S O'Neill, Pantel S Vokonas, David Sparrow, and Joel Schwartz. Effects of air pollution on heart rate variability: the VA normative aging study. *Environmental Health Perspectives*, 113(3):304–309, 2005.
- [34] Lie Hong Chen, Synnove F Knutsen, David Shavlik, W Lawrence Beeson, Floyd Petersen, Mark Ghamsary, and David Abbey. The association between fatal coronary heart disease and ambient particulate air pollution: are females at greater risk? *Environmental Health Perspectives*, 113(12):1723–1729, 2005.
- [35] Justin Barclay, Graham Hillis, and Jon Ayres. Air pollution and the heart: cardiovascular effects and mechanisms. *Toxicological Reviews*, 24:115–123, 2005.
- [36] Pierre R Band, Nhu D Le, Raymond Fang, Michele Deschamps, Andrew J Coldman, Richard P Gallagher, and Joanne Moody. Cohort study of Air Canada pilots: mortality, cancer incidence, and leukemia risk. *American Journal of Epidemiology*, 143(2):137–143, 1996.
- [37] David E Abbey, Naomi Nishino, William F McDonnell, Raoul J Burchette, Synnøve F Knutsen, W Lawrence Beeson, and Jie X Yang. Long-term inhalable particles and other air pollutants related to mortality in nonsmokers. *American Journal of Respiratory and Critical Care Medicine*, 159(2):373–382, 1999.
- [38] Michael Jerrett, Richard Burnett, Renjun Ma, Bruce Newbold, George Thurston, and Daniel Krewski. A cohort study of air pollution and mortality in Los Angeles. *Epidemiology*, 15(4):S46, 2004.

- [39] Edward L Korn and Alice S Whittemore. Methods for analyzing panel studies of acute health effects of air pollution. *Biometrics*, pages 795–802, 1979.
- [40] Alice S Whittemore and Edward L Korn. Asthma and air pollution in the Los Angeles area. *American Journal of Public Health*, 70(7):687–696, 1980.
- [41] Justin L Barclay, Brian G Miller, Smita Dick, Martine Dennekamp, Isobel Ford, Graham S Hillis, Jon G Ayres, and Anthony Seaton. A panel study of air pollution in subjects with heart failure: negative results in treated patients. *Occupational and Environmental Medicine*, 66(5):325–334, 2009.
- [42] Ubiratan de Paula Santos, Alfésio Luís Ferreira Braga, Dante Marcelo Artigas Giorgi, Luiz Alberto Amador Pereira, César Jose Grupi, Chin An Lin, Marcos Antonio Bussacos, Dirce Maria Trevisan Zanetta, Paulo Hilário do Nascimento Saldiva, and Mario Terra Filho. Effects of air pollution on blood pressure and heart rate variability: a panel study of vehicular traffic controllers in the city of Sao Paulo, Brazil. *European Heart Journal*, 26(2):193–200, 2005.
- [43] Katharina Hildebrandt, Regina Rückerl, Wolfgang Koenig, Alexandra Schneider, Mike Pitz, Joachim Heinrich, Victor Marder, Mark Frampton, Günter Oberdörster, H Erich Wichmann, et al. Short-term effects of air pollution: a panel study of blood markers in patients with chronic pulmonary disease. *Particle and Fibre Toxicology*, 6:1–13, 2009.
- [44] Nelly D Saenen, Eline B Provost, Mineke K Viaene, Charlotte Vanpoucke, Wouter Lefebvre, Karen Vrijens, Harry A Roels, and Tim S Nawrot. Recent versus chronic exposure to particulate matter air pollution in association with neurobehavioral performance in a panel study of primary schoolchildren. *Environment International*, 95:112–119, 2016.
- [45] Susanna Lagorio, Francesco Forastiere, Riccardo Pistelli, Ivano Iavarone, Paola Michelozzi, Valeria Fano, Achille Marconi, Giovanni Ziemacki, and Bart D

- Ostro. Air pollution and lung function among susceptible adult subjects: a panel study. *Environmental Health*, 5:1–12, 2006.
- [46] Howard Kipen, David Rich, Wei Huang, Tong Zhu, Guangfa Wang, Min Hu, Shou-en Lu, Pamela Ohman-Strickland, Ping Zhu, Yuedan Wang, et al. Measurement of inflammation and oxidative stress following drastic changes in air pollution during the Beijing Olympics: a panel study approach. *Annals of the New York Academy of Sciences*, 1203(1):160–167, 2010.
- [47] JJK Jaakkola. Case-crossover design in air pollution epidemiology. *European Respiratory Journal*, 21(40 suppl):81s–85s, 2003.
- [48] Malcolm MacLure. The case-crossover design: a method for studying transient effects on the risk of acute events. *American Journal of Epidemiology*, 133(2):144–153, 1991.
- [49] Lucas M Neas, Joel Schwartz, and Douglas Dockery. A case-crossover analysis of air pollution and mortality in Philadelphia. *Environmental Health Perspectives*, 107(8):629–631, 1999.
- [50] Jong-Tae Lee and Joel Schwartz. Reanalysis of the effects of air pollution on daily mortality in Seoul, Korea: a case-crossover design. *Environmental Health Perspectives*, 107(8):633–636, 1999.
- [51] Harvey Checkoway, Drew Levy, Lianne Sheppard, Joel Kaufman, Jane Koenig, and David Siscovick. A case-crossover analysis of fine particulate matter air pollution and out-of-hospital sudden cardiac arrest. *Research Report-Health Effects Institute*, 2000.
- [52] Hyo Young June Im, Sang Yun Lee, Ki Jung Yun, Young Su Ju, Dae Hee Kang, and Soo Hon Cho. A case-crossover study between air pollution and hospital emergency room visits by asthma attack. *Korean Journal of Occupational and Environmental Medicine*, 12(2):249–257, 2000.

- [53] Daniela D'Ippoliti, Francesco Forastiere, Carla Ancona, Nera Agabiti, Danilo Fusco, Paola Michelozzi, and Carlo A Perucci. Air pollution and myocardial infarction in Rome: a case-crossover analysis. *Epidemiology*, 14(5):528–535, 2003.
- [54] Francesco Forastiere, Massimo Stafoggia, Sally Picciotto, Tom Bellander, Daniela D'Ippoliti, Timo Lanki, Stephanie von Klot, Fredrik Nyberg, Pentti Paatero, Annette Peters, Juha Pekkanen, Jordi Sunyer, and CA Perucci. A case-crossover analysis of out-of-hospital coronary deaths and air pollution in Rome, Italy. *American Journal of Respiratory and Critical Care Medicine*, 172(12):1549–1555, 2005.
- [55] Roger Zemek, Mieczysław Szyszkowicz, and Brian H Rowe. Air pollution and emergency department visits for otitis media: a case-crossover study in Edmonton, Canada. *Environmental Health Perspectives*, 118(11):1631–1636, 2010.
- [56] Takashi Yorifuji, Etsushi Suzuki, and Saori Kashima. Hourly differences in air pollution and risk of respiratory disease in the elderly: a time-stratified case-crossover study. *Environmental Health*, 13:1–11, 2014.
- [57] Cheryl S Pirozzi, Barbara E Jones, James A VanDerslice, Yue Zhang, Robert Paine III, and Nathan C Dean. Short-term air pollution and incident pneumonia. a case-crossover study. *Annals of the American Thoracic Society*, 15(4):449–459, 2018.
- [58] Jerry A Hausman, Bart D Ostro, and David A Wise. Air pollution and lost work. *NBER Working Paper*, (w1263), 1984.
- [59] Bart D Ostro and Susy Rothschild. Air pollution and acute respiratory morbidity: an observational study of multiple pollutants. *Environmental Research*, 50(2):238–247, 1989.
- [60] Bart David Ostro. Estimating the risks of smoking, air pollution, and passive smoke on acute respiratory conditions. *Risk Analysis*, 9(2):189–196, 1989.

- [61] Toshiro Tango. Effect of air pollution on lung cancer: a Poisson regression model based on vital statistics. *Environmental Health Perspectives*, 102(suppl 8):41–45, 1994.
- [62] Joel Schwartz. Particulate air pollution and daily mortality in Detroit. *Environmental Research*, 56(2):204–213, 1991.
- [63] Arnoud P Verhoeff, Gerard Hoek, Joel Schwartz, and Joop H van Wijnen. Air pollution and daily mortality in Amsterdam. *Epidemiology*, 7(3):225–230, 1996.
- [64] Lianne Sheppard and Doris Damian. Estimating short-term pm effects accounting for surrogate exposure measurements from ambient monitors. *Environmetrics: The official journal of the International Environmetrics Society*, 11(6):675–687, 2000.
- [65] Claudia Spix, H Ross Anderson, Joel Schwartz, Maria Angela Vigotti, Alain Letertre, Judith M Vonk, Giota Touloumi, Franck Balducci, Tomasz Piekarski, Ljuba Bacharova, Aurelio Tobias, antti Pönkä, and Klea Katsouyanni. Short-term effects of air pollution on hospital admissions of respiratory diseases in Europe: a quantitative summary of APHEA study results. *Archives of Environmental Health: An International Journal*, 53(1):54–64, 1998.
- [66] Klea Katsouyanni, Giota Touloumi, Evangelia Samoli, Alexandros Gryparis, Alain Le Tertre, Yannis Monopolis, Giuseppe Rossi, Denis Zmirou, Ferran Ballester, Azedine Boumghar, Hugh Ross Anderson, Bogdan Wojtyniak, Anna Paldy, Rony Braunstein, Juha Pekkanen, Christian11 Schindler, and Joel Schwartz. Confounding and effect modification in the short-term effects of ambient particles on total mortality: results from 29 European cities within the APHEA2 project. *Epidemiology*, 12(5):521–531, 2001.
- [67] Ariana Zeka and Joel Schwartz. Estimating the independent effects of multiple pollutants in the presence of measurement error: an application of a

- measurement-error-resistant technique. *Environmental Health Perspectives*, 112(17):1686–1690, 2004.
- [68] Antonella Zanobetti, Matt P Wand, Joel Schwartz, and Louise M Ryan. Generalized additive distributed lag models: quantifying mortality displacement. *Biostatistics*, 1(3):279–292, 2000.
- [69] Francesca Dominici, Aidan McDermott, Scott L Zeger, and Jonathan M Samet. Airborne particulate matter and mortality: timescale effects in four US cities. *American Journal of Epidemiology*, 157(12):1055–1065, 2003.
- [70] Wesley S Burr, Hwashin H Shin, and Glen Takahara. Synthetically lagged models. *Statistics and Probability Letters*, 144:37–43, 2019.
- [71] Guowen Huang, Marta Blangiardo, Patrick E Brown, and Monica Pirani. Long-term exposure to air pollution and COVID-19 incidence: a multi-country study. *Spatial and Spatio-temporal Epidemiology*, 39:100443, 2021.
- [72] Wesley S Burr, Glen Takahara, and Hwashin H Shin. Bias correction in estimation of public health risk attributable to short-term air pollution exposure. *Environmetrics*, 26(4):298–311, 2015.
- [73] Dimitris Evangelopoulos, Klea Katsouyanni, Joel Schwartz, and Heather Walton. Quantifying the short-term effects of air pollution on health in the presence of exposure measurement error: a simulation study of multi-pollutant model results. *Environmental Health*, 20:1–13, 2021.
- [74] Francesca Dominici, Antonella Zanobetti, Joel Schwartz, Danielle Braun, Ben Sabath, and Xiao Wu. Assessing adverse health effects of long-term exposure to low levels of ambient air pollution: implementation of causal inference methods. *Research Reports: Health Effects Institute*, 2022, 2022.
- [75] Hwashin Hyun Shin, Rajendra Prasad Parajuli, Aubrey Maquiling, and Marc Smith-Doiron. Temporal trends in associations between ozone and circulatory

- mortality in age and sex in Canada during 1984–2012. *Science of the Total Environment*, 724:137944, 2020.
- [76] Government of Canada Open Data Portal. National Air Pollution Surveillance Program. <https://data-donnees.az.ec.gc.ca/data/air/monitor/national-air-pollution-surveillance-naps-program/?lang=en>, 2024.
- [77] Government of Canada Open Data Portal. National Air Pollution Surveillance Program. <https://open.canada.ca/data/en/dataset/1b36a356-defd-4813-acea-47bc3abd859b>, 2024.
- [78] Government of Canada Open Data Portal. Historical climate data. <https://climate.weather.gc.ca/>, 2024.
- [79] Roger D Peng and Francesca Dominici. Statistical methods for environmental epidemiology with R. *R: a case study in air pollution and health*, 2008.
- [80] Roger D Peng and Leah J Welty. The nmmapsdata package. *R news*, 4(2):10–14, 2004.
- [81] National Weather Service. National Oceanic and Atmospheric Administration Online Weather Data. <https://www.weather.gov/wrh/Climate?wfo=lot>, 2024.
- [82] Shannon Jarvis. Particulate matter component analyses in relation to public health in Canada. Master’s thesis, Trent University (Canada), 2023.
- [83] Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 69(6):066138, 2004.
- [84] Claude Elwood Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, 1948.
- [85] David N Reshef, Yakir A Reshef, Hilary K Finucane, Sharon R Grossman, Gilean McVean, Peter J Turnbaugh, Eric S Lander, Michael Mitzenmacher,

- and Pardis C Sabeti. Detecting novel associations in large data sets. *Science*, 334(6062):1518–1524, 2011.
- [86] David N Reshef, Yakir A Reshef, Michael Mitzenmacher, and Pardis C Sabeti. Cleaning up the record on the maximal information coefficient and equitability. *Proceedings of the National Academy of Sciences*, 111(33):E3362–E3363, 2014.
- [87] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020.
- [88] Davide Albanese, Michele Filosi, Roberto Visintainer, Samantha Riccadonna, Giuseppe Jurman, and Cesare Furlanello. Minerva and minepy: a C engine for the MINE suite and its R, Python and MATLAB wrappers. *Bioinformatics*, 29(3):407–408, 2013.
- [89] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [90] MI Jordan. Serial order: a parallel distributed processing approach. technical report, june 1985-march 1986. Technical report, California Univ., San Diego, La Jolla (USA). Inst. for Cognitive Science, 1986.
- [91] Jeffrey L Elman. Finding structure in time. *Cognitive Science*, 14(2):179–211, 1990.

- [92] John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, 1982.
- [93] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, 1994.
- [94] Barak A Pearlmutter. Gradient calculations for dynamic recurrent neural networks: A survey. *IEEE Transactions on Neural networks*, 6(5):1212–1228, 1995.
- [95] Kanchan M Tarwani and Swathi Edem. Survey on recurrent neural network in natural language processing. *International Journal of Engineering Trends and Technology*, 48(6):301–304, 2017.
- [96] Alex Graves, Navdeep Jaitly, and Abdel-rahman Mohamed. Hybrid speech recognition with deep bidirectional LSTM. In *2013 IEEE Workshop on Automatic Speech Recognition and Understanding*, pages 273–278. IEEE, 2013.
- [97] Shashi Pal Singh, Ajai Kumar, Hemant Darbari, Lenali Singh, Anshika Rastogi, and Shikha Jain. Machine translation using deep learning: An overview. In *2017 International Conference on Computer, Communications and Electronics (comptelix)*, pages 162–167. IEEE, 2017.
- [98] Brian S Freeman, Graham Taylor, Bahram Gharabaghi, and Jesse Thé. Forecasting air quality time series using deep learning. *Journal of the Air and Waste Management Association*, 68(8):866–886, 2018.
- [99] Hussain Alibrahim and Simone A Ludwig. Hyperparameter optimization: Comparing genetic algorithm against grid search and bayesian optimization. In *2021 IEEE Congress on Evolutionary Computation (CEC)*, pages 1551–1559. IEEE, 2021.

- [100] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(2), 2012.
- [101] A Helen Victoria and Ganesh Maragatham. Automatic tuning of hyperparameters using bayesian optimization. *Evolving Systems*, 12(1):217–223, 2021.

# APPENDICES

## Functions Used in LSTM network

The sigmoid, softmax and tanh functions used in the [LSTM](#) model is as Eqns. (1-3).

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

$$SF(i) = \frac{e^i}{\sum_j e^j} \quad (2)$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (3)$$