

# Investigation of NO<sub>x</sub> (in ppb), O<sub>3</sub> (in ppb), and PM<sub>2.5</sub> (in µg/m<sup>3</sup>) trends in Ontario from 2003 to 2022

Mohamad Damaj - 10109, Moham An, Jo Zhu

2025-03-18

## 1. Description of Datasets

The datasets analyzed in this report are three, one for each of NO<sub>x</sub> (in ppb), PM<sub>2.5</sub> (in µg/m<sup>3</sup>), and O<sub>3</sub> (in ppb), all containing the same variables. These datasets contain records of the pollutant levels from 2003 (considered the earliest of all three) to 2022. Each dataset contains detailed information of pollutant levels in the form of the following variables:

|    |      |                   |                   |                   |                   |
|----|------|-------------------|-------------------|-------------------|-------------------|
| ## | [1]  | "Year"            | "Station Number"  | "City"            | "Location"        |
| ## | [5]  | "Type"            | "Valid Hour"      | "10th Percentile" | "30th Percentile" |
| ## | [9]  | "50th Percentile" | "70th Percentile" | "90th Percentile" | "99 Percentile"   |
| ## | [13] | "Mean"            | "1-Hour Maximum"  | "24-Hour Maximum" |                   |

1. Year: The calendar year when the pollutant data was recorded.
2. Station Number: A unique numerical identifier for the monitoring station.
3. City: The name of the city where the pollutant was measured.
4. Location: A more specific description of the monitoring location (street, area).
5. Type: The type of monitoring station or measurement protocol used.
6. Valid Hour: The number of valid hourly measurements recorded for that pollutant in the given year.
7. Percentiles: Percentile statistics showing the distribution of pollutant levels (in ppb, µg/m<sup>3</sup>, ppb) throughout the year.
8. Mean: The annual average concentration of the pollutant (in ppb, µg/m<sup>3</sup>, ppb).
9. 1-Hour Maximum: The highest recorded concentration within a one-hour period during the year.
10. 24-Hour Maximum: The highest average concentration recorded over a 24-hour period in that year.

## 2. Background of the Data

The data used in this analysis was collected across various cities in Ontario, throughout the period from 2003 to 2022. These measurements were compiled and made available by the Ontario Ministry of the Environment, Conservation and Parks.

The information is publically available and is helpful in assisting researchers, policymakers, and public health officials to investigate air quality trends, evaluate the impact of environmental regulations, and inform policy decisions. It provides a long-term view of pollutant levels in various cities in Ontario, allowing for improving environmental and public health outcomes in Ontario.

### 3. Overall Research Question

The aim of this paper is to investigate the spatial and temporal trends of air pollutants ( $\text{NO}_x$ ,  $\text{O}_3$ , and  $\text{PM}_{2.5}$ ) across various Ontario cities from 2003 to 2022.

1. How have annual mean concentrations (in ppb,  $\mu\text{g}/\text{m}^3$ , ppb) for each pollutant changed from 2003 to 2022, and which years show the most significant shifts?
2. Which cities consistently rank among the highest (or lowest) in terms of average pollutant levels, and do these rankings shift over time?
3. How do the four pollutants correlate with each other across different cities and years, and what might this indicate about broader air quality patterns?
4. How do the four pollutants project into future years based on current trends?

### 4. Tables

#### 4.1 The Top 10 Cities With the Highest Pollutant Concentration

| NO <sub>x</sub> Table (in ppb) |       | O <sub>3</sub> Table (in ppb) |       | PM <sub>2.5</sub> Table (in $\mu\text{g}/\text{m}^3$ ) |       |
|--------------------------------|-------|-------------------------------|-------|--|-------|
| City                           | Conc. | City                          | Conc. | City   | Conc. |
| Toronto West                   | 31.32 | Port Stanley                  | 32.84 | Sarnia   | 9.57  |
| Toronto East                   | 21.82 | Tiverton                      | 31.93 | Windsor West   | 9.03  |
| Toronto Downtown               | 20.41 | Grand Bend                    | 31.28 | Hamilton Downtown                                      | 8.91  |
| Toronto North                  | 19.67 | Parry Sound                   | 30.56 | Etobicoke West   | 8.44  |
| Hamilton Downtown              | 19.36 | Chatham                       | 29.88 | Windsor Downtown                                       | 8.36  |
| Windsor Downtown               | 18.00 | Belleville                    | 29.63 | Hamilton West  | 8.06  |
| Burlington                     | 17.74 | Kingston                      | 29.61 | Hamilton Mountain                                      | 7.91  |
| Hamilton West                  | 17.63 | Newmarket                     | 29.13 | Toronto West   | 7.86  |
| Windsor West                   | 16.89 | Hamilton Mountain             | 28.61 | Toronto North  | 7.59  |
| Brampton                       | 15.95 | Peterborough                  | 28.55 | Kitchener  | 7.52  |

Table 1: Top 10 Cities with Highest Concentration of Each Pollutant

The table above is three separate tables each one for a respective pollutant, they are sorted in descending Mean and only the top 10 are being shown based on the overall Mean of the City, over all years from 2003 to 2022; as such, the cities with the high concentration will be listed first. Therefore, we can draw a few key observations about pollutant concentrations across these Ontario cities:

- The highest mean  $\text{NO}_x$  readings (31.32 ppb) appear at Toronto West, followed closely by other Toronto stations (Toronto East, Toronto Downtown) and industrial/urban areas like Hamilton and Windsor.
- Port Stanley shows the highest  $\text{O}_3$  levels (32.84 ppb), with other high concentrations at Tiverton, Grand Bend, and Parry Sound—generally smaller or semi-rural communities.
- Sarnia tops the  $\text{PM}_{2.5}$  list (9.57  $\mu\text{g}/\text{m}^3$ ), followed by Windsor (West and Downtown) and Hamilton stations, reflecting the influence of industrial facilities and cross-border pollution.

## 4.2 The Top 10 Regions and Years with Highest Concentration of Each Pollutant

In order to view the pollutant concentration changes by region, we grouped the cities by region based on the map of Ontario from the Ministry of Natural Resources and Forestry, and calculated the mean of the cities in each region grouped by year.

| NO <sub>x</sub> Table (in ppb) |                  |       | O <sub>3</sub> Table (in ppb) |                 |       | PM <sub>2.5</sub> Table (in µg/m <sub>3</sub> ) |                 |       |
|--------------------------------|------------------|-------|-------------------------------|-----------------|-------|---|-----------------|-------|
| Year                           | Region           | Conc. | Year                          | Region          | Conc. | Year  | Region          | Conc. |
| 2003                           | Central Ontario  | 31.45 | 2010                          | Western Ontario | 29.54 | 2005  | Western Ontario | 9.48  |
| 2005                           | Central Ontario  | 28.17 | 2022                          | Western Ontario | 29.49 | 2014  | Western Ontario | 9.27  |
| 2003                           | Western Ontario  | 26.98 | 2007                          | Eastern Ontario | 29.46 | 2003  | Western Ontario | 8.93  |
| 2004                           | Central Ontario  | 26.94 | 2010                          | Eastern Ontario | 29.20 | 2015  | Western Ontario | 8.79  |
| 2006                           | Central Ontario  | 23.80 | 2021                          | Western Ontario | 29.10 | 2013  | Western Ontario | 8.72  |
| 2007                           | Central Ontario  | 21.68 | 2018                          | Eastern Ontario | 28.97 | 2005  | Central Ontario | 8.56  |
| 2004                           | Western Ontario  | 21.55 | 2012                          | Western Ontario | 28.82 | 2004  | Western Ontario | 8.44  |
| 2008                           | Central Ontario  | 20.07 | 2016                          | Western Ontario | 28.79 | 2014  | Central Ontario | 8.33  |
| 2005                           | Western Ontario  | 20.00 | 2013                          | Western Ontario | 28.57 | 2007  | Western Ontario | 8.30  |
| 2004                           | Northern Ontario | 19.04 | 2008                          | Eastern Ontario | 28.54 | 2003  | Central Ontario | 8.26  |

Table 2: Top 10 Regions and Years with Highest Concentration of Each Pollutant

This table follows a similar format to the Table 1, listing the highest concentrations first. The concentrations of NO<sub>x</sub> (in ppb) and PM<sub>2.5</sub> (in µg/m<sub>3</sub>) were highest in the early to mid-2000s, particularly in Western and Central Ontario, and have since declined. Meanwhile, O<sub>3</sub> (in ppb) levels are highest in Western Ontario dominating the top 10 with Northern Ontario not occupying a single spot. The O<sub>3</sub> concentrations show minor decrease, likely reflecting O<sub>3</sub> long lifespan in the atmosphere.

## 4.3 The Top 10 Years with Highest Concentration of Each Pollutant

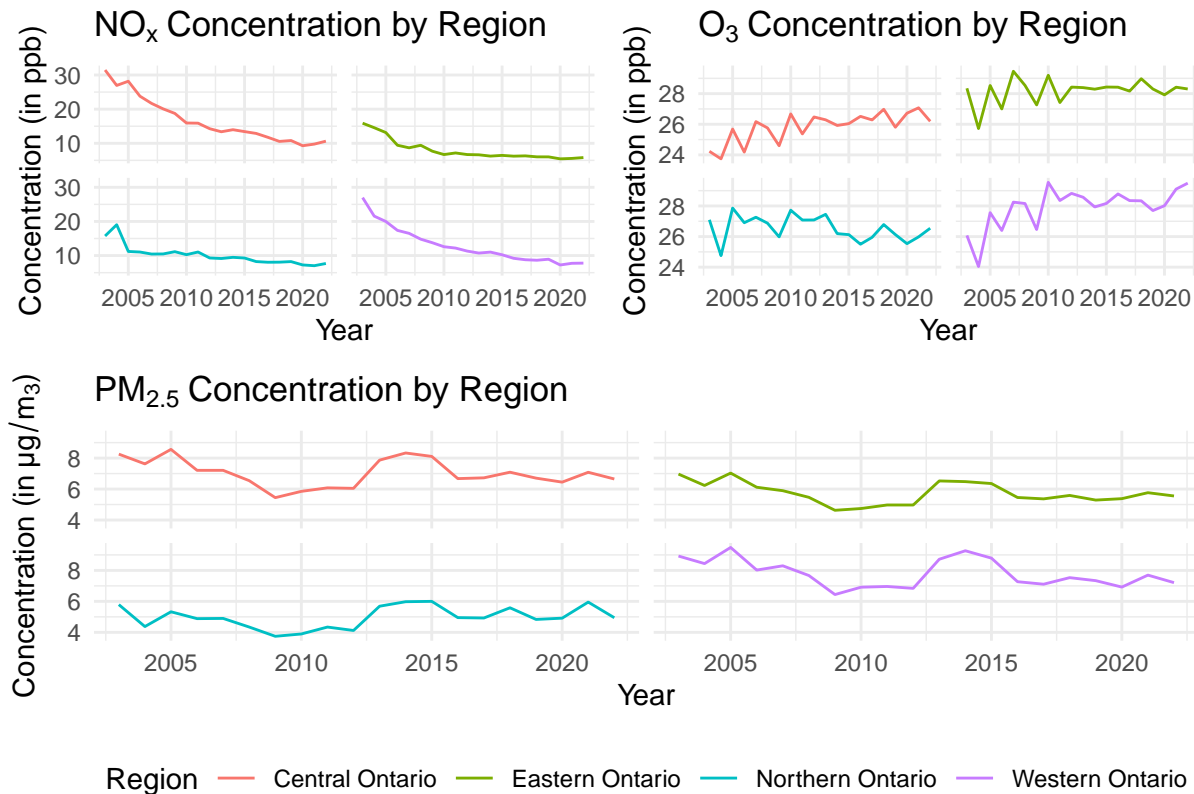
| NO <sub>x</sub> Table (in ppb) |       | O <sub>3</sub> Table (in ppb) |       | PM <sub>2.5</sub> Table (in µg/m <sub>3</sub> ) |       |
|--------------------------------|-------|-------------------------------|-------|---|-------|
| Year                           | Conc. | Year                          | Conc. | Year  | Conc. |
| 2003                           | 25.08 | 2010                          | 28.41 | 2005  | 8.36  |
| 2004                           | 23.42 | 2021                          | 27.97 | 2003  | 8.04  |
| 2005                           | 22.43 | 2018                          | 27.88 | 2014  | 7.96  |
| 2006                           | 17.87 | 2012                          | 27.85 | 2015  | 7.73  |
| 2007                           | 15.80 | 2007                          | 27.83 | 2013  | 7.61  |
| 2008                           | 15.04 | 2022                          | 27.83 | 2004  | 7.30  |
| 2009                           | 13.91 | 2013                          | 27.73 | 2006  | 7.04  |
| 2011                           | 12.36 | 2016                          | 27.67 | 2007  | 7.03  |
| 2010                           | 12.35 | 2017                          | 27.43 | 2021  | 6.92  |
| 2012                           | 11.23 | 2008                          | 27.41 | 2018  | 6.78  |

Table 3: Top 10 Years with the Highest Concentration

The table reveals that NO<sub>x</sub> levels peaked in 2003 – the first year of the dataset – and has steadily decreased since. Additionally, PM<sub>2.5</sub> concentrations were highest in the mid-2000s before decreasing, and O<sub>3</sub> levels, which peaked around 2010–2012, have remained relatively high before a minor decrease in later years.

## 5. Graphs

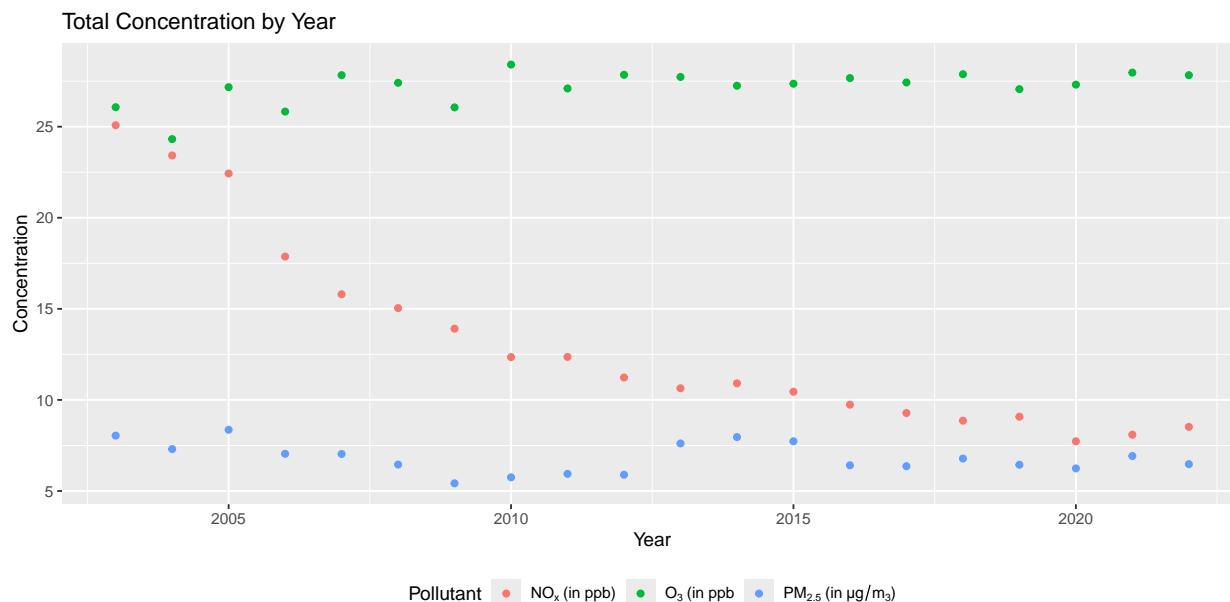
### 5.1 Line Charts of Yearly Concentration by Region



The Line Charts of Yearly Concentration by Region display the concentration of each pollutant faceted by region over 2003 till 2022. The x-axes represent the year and the y-axes represent the concentration of the pollutant at that year in its respective unit.

The graphs reveal a notable improvement in air quality over the years, with clear declines in both NO<sub>x</sub> and PM<sub>2.5</sub> levels across all regions, suggesting that emission controls and cleaner technologies have had a significant impact. NO<sub>x</sub> concentrations show a consistent downward trend from higher levels in the early years, converging towards lower values by 2020, while PM<sub>2.5</sub> levels exhibit a marked decrease, with the initial difference between regions decreasing over time. In contrast, O<sub>3</sub> concentrations appear relatively stable – with some fluctuations – which reflects the complex nature of ozone formation that depends on multiple factors such as sunlight, temperature, and precursor emissions.

## 5.2 Scatter Plot of Overall Pollutant Concentration by Year



This Scatterplot of Overall Pollutant Concentration by Region displays the overall concentration, calculated by taking the mean of each year for all cities in the dataset, of each pollutant over 2003 till 2022. The x-axes represent the year and the y-axes represent the concentration of the pollutant at that year in its respective unit.

The plot reveals a downward trends of NO<sub>x</sub> until 2020, with concentrations appearing to drop from around 25 ppb to near 8 ppb in 2020. This substantial decrease suggests that measures aimed at reducing emissions through regulations and increased standards have been effectively implemented over the years. On the other hand, PM<sub>2.5</sub> and O<sub>3</sub> do not show a steady downward trend; instead, their concentrations fluctuate over the period analyzed. These fluctuations show the challenges in the efforts of controlling pollutants that are not directly emitted but are produced by chemical interactions in the atmosphere.

## 5.3 Heatmap

## 6. Hypothesis Testing

### Estimating the Average Concentration of pollutants in 2022

To estimate the average concentration of pollutants in 2022, we can use bootstrapping to find a 95% confidence interval of the pollutant concentration of the cities.

In order to use bootstrapping to find a 95% confidence interval, we will repeated take samples of the pollutants filtered for only the year 2022 and calculate the mean of each sample. This process of re-sampling will be repeated 1000 times to produce a sampling distribution for each pollutant. From the sample distribution, we can find 2.5% and 97.5% quantile.

```
## 95% confidence interval for NOX is [ 7.219514 , 10.0246 ]
## 95% confidence interval for O3 is [ 27.10182 , 28.6084 ]
## 95% confidence interval for PM25 is [ 6.140158 , 6.821632 ]
```

## 7. Bootstrapping

### Estimating the Average Concentration of Pollutants in 2022

To estimate the average concentration of pollutants in 2022, we can use bootstrapping to find a 95% confidence interval of the mean pollutants of the cities.

In order to use bootstrapping to find a 95% confidence interval, we will repeated take samples of the pollutants filtered for only the year 2022 and calculate the mean of each sample. This process of re-sampling will be repeated 1000 times to produce a sampling distribution for each pollutant. From the sample distribution, we can find 2.5% and 97.5% quantile.

```
## 1 ) Mean NOX concentration of: 8.523611 , 95% confidence interval for NOX is [ 7.219514 , 10.0246 ]
## 2 ) Mean O3 concentration of: 27.83395 , 95% confidence interval for O3 is [ 27.10182 , 28.6084 ]
## 3 ) Mean PM25 concentration of: 6.474474 , 95% confidence interval for PM25 is [ 6.140158 , 6.821632 ]
```

Based on the result, we have:

- 1) A mean  $\text{NO}_x$  concentration of 8.523611. A 95% confidence interval of [7.219514 , 10.0246] for the concentration of  $\text{NO}_x$ . This means that out of 100 samples, 95 samples will have a mean  $\text{NO}_x$  concentration between [7.219514 , 10.0246].
- 2) A mean  $\text{O}_3$  concentration of 27.83395. A 95% confidence interval of [27.10182 , 28.6084] for the concentration of  $\text{O}_3$ . This means that out of 100 samples, 95 samples will have a mean  $\text{O}_3$  concentration between [27.10182 , 28.6084].
- 3) A mean  $\text{PM}_{2.5}$  concentration of 6.474474. A 95% confidence interval of [6.140158 , 6.821632] for the concentration of  $\text{PM}_{2.5}$ . This means that out of 100 samples, 95 samples will have a mean  $\text{PM}_{2.5}$  concentration between [6.140158 , 6.821632].

## 8. Regression

### Analyzing the Relationship Between Pollutants and Year

In order to analyze the relationship of a pollutant with other pollutants and the year, we would treat the year and two pollutants as continuous variables to fit a regression. For example, if we were examining the relationship between the overall yearly  $\text{NO}_x$  concentration and the year, the  $\text{O}_3$  concentration and the overall yearly  $\text{PM}_{2.5}$  concentration, we will take the logarithm of the  $\text{NO}_x$  concentration as the dependent variable and the year, the overall yearly  $\text{O}_3$  concentration and the overall yearly  $\text{PM}_{2.5}$  concentration all as the independent variable.

```
##
## Call:
## lm(formula = formula(log(Conc.NOX) ~ Year + Conc.PM25 + Conc.O3),
##     data = Combined_Yearly_Mean)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.137814 -0.034320 -0.005249  0.022683  0.173198
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept) 105.538569    7.485517   14.099 1.93e-10 ***
## Year        -0.050618    0.003861  -13.110 5.65e-10 ***
## Conc.PM25    0.062989    0.023487    2.682  0.0164 *
## Conc.O3      -0.058474    0.023494   -2.489  0.0242 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.08013 on 16 degrees of freedom
## Multiple R-squared:  0.9584, Adjusted R-squared:  0.9506
## F-statistic: 122.9 on 3 and 16 DF,  p-value: 2.925e-11
```

From the linear model, it suggests that there is a strong relationship between the logarithm of the  $\text{NO}_x$  concentration and the year, the  $\text{PM}_{2.5}$  concentration and the  $\text{O}_3$  concentration. In particular, by analyzing the p-values of each parameter, it can be seen that the year parameter is especially significant and that the mean  $\text{PM}_{2.5}$  and  $\text{O}_3$  levels have moderate significance on the logarithm of the mean  $\text{NO}_x$  concentration in comparison. The coefficient of determination ( $R^2$ ) of 0.9584 suggests that 95.84% of the variables in the logarithm of the  $\text{NO}_x$  concentration can be explained by the model.

## Interpreting regression parameters

Here, the regression parameters refer to the coefficients of the independent variables of the model.

### 1) Intercept:

- The intercept coefficient expresses the value of the dependent variable (logarithm of the  $\text{NO}_x$  concentration) when all other features are equal to 0. In this case, an intercept of  $\sim 105.54$  suggests that the logarithm of the  $\text{NO}_x$  concentration is equal to  $\sim 105.54$  when all independent variables are equal to 0.

### 2) Year, Mean $\text{PM}_{2.5}$ Concentration, and Mean $\text{O}_3$ Concentration:

- The coefficient for Year expresses the linear effect of the Year parameter; specifically, a coefficient of -0.050621 suggests that as the Year increases by 1, the logarithm of the  $\text{NO}_x$  concentration decreases by 0.050621. Similarly, a coefficient for the  $\text{PM}_{2.5}$  level of 0.062955 and a coefficient for the  $\text{O}_3$  level of -0.058456 suggests that as the  $\text{PM}_{2.5}$  concentration increases by 1 and the  $\text{O}_3$  level increases by 1 there will be an increase in the logarithm of the  $\text{NO}_x$  level by 0.062955 and a decrease in the logarithm of the  $\text{NO}_x$  level by 0.058456, respectively.

Overall, this model suggests a strong linear relationship between the logarithm of the mean  $\text{NO}_x$  level and the year, the mean  $\text{PM}_{2.5}$  level and the mean  $\text{O}_3$  level. In other words, this model also explains that the  $\text{NO}_x$  concentration has an exponential relationship with the features used, in particular an exponential decay relationship with the year.

## 8 Cross Validation

Previously, we analyzed the regression parameters to see how well our model performed on our data set. In this case, we will further analyze the relationships of pollutants with each other and the year by performing cross validation.

Specifically, we will analyze the relationship between  $\text{PM}_{2.5}$  and the other pollutants alongside the year by using k-fold cross validation. To perform k-fold cross validation, we will split the data set into k even pieces and use each  $i^{\text{th}}$  fold to check the validity of our model that is trained on the remaining  $k - 1$  folds by taking predictions using the data in the  $i^{\text{th}}$  fold. This process will be repeated for each fold of the data set, a mean squared error (MSE) will be calculated for each fold, and an average of the MSE's will be taken at the end.

## The average MSE is: 0.5452378

## Conclusion on results

The average MSE from the k-fold cross validation was 0.5452296 which is relatively low and maybe indicate the moderate predicting power of the linear model where the  $\text{PM}_{2.5}$  level was the dependent variable with year,  $\text{NO}_x$  concentration and the  $\text{O}_3$  concentration were the independent variables. This suggests that there is likely a linear relationship between the  $\text{PM}_{2.5}$  concentration, the year,  $\text{NO}_x$  concentration and  $\text{O}_3$  concentration.