

Data Mining Air Quality Data for Athens, Greece

Melina Efraimidou¹, Maria Kanaki¹, Ioannis Athanasiadis²,
Pericles Mitkas³, Kostas Karatzas⁴

Abstract

Urban air quality management and information systems are required to include advanced capabilities of quick, effective and easy to operate environmental data analysis applications, for information extraction and analysis and for the support of decision making. These systems are based on the need of city authorities and national governments to establish a framework which enables them to take actions, in order to ensure that air quality is improved and relevant standards are maintained in urban areas. In this context, quantitative data-driven decision support models are challenged by the difficulties in handling dynamic and uncertain features of real-world environmental systems. In addition, conditions for environmental management keep changing with time, demanding periodically updated decision support. These properties can be realized by learning from environmental data, using knowledge discovery techniques. In the present paper, data mining techniques are applied for data analysis and for the construction of forecasting modules towards decision making, on the basis of selected air quality information for Athens, Greece. Conclusions are drawn concerning the performance of algorithms, and for the research to be conducted in the future.

1. Introduction

Urban air quality information originates either from observations or from mathematical tools-models and estimations. While the former correspond to the current status of air quality, and may be directly interpreted in terms of human health risk and eco-system degradation potential or effect, the latter provide forecasting capabilities in advance, thus offering decision makers with the opportunity to take preventive measures that would "smooth" or alter the results of a forecasted "episode" or even "crisis". Earlier research work has dealt with using knowledge discovery techniques mainly for air quality associated incident forecasting. Several models have been built for predicting incidents that may occur in the near future. For instance, conventional statistical regression models (Bordignon et. al., 2002, Huang and Smith, 1999, Kim and Goldmann, 2001) and time-series analysis have been applied to predict ozone levels (Zolghadri and Henry, 2004). Neural networks have been used for short-term ozone prediction (Yi and Prybutok, 1996, Shiva Nagendra and Khare, 2006.), while case-based reasoning (Kalapanidas and Avouris, 2001) and classification and regression trees (Barrero et. al., 2006) have been employed for predicting air pollutant concentrations. In the present paper several data mining algorithms are applied, for analyzing air quality information and for forecasting the maximum ozone concentration levels on a daily basis in a dense urban area.

¹ MSc Programme on Advanced Systems on Computers and Telecommunications, Dept. of Electrical and Computer Engineering, Aristotle University, 541 24 Thessaloniki, Greece

² IDSIA - Istituto Dalle Molle di Studi sull'Intelligenza Artificiale Galleria 2, CH-6928 Manno, Switzerland

³ Dept. of Electrical and Computer Engineering, Aristotle University, 541 24 Thessaloniki, Greece, e-mail: mitkas@eng.auth.gr

⁴ Dept. of Mechanical Engineering, Aristotle University, Box 484, 54124 Thessaloniki, Greece, e-mail: kkara@eng.auth.gr

2. Data mining air quality data

2.1 The Athens air quality data set

Legislative acts, as the US Clean Air Act and the 96/62/EU framework directive for urban air management (and the accompanying daughter directives) have delimited certain thresholds for characterizing the quality for the atmospheric environment. In this background, qualitative information (like air quality indicators) is identified as being of great importance for an operational Air Quality Management System. On this basis, the current paper presents first results concerning the usage of data mining techniques for the quantitative and qualitative analysis of air quality data, and the prediction (forecasting) of pollution levels, using maximum hourly ozone concentration of a day as the parameter of investigation. For this purpose, measurement data, coming from the monitoring network operated by the Ministry of Environmental Physical Planning and Public Works were taken under consideration. The available time series of hourly values covered the period from 1/1/1999 to 1/7/2002, for the 15 stations indicated in Figure 1, and included a number of variables (air pollutants plus basic meteorological parameters) per station.

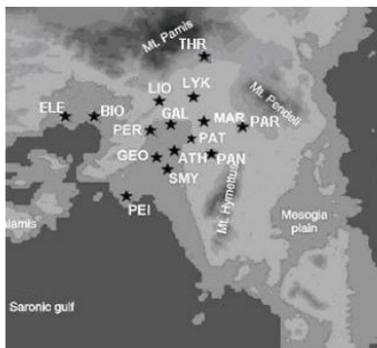


Fig. 1:
Locations and short names of the Athens air quality monitoring station network.

2.2 Data mining the Athens air quality data set

For all stations, 69% of the values were used for training, and 31% for forecasting validation. Aiming at investigating the operational performance of state-of-the-art data mining methods and tools towards ozone forecasting, a number of mining methods were employed. Contrary to the conventional statistical approaches, data mining algorithms used in this paper utilize other criteria/functions, as the information gain and entropy, distance-metrics, or memory-based learning, for encapsulating data-driven knowledge and ultimately drawing conclusions. Thus, different types of algorithms for classification were applied.

Among them, instance-based learners are used, as iBK, Kstar, Nnge (Nearest Neighbor With Generalization), rule-based classifiers, as Conjunctive Rules, OneR, decision trees (C4.5 or J48 algorithm, Decisionstump), along with Bayesian Classifiers (BayesNet, NaiveBayes), Neural networks (Multilayer Perceptron) and Fuzzy-Lattice Reasoning (FLR). All these algorithms are implemented in WEKA (The Waikato Environment for Knowledge Analysis, WEKA 2004). The WEKA platform (Witten and Frank, 1999) was used (notably this is an open source software) for the data mining experiments described below.

3. Results and discussion

In order to perform the experiments concerning the forecasting capabilities of the selected algorithms, two sets of ozone limit values were applied: the one resulting from the EU experience, practice and relevant legislation, and the other resulting from the statistical analysis of the available datasets (Table 1).

Ozone concentration ($\mu\text{g}/\text{m}^3$)	Clarification	Ozone concentration ($\mu\text{g}/\text{m}^3$)	Clarification
0-44	very low	0-29	very low
45-89	low	30-59	Low
90-134	medium	60-99	low – medium
135-179	high	100-139	Medium
180-239	information threshold	140-179	High
>240	alert threshold	180-239	information threshold
		>240	alert threshold

Table 1: Ozone concentration categorisation. EU oriented (left, EU), and statistical analysis oriented (right, AUTH)

The main goal of the experiments performed was to investigate the effectiveness of the algorithms in ozone forecasting for the next 8, 24, 48 and 72 hours at each station. For this purpose, a number of success criteria were used, among which the total weighted percentage of successful prediction is selected to be presented hereafter. For D_i the number of available values (samples) for station i and P_i the percentage of successful prediction for station i , the total weighted percentage of successful prediction (TWSPSP) is calculated as follows:

$$\text{TWSPSP} = (\sum D_i \cdot P_i) / (\sum D_i), i=1 \text{ to } 15 \text{ (i.e. for all 15 monitoring stations studied)}$$

The application of the method in all stations studied resulted in high scores for the TWSPSP. Although many experiments have been conducted, only the best (concerning performance) results are presented in Table 2, for reasons of brevity. These results correspond to algorithm J48, which is the C4.5 decision tree learning algorithm invented by Ross Quinlan from Sydney University (http://www.aiknomic.com/categorical_learners.html). Overall, the best performance when it comes to successful forecasting derives from the application of the J48 algorithm, followed by Nnge and the MLP Neural Networks. The rest of the investigated algorithms perform within the same value spectrum.

8EU - 8AUTH	24EU - 24AUTH	48EU - 48AUTH	72EU - 72AUTH
80% - 75,4%	86,5% - 85,4%	95,5% - 86,8%	96,6% - 89,3%

Table 2:

Total weighted percentage of successful prediction of the J48 algorithm for various time frames of forecast and on the basis of two different ozone value categorisation schemes (EU and AUTH, ref Table 1).

4. Conclusions

A number of data mining algorithms have been applied for the construction of forecasting models concerning maximum per day hourly ozone concentration values, for a total of 15 monitoring sites in Athens,

Greece. Forecasting results up to 95% average across all stations, demonstrate a good performance that should be considered for air quality forecasting modules applied at an operational basis. This implies the potential added value of applying classification techniques for air quality forecasting and assessment, and that data mining techniques could become the medium for overcoming the obstacles of conventional models. Future work may be directed towards a thorough comparison of classical approaches based on statistics and time series analyses, with those of data mining using classification techniques. Future activities may include the investigation of interrelations between monitoring sites and between pollutants and meteorological parameters, thus aiming to reveal "hidden" interconnection, periodicities and other qualitative and quantitative characteristics within the environmental area of interest..

Bibliography

- Barrero, M.A., Grimalt, J.O., & Cant, o'n L. (2006): Prediction of daily ozone concentration maxima in the urban atmosphere. *Chemometrics and Intelligent Laboratory Systems*, **80**, 67–76.
- Bordignon, S., Gaetan, C. & Lisi, F. (2002): Nonlinear models for ground-level ozone forecasting. *Statistical Methods and Applications*, Vol. 11, p. 227-246
- Huang, L. S. & Smith, R. L. (1997): Meteorologically-dependent trends in urban ozone. Technical Report 72, National Institute of Statistical Sciences.
- Kalapanidas, E. & Avouris, N. (2001): Short-term air quality prediction using a case based classifier. *Environmental Modelling and Software*, 16(3), p. 263-272.
- Kim, H.Y. & Guldman, J.M. (2001): Modeling air quality in urban areas: A cell-based statistical approach. *Geographical Analysis*, 33
- Shiva Nagendra, S.M. & Khare, M.(2006): Artificial neural network approach for modelling nitrogen dioxide dispersion from vehicular exhaust emissions. *Ecological Modelling*, **190**, 99–115.
- Yi, J. & Prybutok, V. R. (1996): A neural network model for the prediction of daily maximum ozone concentration in an industrialized urban area. *Environmental Pollution*, **92**, 349-357.
- WEKA, The Waikato Environment for Knowledge Analysis, www.cs.waikato.ac.nz/ml/weka, Version 3.4, 2004.
- Witten, I. & Frank, E. (1999): *Data Mining, Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann Publishers.
- Zolghadri, A. & Henry, D. (2004): Minmax statistical models for air pollution time series. Application to Ozone time series data measured in Bordeaux. *Environmental Monitoring and Assessment* **98**, 275–294.