

## Prediction of the Daily Mean PM<sub>10</sub> Concentrations Using Linear Models

J. C. M. Pires, F. G. Martins, S. I. V. Sousa, M. C. M. Alvim-Ferraz, M. C. Pereira  
LEPAE, Departamento de Engenharia Química, Faculdade de Engenharia, Universidade do Porto,  
Rua Dr. Roberto Frias s/n, 4200-465 Porto, Portugal.

---

**Abstract:** The performance of five linear models to predict the daily mean PM<sub>10</sub> concentrations was compared. The linear models proposed were: (i) multiple linear regression; (ii) principal component regression; (iii) independent component regression; (iv) quantile regression; and (v) partial least squares regression. The study was based on data from an urban site in Oporto Metropolitan Area and the analysed period was from January 2003 to December 2005. The linear models were evaluated with two datasets of different sizes belonging to the analysed period. Environmental data (SO<sub>2</sub>, CO, NO, NO<sub>2</sub> and PM<sub>10</sub> concentrations) and meteorological data (temperature, relative humidity and wind speed) were used as PM<sub>10</sub> predictors. During the training step, quantile regression presented the lowest residual errors for the two datasets. Independent component regression was the worst model using the larger dataset. Multiple linear regression, principal component regression and partial least squares regression presented similar results for both datasets. During the test step, independent component regression and quantile regression showed bad performance, while multiple linear regression, principal component regression and partial least squares regression presented similar results using the larger dataset. For the smaller dataset, the models that remove the correlation of the variables (principal component regression, independent component regression and partial least squares regression) presented better results than multiple linear regression and quantile regression. Independent component regression was the linear model with the lowest value of residual error. Concluding, the dataset size is also an important parameter for the evaluation of the models concerning the prediction of variables. The prediction of the daily mean PM<sub>10</sub> concentrations was more efficient when using independent component regression for the smaller dataset and partial least squares regression for the larger datasets.

**Key words:** PM<sub>10</sub> concentrations, multiple linear regression, principal component regression, independent component regression, quantile regression, partial least squares regression

---

### INTRODUCTION

Atmospheric particulate matter is made up of solid and liquid particles suspended in the atmosphere. They are emitted by: (i) natural (volcanic eruptions, seismic activity, forest fires, winds of great intensity or natural particle transport from the dry regions); and (ii) anthropogenic sources (all types of combustion and some industrial processes). In Europe, particulate matter is one of the most important air pollutants responsible for loss of human health<sup>[1]</sup>. In the last decade, several studies about health effects of particulate matter were published<sup>[2-5]</sup>. Long exposure to PM<sub>10</sub> (particles with diameter smaller than 10 µm) and to PM<sub>2.5</sub> (particles with diameter smaller than 2.5 µm) has been associated with respiratory and cardiovascular diseases. Recent research seems to indicate that

particles with smaller sizes are associated with childhood morbidity and mortality<sup>[5]</sup>.

The selection of the modelling techniques must consider some features, such as, complexity, flexibility, accuracy and speed of computation. The interpretability is also a very important characteristic of a model<sup>[6]</sup>. Without interpretability, the model is only used for prediction. In many situations, this type of model is enough. But in some cases, it is relevant to know the correlations between input variables and predictive variables. Furthermore, an interpretable model provides a sense of confidence. Comparing to the nonlinear models, linear models are simple and interpretable, taking less computation time.

Recently, several attempts have been made to model PM<sub>10</sub> concentrations using different models. G. Corani<sup>[7]</sup> tried to predict this pollutant using feed-

forward neural networks, pruned neural networks (nonlinear approaches) and lazy learning (local linear modelling approach). Comparing these three methodologies, lazy learning presents slightly better results than the other methods. Perez et al.<sup>[8]</sup> developed a neural network (nonlinear approach) to predict the maximum of the 24-h moving average of PM<sub>10</sub> concentration on the next day. This method was compared with linear perceptron (linear approach) and presented slightly better performance. Fuller et al.<sup>[9]</sup> used an empirical model to predict concentrations of PM<sub>10</sub> at background and roadside locations. The method was based on the regression analysis between PM<sub>10</sub> and NO<sub>x</sub>. The model accurately predicted daily mean PM<sub>10</sub> concentrations but presented some limitations. For example, it depends of the existence of a consistent relationship between PM<sub>10</sub> and NO<sub>x</sub> emissions. Thus, the studies that compared linear and nonlinear models<sup>[7,8]</sup> did not find a significant difference in the results obtained by the different methodologies.

This paper aims to analyse the performance of linear models, obtained by different methodologies, to predict the next day daily mean PM<sub>10</sub> concentrations. The considered models were: (i) multiple linear regression; (ii) principal component regression; (iii) independent component regression; (iv) quantile regression; and (v) partial least squares regression. The explanatory variables were meteorological data (daily means of temperature, relative humidity and wind speed) and environmental data (daily means of CO, SO<sub>2</sub>, NO, NO<sub>2</sub> and PM<sub>10</sub> concentrations of the previous day).

## MODELS

**Multiple Linear Regression:** Multiple linear regression (MLR) attempts to model the relationship between two or more explanatory variables and a response variable, by fitting a linear equation to the observed data. The dependent variable (y) is calculated by:

$$y = \hat{\beta}_0 + \sum_{i=1}^k \hat{\beta}_i x_i + \varepsilon \quad (1)$$

where  $x_i$  ( $i=1, \dots, k$ ) are the explanatory variables,  $\hat{\beta}_i$  ( $i=0, \dots, k$ ) are the regression coefficients, and  $\varepsilon$  is the error associated with the regression and assumed to be normally distributed with both expectation value zero and constant variance<sup>[10]</sup>.

The predicted value given by the regression model ( $\hat{y}$ ) is calculated by:

$$\hat{y} = \hat{\beta}_0 + \sum_{i=1}^k \hat{\beta}_i x_i \quad (2)$$

The most common method to estimate the regression parameters  $\hat{\beta}_i$  is the minimization of the sum of square errors (SSE). The equation is as follows:

$$\hat{\beta}_i = \arg \min \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3)$$

**Principal Component Regression:** Principal component regression (PCR) is a method that combines linear regression and principal component analysis (PCA). PCA creates new variables, the principal components (PC), that are orthogonal and uncorrelated. These variables are linear combinations of the original variables. The PC are ordered in such a way that the first component has the largest fraction of the original data variability<sup>[11-13]</sup>. To evaluate the influence of each variable in the PC, varimax rotation is generally used to obtain the rotated factor loadings that represent the contribution of each variable in a specific PC. PCR establishes a relationship between the output variable (y) and the selected PC obtained from the explanatory variables ( $x_i$ ).

**Independent Component Regression:** Independent component regression (ICR) is a method that combines linear regression and independent component analysis (ICA). In ICA, the input variables are considered linear combinations of latent variables. These latent variables are considered nongaussian and independent<sup>[14]</sup>. Therefore, they are called independent components (IC). PCA and ICA are considered linear representation models. While PCA determines the orthogonal variables (PC), ICA tries to find independent variables (IC). ICR establishes a relationship between the output variable (y) and the selected IC obtained from the explanatory variables ( $x_i$ )

**Quantile Regression:** Quantile regression (QR) was introduced by Koenker and Bassett<sup>[15]</sup> and can be seen as a natural extension of the least squares estimation of conditional mean models. This method presents some advantages when compared with ordinary least squares regression. For example, it allows the examination of the entire distribution of the variable of interest rather than a single measure of the central tendency of its distribution. It can also provide information about any linear or nonlinear relationships between the dependent variable and the explanatory variables without an a priori knowledge of the type of (potential) nonlinearities. Thus, it is more flexible to model data with

heterogeneous conditional distribution. To describe the quantile function, a random variable  $Y$  with the distribution function  $F(y) = \Pr(Y \leq y)$  is considered. The quantile function  $Q(\tau)$  with  $\tau \in [0, 1]$  is defined as follows:

$$Q(\tau) = \inf \{y : F(y) \geq \tau\} \quad (4)$$

The median is  $Q(1/2)$ , the first quartile is  $Q(1/4)$  and the first decile is  $Q(1/10)$ . The median regression minimizes a sum of absolute errors. The remaining conditional quantile functions are estimated by minimizing an asymmetrically weighted sum of absolute errors:

$$\hat{Q}(\tau) = \arg \min_a \left\{ \sum_{i: y_i \geq a} \tau |y_i - a| + \sum_{i: y_i < a} (1-\tau) |y_i - a| \right\} \quad (5)$$

Equation 6 presents a way to calculate the model parameters, considering quantile approach and the regression equation given by Equation 2:

$$\hat{\beta}(\tau) = \arg \min_{\beta(\tau)} \left\{ \sum_{i: y_i \geq \hat{y}_i} \tau |y_i - \hat{y}_i| + \sum_{i: y_i < \hat{y}_i} (1-\tau) |y_i - \hat{y}_i| \right\} \quad (6)$$

**Partial Least Squares Regression:** Partial least squares regression (PLSR) is probably the least restrictive of the various multivariate extensions of the MLR model<sup>[16]</sup>. This flexibility allows it to be used in situations where the use of traditional multivariate methods is severely limited, such as when there are fewer observations than predictor variables. In these situations, the MLR approach is not feasible due to the multicollinearity between the explanatory variables.

PLSR is based on linear transition from a large number of original descriptors to a new variable space based on small number of orthogonal factors (latent variables - LV)<sup>[16]</sup>. In other words, factors are orthogonal and linear combinations of original descriptors. Unlike some similar approaches (e.g. PCR), LV are chosen in such a way that provides maximum correlation with dependent variable. Thus, PLSR model contains the smallest necessary number of factors.

PLSR decomposes both explanatory variables (X) and the output variables (Y) as a product of a common set of orthogonal factors and a set of specific loadings. The complete procedure is as follows<sup>[16]</sup>:

**Step 1:** Normalization of X and Y:  $X_0 = X / \|X\|$  and  $Y_0 = Y / \|Y\|$ ;

**Step 2:** Definition of the vector  $u$  with random values;

**Step 3:** Estimation of the X weights:  $w = \frac{X_k^T u}{\|X_k^T u\|}$ ;

**Step 4:** Estimation of the X factor scores (LV):  $t = \frac{X_k w}{\|X_k w\|}$ ;

**Step 5:** Estimation of the Y weights:  $c = \frac{Y_k^T t}{\|Y_k^T t\|}$ ;

**Step 6:** Estimation of the Y factor scores:  $u = Y_k c$ ;

**Step 7:** Repetition of the steps 3 to 6 until the convergence of  $t$ ;

**Step 8:** Determination of the value of  $b$  used to predict Y from  $t$ :  $b = t^T u$ ;

**Step 9:** Determination of the X factor loadings:  $p = X_k^T t$ ;

**Step 10:** Elimination of the effect of  $t$  from X and Y:  $X_{k+1} = X_k - tp^T$  and  $Y_{k+1} = Y_k - btc^T$ ;

**Step 11:** Repetition of the steps 2 to 10 until the determination of a selected number of LV.

Each vector  $t$ ,  $u$ ,  $w$ ,  $c$  and  $p$  is stored in the columns of the correspondent matrices (T, U, W, C and P) and the scalar  $b$  is stored in a diagonal matrix (B). If  $X_k$  is a null matrix, all latent variables were determined.

The prediction of the dependent variable is done by  $\hat{Y} = TBC^T = XB_{PLS}$ . If all latent variables are used, the results of PLSR are similar to that obtained by PCR.

**Regression parameters validation:** It is important to know which explanatory variables are relevant to predict the dependent variable. For the studied models, PLSR is the only one that includes this step in its procedure. For MLR, PCR and ICR, the significance of each regression parameter in the models was evaluated through the calculation of their confidence interval. The parameter  $\hat{\beta}_i$  is statistically significant if<sup>[17]</sup>:

$$|\hat{\beta}_i| > \frac{t_{n-k-1}^{\alpha/2} \hat{\sigma}}{\sqrt{Sxx_i}} \quad (7)$$

where  $t$  is the Student t distribution,  $n$  is the number of points,  $k$  is the number of parameters,  $\alpha$  is the significance level,  $\hat{\sigma}$  is the standard deviation given by

$$\sqrt{\frac{SSE}{n-k-1}} \text{ and } Sxx_i \text{ is the sum of squares related to } x_i$$

$$\text{given by } \sum_{j=1}^n (x_{i,j} - \bar{x}_i)^2 .$$

For QR, bootstrap estimates of standard error (at the 95% confidence level) were calculated by randomly sampling each dataset with replacement (1000 times).

**Performance indexes:** The linear models were compared through the calculation of the following statistical parameters: mean bias error (MBE), mean absolute error (MAE), root mean squared error (RMSE), index of agreement ( $d_2$ ), that are commonly referred in literature<sup>[18, 19]</sup>.

MBE indicates if the observed values are over or under estimated. MAE and RMSE measure residual errors, which give a global idea of the difference between the observed and the modelled values. The values of  $d_2$  compare the difference between the mean, the predicted and the observed variables, indicating the degree of error free for the predictions<sup>[18, 19]</sup>.

## DATA

This study aims to predict the daily mean  $PM_{10}$  concentration of the next day. The concentrations of pollutants ( $SO_2$ , CO, NO,  $NO_2$  and  $PM_{10}$ ) were recorded in an urban site (Matosinhos) with traffic influences situated in Oporto Metropolitan Area, Northern Portugal.  $SO_2$  concentrations were obtained by the ultraviolet fluorescence method; CO concentrations were measured through nondispersive infrared spectrometric;  $NO_2$  was obtained through the chemiluminescence method;  $PM_{10}$  concentrations were obtained through the beta radiation attenuation method. These equipments were submitted to a rigid maintenance program and calibrated periodically. Measurements were continuously made and hourly average concentrations (in micrograms per cubic meter) were registered. Meteorological variables were measured on the left edge of Douro River at an altitude of 90 m approximately. These variables are hourly means of: air temperature (T), relative humidity (RH) and wind speed (WS). Daily average values for these variables were calculated and used if more than 75% of hourly values were available. The period of measurement was from January 2003 to December 2005. To evaluate the influence of the dataset size in the performance of the linear models, two datasets with different sizes were considered: (i) dataset 1 (DS1) considered 2003, 2004 and the first three trimesters of 2005 as training period and the last trimester of 2005 as test period; and (ii) dataset 2 (DS2), January 2003 to May 2003 was considered as training period and June 2003 was considered as test period. Additionally, QR and PLSR models required a validation period. For DS1, the last trimester of 2004 was considered as a

validation period. For DS2, May 2003 was considered as a validation period. The explanatory variables were standardized (zero mean and unit standard deviation).

## RESULTS AND DISCUSSION

Table 1 presents the correlation coefficients between pollutants and meteorological variables for the two analysed datasets (DS1 – upper triangular matrix; DS2 – lower triangular matrix). These coefficients provide a measure of linear relationship between two variables. Values in bold correspond to statistically significant coefficients<sup>[20]</sup>. For the evaluation of their statistical significance, the critical correlation coefficient was calculated (with a significance level of 0.05). A correlation coefficient is significant if its absolute value is greater than the critical value. In this preliminary analysis, concerning the correlations with the next day  $PM_{10}$  concentration ( $PM_{10(d+1)}$ ), only temperature was not considered statistically significant for DS2. The explanatory variables that presented highest correlation with  $PM_{10(d+1)}$  were CO, NO,  $NO_2$  and  $PM_{10}$  concentrations (for both datasets). In both datasets, as it was expected,  $PM_{10(d+1)}$  concentration had positive correlations with  $SO_2$ , CO, NO and  $NO_2$ . Some combustion and industrial processes (important  $PM_{10}$  sources) also increase the emission of these pollutants.  $PM_{10(d+1)}$  concentration had negative correlation with RH and WS. In wet weather, the particles in suspension can stick on the surface and can be whirled up into the air in dry weather. For high values of WS, there is an efficient dispersion of pollutants. Thus, high values of WS correspond to low pollutant concentrations.  $PM_{10(d+1)}$  concentration had also positive correlation with its concentration of the previous day.

For the statistical evaluation of the regression parameters, a t-test (significance level of 0.05) was performed for MLR, PCR and ICR. Considering the statistically significant parameters, new regressions were performed. Table 2 shows the statistically significant regression parameters for MLR, PCR and ICR. For MLR, the parameters  $\beta_1$  to  $\beta_8$  correspond to  $SO_2$ , CO, NO,  $NO_2$ , T, RH, WS and  $PM_{10}$ , respectively. The variables considered important for the prediction of the next day  $PM_{10}$  concentration were: (i) CO,  $NO_2$ , T, RH and  $PM_{10}$  for DS1; and (ii) CO, NO,  $NO_2$  and T for DS2. For PCR, the parameters  $\beta_i$  ( $i=1, 8$ ) correspond to each PC. All PC were used as input variables in PCR. However, the important predictors of  $PM_{10}$  concentration considered by this method were: (i) PC1 to PC6 for DS1; and (ii) PC1, PC2, PC4 and PC8 for DS2. Table 3 shows the rotated factor loadings for DS1 and DS2 that represent the contribution of each variable

Table 1: Correlation coefficients between pollutants and meteorological variables for the two analysed datasets (DS1 – upper triangular matrix; DS2 – lower triangular matrix).

	SO <sub>2</sub>	CO	NO	NO <sub>2</sub>	T	RH	WS	PM <sub>10</sub>	PM <sub>10(d+1)</sub>
SO <sub>2</sub>	1	<b>0.171</b>	<b>0.329</b>	<b>0.430</b>	<b>0.297</b>	<b>-0.163</b>	<b>-0.155</b>	<b>0.371</b>	<b>0.276</b>
CO	<b>0.510</b>	1	<b>0.860</b>	<b>0.786</b>	<b>-0.374</b>	<b>-0.086</b>	<b>-0.343</b>	<b>0.638</b>	<b>0.486</b>
NO	<b>0.575</b>	<b>0.917</b>	1	<b>0.798</b>	<b>-0.211</b>	-0.062	<b>-0.328</b>	<b>0.565</b>	<b>0.439</b>
NO <sub>2</sub>	<b>0.619</b>	<b>0.882</b>	<b>0.832</b>	1	-0.043	<b>-0.318</b>	<b>-0.283</b>	<b>0.734</b>	<b>0.594</b>
T	-0.089	<b>-0.433</b>	<b>-0.366</b>	<b>-0.295</b>	1	<b>-0.251</b>	-0.022	<b>0.168</b>	<b>0.153</b>
RH	<b>-0.333</b>	-0.138	-0.167	<b>-0.338</b>	<b>-0.270</b>	1	<b>-0.278</b>	<b>-0.275</b>	<b>-0.323</b>
WS	<b>-0.221</b>	<b>-0.532</b>	<b>-0.447</b>	<b>-0.482</b>	<b>0.188</b>	-0.065	1	<b>-0.234</b>	<b>-0.150</b>
PM <sub>10</sub>	<b>0.436</b>	<b>0.773</b>	<b>0.684</b>	<b>0.778</b>	-0.012	<b>-0.289</b>	<b>-0.401</b>	1	<b>0.685</b>
PM <sub>10(d+1)</sub>	<b>0.352</b>	<b>0.640</b>	<b>0.551</b>	<b>0.675</b>	-0.052	<b>-0.337</b>	<b>-0.289</b>	<b>0.669</b>	1

Values in bold correspond to statistically significant correlation coefficients

Table 2. Statistically significant regression parameters for MLR, PCR and ICR using both datasets (DS1 and DS2).

		β <sub>0</sub>	β <sub>1</sub>	β <sub>2</sub>	β <sub>3</sub>	β <sub>4</sub>	β <sub>5</sub>	β <sub>6</sub>	β <sub>7</sub>	β <sub>8</sub>
MLR			SO <sub>2</sub>	CO	NO	NO <sub>2</sub>	T	RH	WS	PM <sub>10</sub>
	DS1	42.12		2.36		2.88	1.98	-2.74		10.82
	DS2	40.89		16.17	-8.17	10.54	5.87			
PCR			PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8
	DS1	42.12	7.35	4.35	1.41	4.82	4.87	5.19		
	DS2	40.89	7.62	3.90		7.23				-13.49
ICR			IC1	IC2	IC3	IC4	IC5	IC6	IC7	IC8
	DS1	26.39	7.57		3.14		-2.95	1.65	2.49	3.19
	DS2	6.11	7.50	8.15	6.42		-10.83	-3.55		4.22

Table 3. Varimax rotated loadings using both datasets (DS1 and DS2).

		PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8
DS1	SO <sub>2</sub>	-0.171	0.159	-0.061	<b>0.960</b>	0.065	-0.121	-0.035	0.002
	CO	<b>-0.842</b>	-0.277	-0.176	0.020	0.050	-0.318	-0.008	-0.280
	NO	<b>-0.952</b>	-0.090	-0.132	0.162	-0.003	-0.144	0.034	0.139
	NO <sub>2</sub>	<b>-0.742</b>	-0.014	-0.129	0.234	0.223	-0.355	-0.449	-0.004
	T	0.177	<b>0.957</b>	-0.033	0.160	0.132	-0.088	0.000	0.010
	RH	0.076	-0.127	-0.160	-0.065	<b>-0.967</b>	0.106	0.033	0.004
	WS	0.194	-0.031	<b>0.962</b>	-0.062	0.161	0.077	0.021	0.008
	PM <sub>10</sub>	-0.440	0.129	-0.094	0.161	0.141	<b>-0.855</b>	-0.053	-0.011
DS2	SO <sub>2</sub>	-0.281	-0.028	-0.076	<b>0.929</b>	0.167	-0.146	0.045	0.006
	CO	<b>-0.753</b>	-0.276	-0.275	0.209	0.053	-0.429	0.081	0.211
	NO	<b>-0.878</b>	-0.186	-0.190	0.280	0.065	-0.258	0.000	-0.096
	NO <sub>2</sub>	<b>-0.607</b>	-0.180	-0.250	0.327	0.225	-0.439	0.432	0.013
	T	0.204	<b>0.964</b>	0.068	-0.029	0.150	-0.016	-0.027	-0.005
	RH	0.082	-0.147	-0.050	-0.150	<b>-0.967</b>	0.107	-0.031	-0.001
	WS	0.225	0.070	<b>0.955</b>	-0.074	0.055	0.151	-0.037	-0.006
	PM <sub>10</sub>	-0.442	0.064	-0.181	0.162	0.141	<b>-0.848</b>	0.046	-0.004

Values in bold indicate the variables that most influence each principal component.

in a specific PC. In both datasets, PC1 is heavily loaded by CO, NO and NO<sub>2</sub> and PC2 to PC6 had greater contributions of T, WS, SO<sub>2</sub>, RH and PM<sub>10</sub>, respectively. Thus, the most important original variables selected by PCR were: (i) all variables for DS1; and (ii) SO<sub>2</sub>, CO, NO, NO<sub>2</sub> and T for DS2. Similarly to PCR, with ICR the parameters β<sub>i</sub> (i=1, 8) correspond to each IC. All IC were used as input variables in ICR. However, the important predictors of PM<sub>10</sub> concentration considered by this method were: (i)

IC1, IC3, IC5, IC6, IC7 and IC8 for DS1; and (ii) IC1, IC2, IC3, IC5, IC6 and IC8 for DS2. Table 4 presents the correlation matrix between the original variables and IC. These correlation values showed the importance of each original variable on the prediction of PM<sub>10</sub> concentration. Thus, the original variables considered relevant for prediction of PM<sub>10</sub> concentration were: (i) SO<sub>2</sub>, CO, NO, NO<sub>2</sub>, WS and PM<sub>10</sub> for DS1; and (ii) CO, NO, NO<sub>2</sub>, T, RH, WS and PM<sub>10</sub> for DS2.

Table 4. Correlation matrix between the original variables and the IC using both datasets (DS1 and DS2).

		IC1	IC2	IC3	IC4	IC5	IC6	IC7	IC8
DS1	SO <sub>2</sub>	<b>0.329</b>	0.038	0.012	<b>-0.129</b>	0.014	<b>0.923</b>	<b>0.107</b>	<b>0.096</b>
	CO	<b>0.410</b>	<b>0.203</b>	0.005	<b>0.486</b>	<b>-0.244</b>	0.003	<b>0.288</b>	<b>0.642</b>
	NO	<b>0.428</b>	<b>0.146</b>	0.041	<b>0.195</b>	<b>0.138</b>	<b>0.110</b>	<b>0.235</b>	<b>0.818</b>
	NO <sub>2</sub>	<b>0.807</b>	<b>0.170</b>	<b>0.167</b>	<b>0.228</b>	0.027	<b>0.128</b>	<b>0.323</b>	<b>0.346</b>
	T	<b>0.291</b>	0.040	-0.053	<b>-0.901</b>	<b>-0.235</b>	<b>0.120</b>	<b>-0.081</b>	<b>-0.152</b>
	RH	<b>-0.521</b>	<b>0.776</b>	-0.011	<b>0.102</b>	<b>0.212</b>	0.010	<b>-0.256</b>	<b>0.071</b>
	WS	<b>0.133</b>	<b>-0.543</b>	-0.018	<b>0.127</b>	0.037	<b>-0.065</b>	<b>-0.798</b>	<b>-0.172</b>
	PM <sub>10</sub>	<b>0.526</b>	<b>0.141</b>	<b>0.582</b>	0.010	<b>-0.426</b>	<b>0.152</b>	<b>0.163</b>	<b>0.364</b>
DS2	SO <sub>2</sub>	-0.137	<b>0.300</b>	<b>0.389</b>	0.118	<b>-0.470</b>	-0.111	<b>0.700</b>	-0.013
	CO	<b>0.369</b>	<b>0.287</b>	<b>0.600</b>	-0.102	<b>-0.636</b>	0.065	-0.055	-0.019
	NO	0.018	<b>0.290</b>	<b>0.775</b>	-0.064	<b>-0.550</b>	0.012	-0.088	-0.031
	NO <sub>2</sub>	0.124	<b>0.498</b>	<b>0.331</b>	-0.011	<b>-0.782</b>	-0.043	-0.017	0.110
	T	-0.053	<b>0.201</b>	-0.078	<b>0.475</b>	<b>0.550</b>	-0.143	0.097	<b>0.621</b>
	RH	-0.060	<b>-0.616</b>	0.000	<b>-0.465</b>	0.035	<b>0.591</b>	-0.023	<b>0.223</b>
	WS	<b>-0.232</b>	<b>-0.661</b>	<b>-0.241</b>	<b>0.551</b>	0.058	<b>-0.364</b>	-0.057	-0.082
	PM <sub>10</sub>	<b>0.385</b>	<b>0.306</b>	<b>0.431</b>	-0.153	<b>-0.474</b>	<b>-0.384</b>	-0.017	<b>0.421</b>

Values in bold correspond to statistically significant correlation coefficients

Table 5. Statistically significant regression parameters for QR using both datasets (DS1 and DS2).

	$\tau$	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	$\beta_5$	$\beta_6$	$\beta_7$	$\beta_8$
			SO <sub>2</sub>	CO	NO	NO <sub>2</sub>	T	RH	WS	PM <sub>10</sub>
DS1	0.10	23.7		5.39			3.72	-3.50		4.6
	0.30	33.6					0.64	-2.73		12.3
	0.50	40.8					0.47	-2.41		15.7
	0.70	49.0				2.25		-2.32		16.3
	0.90	62.3				7.33				16.7
DS2	0.10	22.7	SO <sub>2</sub>	CO	NO	NO <sub>2</sub>	T	RH	WS	PM <sub>10</sub>
	0.30	22.9		13.80			2.99	-2.45		
	0.50	33.3								
	0.70	47.2								21.0
	0.90	62.3	-5.85			26.86				

Table 5 shows the statistically significant regression parameters for quantile regression. Five percentiles ( $\tau$ ) were selected. For the evaluation of the statistical significance of the regression parameters, their confidence intervals (significance level of 0.05) were determined by the application of bootstrap (calculated with 1000 times replacements). For DS1, the results showed that RH and PM<sub>10</sub> concentration were the most important explanatory variables. CO concentration and T were important only in low values of  $\tau$ , while NO<sub>2</sub> concentration was relevant in high values of  $\tau$ . For DS2, SO<sub>2</sub>, NO<sub>2</sub> and PM<sub>10</sub> concentrations were important for high values of  $\tau$ , while CO concentration, T and RH were relevant for low values of  $\tau$ .

For PLSR, it is important to determine the number of the LV needed to obtain the best generalisation for the prediction of PM<sub>10</sub> concentrations. The validation period was used to determine the number of LV that presents the lowest value of error and in both datasets, two variables were achieved. Table 6 shows the regression parameters for PLSR ( $B_{PLS}$ ) and the rotated

factor loadings that represent the contribution of each variable in a specific LV. For DS1, LV1 was heavily loaded by the CO, NO, NO<sub>2</sub> and PM<sub>10</sub> concentrations, while LV2 had a great contribution of T. For DS2, LV1 was heavily loaded by the SO<sub>2</sub>, CO, NO, NO<sub>2</sub> and PM<sub>10</sub> concentrations, while LV2 had a great contribution of T.

Table 7 shows the performance indexes for the different linear models for the training period. ICR presented bad performance in the larger dataset (DS1). QR was the linear method that presented better performance using the two datasets.

In the test step and for MLR, PCR and ICR, PM<sub>10</sub> concentrations were determined by the application of the regression equations achieved in the training step. For this step, the application of QR needs to predict the PM<sub>10</sub> percentile for each validation point, following the application of the correspondent regression equation. The PM<sub>10</sub> percentile was determined applying the k-nearest neighbour (k-NN) algorithm. This algorithm was used for classifying objects based on closest examples in the training data. It was based on the

Table 6. Regression parameters for PLSR ( $B_{PLS}$ ) and the rotated factor loadings using both datasets (DS1 and DS2).

	DS1			DS2		
	$B_{PLS}$	LV1	LV2	$B_{PLS}$	LV1	LV2
SO <sub>2</sub>	0.025	0.504	-0.060	-0.010	<b>0.722</b>	-0.273
CO	0.115	<b>0.826</b>	-0.371	0.196	<b>0.960</b>	-0.041
NO	0.041	<b>0.822</b>	-0.431	0.081	<b>0.918</b>	-0.112
NO <sub>2</sub>	0.172	<b>0.936</b>	-0.097	0.189	<b>0.955</b>	-0.003
T	0.149	0.087	<b>0.685</b>	0.096	-0.438	<b>0.700</b>
RH	-0.180	-0.395	-0.602	-0.138	-0.371	-0.459
WS	0.008	-0.309	0.374	0.017	-0.545	0.286
PM <sub>10</sub>	0.327	<b>0.868</b>	0.266	0.287	<b>0.847</b>	0.393

Values in bold indicate the variables that most influence each latent variable.

Table 7: Performance indexes of the different linear models for the training period using both datasets (DS1 and DS2).

		R	MBE	MAE	RMSE	d <sub>2</sub>
DS1	MLR	0.71	0.00	11.99	15.92	0.81
	PCR	0.71	0.00	11.96	15.91	0.81
	ICR	0.43	0.00	15.58	20.45	0.55
	QR	0.79	2.60	9.93	14.21	0.92
	PLSR	0.70	0.00	12.45	16.38	0.81
DS2	MLR	0.72	0.00	12.78	16.55	0.82
	PCR	0.73	0.00	12.51	16.42	0.83
	ICR	0.74	0.00	12.18	16.27	0.83
	QR	0.84	1.24	7.99	13.07	0.94
	PLSR	0.72	0.00	12.78	16.80	0.82

Table 8: Performance indexes of the different linear models for the test period using both datasets (DS1 and DS2).

		R	MBE	MAE	RMSE	d <sub>2</sub>
DS1	MLR	0.74	-1.12	12.66	18.43	0.83
	PCR	0.74	-0.90	12.68	18.37	0.84
	ICR	0.68	-1.90	13.97	20.22	0.71
	QR	0.60	2.16	15.20	21.95	0.86
	PLSR	0.75	-2.07	12.24	18.13	0.83
DS2	MLR	0.70	5.92	11.76	13.84	0.84
	PCR	0.76	-0.25	10.87	12.64	0.86
	ICR	0.79	0.64	10.02	11.83	0.88
	QR	0.72	-3.44	11.14	13.53	0.86
	PLSR	0.77	4.34	10.50	12.48	0.87

Euclidean distance between the correspondent validation point and the training points. The evaluation of the optimal value of  $k$  nearest training samples depends of the dataset. A good value of  $k$  can be achieved using cross-validation. The k-NN algorithm is as follows:

**Step 1:** Selection of the  $k$  value;

**Step 2:** Determination of  $k$  nearest training points from the validation point;

**Step 3:** Determination of the percentile of PM<sub>10</sub> concentration values correspondent to these training points;

**Step 4:** Application of the QR equations correspondent to these percentiles using the validation point:

$$\hat{y}_i = \hat{\beta}_0 + \sum_{i=1}^k \hat{\beta}_i x_i;$$

**Step 5:** Determination of the average of the  $k$  values of  $\hat{y}_i$ ;

**Step 6:** Repetition of the steps 2 to 5 for all validation points;

**Step 7:** Determination of the error associated to the value of  $k$ , based on the difference of the average values calculated above and the true values;

**Step 8:** Repetition of the steps 1 to 7 for different values of  $k$ ;

**Step 9:** Determination of the lowest value of error associated to the optimal value of  $k$ .

For the test step, it was necessary 23 and 5 nearest points for prediction of the percentile of the test points using DS1 and DS2, respectively. For PLSR, PM<sub>10</sub> concentrations were determined through the  $B_{PLS}$  (equation  $\hat{Y} = XB_{PLS}$ ) achieved with the number of LV obtained in the validation step. Table 8 presents the performance indexes achieved in the test step using DS1 and DS2. The results showed that for the larger dataset (DS1), QR and ICR presented the highest values of residual error and the remained models had similar results. For DS2, ICR presented the best performance, achieving the lowest value of RMSE. The other models that remove the correlation between the input variables (PCR and PLS) also presented good results. MLR and QR achieved similar performance indexes.

## CONCLUSIONS

Five linear models were used to predict the daily mean PM<sub>10</sub> concentrations using as predictors air pollutant (SO<sub>2</sub>, CO, NO, NO<sub>2</sub> and PM<sub>10</sub>) concentrations and meteorological parameters (temperature, relative humidity and wind speed) for two datasets with different sizes. PM<sub>10</sub> concentration presented positive correlation with SO<sub>2</sub>, CO, NO and NO<sub>2</sub> concentrations due to the similar emission sources. It was also positively correlated with temperature and negatively with relative humidity and wind speed. In wet weather, the particles in suspension can stick on the surface and can be whirled up into the air in dry weather. For high values of wind speed, there is an efficient dispersion of pollutants.

As the selected models are all linear, they are interpretable and their results were used to determine which variables were important in the prediction of daily mean PM<sub>10</sub> concentrations. All models considered different group of important variables, however CO and NO<sub>2</sub> were always selected.

During the training step, quantile regression presented the lowest residual errors using the two datasets. Multiple linear regression, principal component regression and partial least squares regression presented similar results. Independent component regression presented bad performance in the dataset 1, which may be the result of the dataset size.

During the test step and for the dataset 1, independent component regression and quantile regression showed bad performances, while multiple linear regression, principal component regression and partial least squares regression presented similar results. Using the smaller dataset (dataset 2), the models that remove the correlation of the variables (principal component regression, independent component regression and partial least squares regression) presented better results than multiple linear regression and quantile regression. Independent component regression was the linear model with the lowest value of residual error.

Concluding, the dataset size is also an important parameter for the evaluation of the models concerning the prediction of variables. The prediction of the daily mean PM<sub>10</sub> concentrations was more efficient when using independent component regression for the smaller dataset and partial least squares regression for the larger dataset.

#### ACKNOWLEDGMENTS

Authors are grateful to Comissão de Coordenação da Direcção Regional-Norte and to Instituto Geofísico da Universidade do Porto, for kindly providing the air quality and meteorological data. This work was supported by Fundação para a Ciência e Tecnologia (FCT). J. C. M. Pires also thanks the FCT for the fellowship SFRD/BD/23302/2005.

#### REFERENCES

1. Koelemeijer, R.B.A., C.D. Homan and J. Matthijsen, 2006. Comparison of spatial and temporal variations of aerosol optical thickness and particulate matter over Europe. *Atmos. Environ.*, 40 (27): 5304-5315.
2. Alvim-Ferraz, M.C., M.C. Pereira, J.M. Ferraz, A.M.C. Almeida e Mello and F.G. Martins, 2005. European Directives for Air Quality: Analysis of the New Limits in Comparison with Asthmatic Symptoms in Children Living in the Oporto Metropolitan Area, Portugal. *Hum. Ecol. Risk Assess.*, 11(3): 607-616.
3. Brunekreef, B. and S. T. Holgate, 2002. Air pollution and health. *The Lancet*, 360 (9341): 1233-1242.
4. Hoek, G., B. Brunekreef, S. Goldbohm, P. Fischer and P.A. van der Brand, 2002. Association between mortality and indicators of traffic-related air pollution in the Netherlands: a cohort study. *The Lancet*, 360 (9341): 1203-1209.
5. Kappos, A.D., P. Bruckmann, T. Eikmann, N. Englert, U. Heinrich, P. Hoppe, E. Koch, G. H. M. Krause, W.G. Kreyling, K. Rauchfuss, P. Rombout, V. Schulz-Klemp, W.R. Thiel and H.E. Wichmann, 2004. Health effects of particles in ambient air. *Int. J. Hyg. Environ. Health*, 207 (4): 399-407.
6. Guha, R., D.T. Staton, P.C. Jurs, 2005. Interpreting Computational Neural Network Quantitative Structure-Activity Relationship Models: A Detailed Interpretation of the Weights and Biases. *J. Chem. Inf. Model.*, 45: 1109-1121.
7. Corani, G., 2005. Air quality prediction in Milan: feed-forward neural networks, pruned neural networks and lazy learning. *Ecol. Model.*, 185 (2-4): 513-529.
8. Perez, P. and J. Reyes, 2002. Prediction of maximum of 24-h average of PM<sub>10</sub> concentrations 30 h in advance in Santiago, Chile. *Atmos. Environ.*, 36 (28): 4555-4561.
9. Fuller, G.W., D.C. Carslaw and H.W. Lodge, 2002. An empirical approach for the prediction of daily mean PM<sub>10</sub> concentrations. *Atmos. Environ.*, 36 (9): 1431-1441.
10. Agirre-Basurko, E., G. Ibarra-Berastegi and I. Madariaga, 2006. Regression and multilayer perceptron-based models to forecast hourly O<sub>3</sub> and NO<sub>2</sub> levels in the Bilbao area. *Environ. Modell. Softw.*, 21 (4): 430-446.
11. Abdul-Wahab, S.A., C.S. Bakheit, S.M. Al-Alawi, 2005. Principal component and multiple regression analysis in modelling of ground-level ozone and factors affecting its concentrations. *Environ. Modell. Softw.*, 20 (10): 1263-1271.
12. Wang, S. and F. Xiao, 2004. AHU sensor fault diagnosis using principal component analysis method. *Energy Build.*, 36(2): 147-160.
13. Sousa, S.I.V., F.G. Martins, M.C.M. Alvim-Ferraz and M.C. Pereira, 2007. Multiple linear regression and artificial neural networks based on principal components to predict ozone concentrations. *Environ. Modell. Softw.*, 22 (1): 97-103.

14. Shao, X., W. Wang, Z. Hou and W. Cai, 2006. A new regression method based on independent component analysis. *Talanta*, 69 (3): 676-680.
15. Koenker, R. and G. Basset, 1978. Regression quantiles. *Econometrica*, 44: 33-50.
16. Abdi, H., 2003. Partial Least Squares (PLS) Regression. In: *Encyclopedia of Social Sciences Research Methods* (eds M. Lewis-Beck, A. Bryman and T. Futing) pp. 1-7. Thousand Oaks: Sage.
17. Hayter, A.J., H.P. Wynn and W. Liu, 2005. Slope modified confidence bands for a simple linear regression model. *Stat. Methodol.*, 3 (2): 186-192.
18. Chaloulakou, A., M. Saisana and N. Spyrellis, 2003. Comparative assessment of neural networks and regression models for forecasting summertime ozone in Athens. *Sci. Total Environ.*, 313 (1-3): 1-13.
19. Gardner M.W. and S.R. Dorling, 2000. Statistical surface ozone models: an improved methodology to account for non-linear behaviour. *Atmos. Environ.*, 34 (1): 21-34.
20. Sousa, S.I.V., F.G. Martins, M.C. Pereira and M.C.M. Alvim-Ferraz, 2006. Prediction of ozone concentrations in Oporto city with statistical approaches, *Chemosphere*, 64 (7): 1141-1149.