

Air quality prediction for work environments using ARDL models

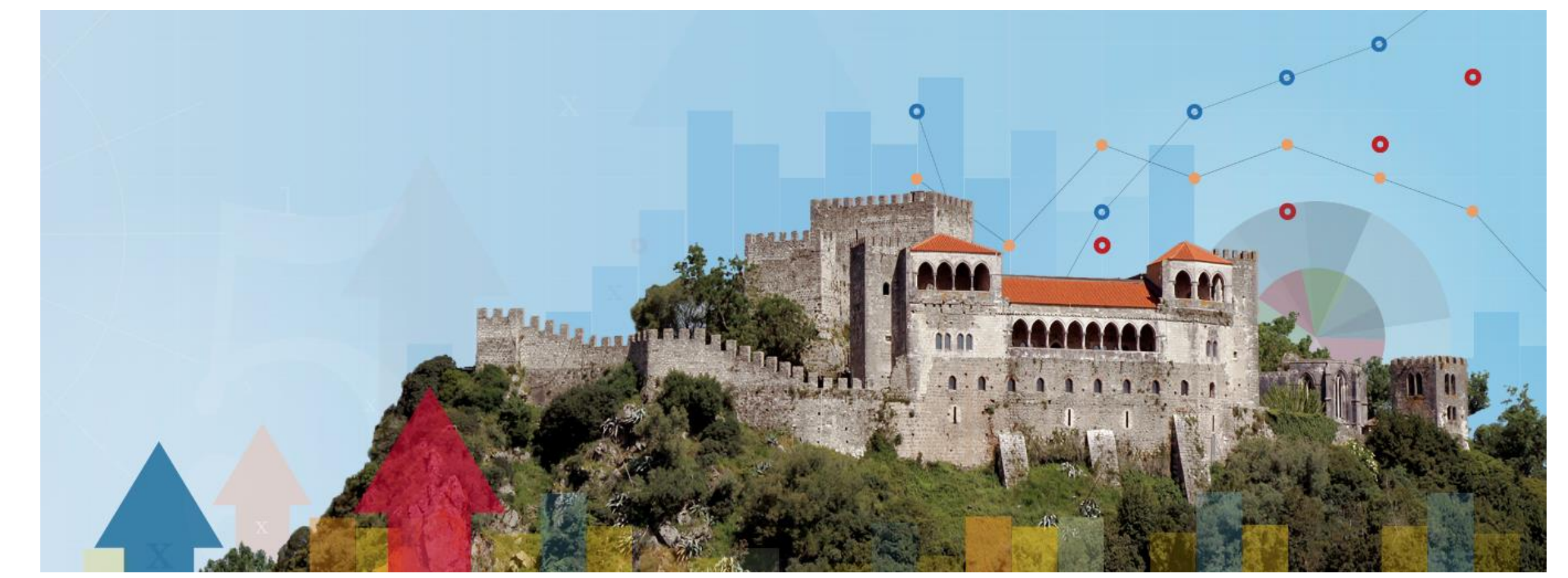
Fernando Batista^{1,2*} e Jorge Siopa¹

1: School of Technology and Management, Polytechnic of Leiria, Portugal

2: CDRSP - Centre for Rapid and Sustainable Product Development, Polytechnic of Leiria

e-mail: fernando.batista@ipleiria.pt, jorge.siopa@ipleiria.pt,

*corresponding author



Introduction

Air quality in urban areas affects the health and well-being of the population. Knowing the importance of this problem, it is necessary in the first phase to monitor all the atmospheric parameters that can affect air quality and in the second phase, to create predictive models. With these predictive models it is possible to analyze the importance of the atmospheric parameters that influence the gases that degrade air quality. Using a Python package for an Autoregressive Distributed Lag (ARDL) analysis of time series, it was possible to characterize this importance.

Aim

In an attempt to assess well-being at Campus 2 of the School of Technology and Management of the Polytechnic of Leiria, air quality monitoring data obtained by the on-site Mobile Air Quality Monitoring Unit was collected and a study was carried out using this data to characterise the importance of certain atmospheric parameters using a predictive time series analysis tool, ARDL.

Method

This ARDL method is widely used in economics to estimate forecasts. The time series in the model are divided into endogenous and exogenous variables. The parameters that characterize the atmospheric conditions are always assumed to be exogenous variables. The gases in the atmosphere that influence air quality are assumed to be endogenous variables or exogenous variables.

The time series obtained for the gases were ozone O_3 , carbon monoxide CO , hydrocarbons (NO , NO_x and NO_2) and particles PM_{10} and $PM_{2.5}$. For the atmospheric conditions we have precipitation (mmH_2O), global radiation (R_{global}), temperature ($Temp$), humidity ($Humid$), atmospheric pressure ($Press$), wind speed ($Vwind$) and direction ($Dwind$). To try to characterize car traffic, a time series of rush hour was estimated ($Rhour$). This time series consisting of "0" and "1" was suggested in the work done in this field. In order to draw more robust conclusions with time series, we need a lot of data from many days. Although in this study we only have data from 17 days, it was possible to reach some conclusions.

Various statistical tests were carried out on the time series to assess their stationarity, such as ADF and KPSS. The relationship importance between the two types of variables was also assessed.

The ARDL method uses previous data on the variables and manages to capture the seasonal effect of the series. Equation shows the various components of the iterative process characterized by lags.

$$Y_t = \delta + \sum_{i=1}^{P-1} \gamma_i S_{[(mod(t,P)+1)=i]} + \sum_{p=1}^A \phi_p Y_{t-p} + \sum_{k=1}^M \sum_{j=0}^{Q_k} \beta_{k,j} X_{k,t-j} + \epsilon_t$$

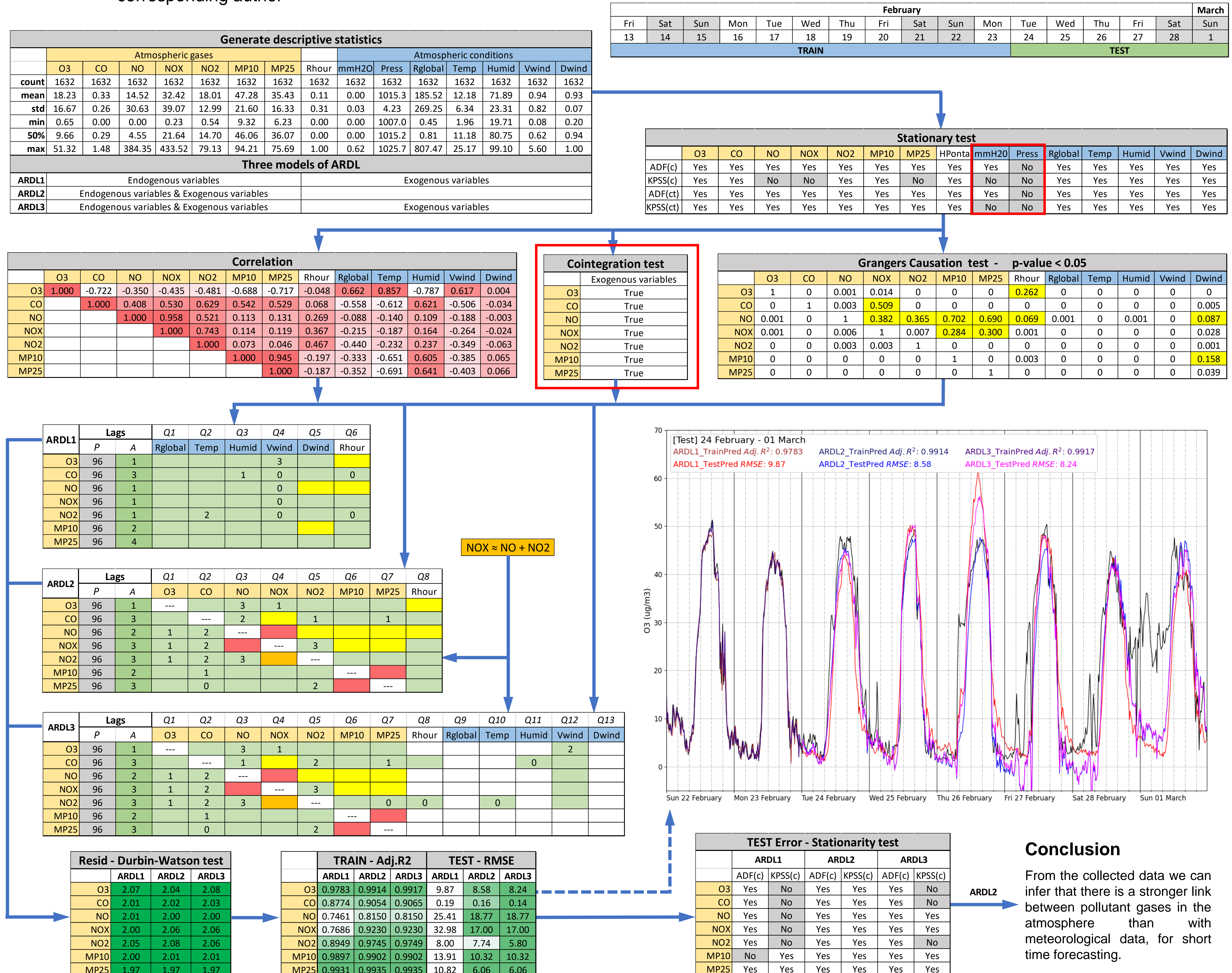
where δ is constant, $\gamma_i S_{[(mod(t,P)+1)=i]}$ capture seasonal shifts, P is the period of the seasonality, A is the lag length of the endogenous variable, M is the number of exogenous variables, X_k , Q_k is included the lag length of X_k and ϵ_t is a white noise shock.

Statistical metrics such as the "adjusted square error", $Adj. R^2$, the "root-mean-square error", $RMSE$ and stationarity methods were used to reach conclusions.

Three types of ARDL (3 models) were carried out, the first with atmospheric conditions as exogenous variables, the second with atmospheric gases as exogenous variables and the third was a combination of atmospheric gases and the atmospheric conditions identified in the first model. The three models performed well on the training data, but when the RMSE and the stationarity of the predictive residuals were evaluated, model 2, which only has gas variables, had the best fit and was the only one to have stationary predictive residuals.

Tools

IDE: PyCharm, Python package: Numpy, Pandas, Sklearn, Statsmodels and Matplotlib.



Conclusion

From the collected data we can infer that there is a stronger link between pollutant gases in the atmosphere than with meteorological data, for short time forecasting.