

Automatic covariates selection in dynamic regression models with application to COVID-19 evolution

Ana X. Ezquerro¹, Germán Aneiros², and Manuel Oviedo de la Fuente³

¹ University of A Coruña, ana.ezquerro@udc.es

² Grupo MODES, Departamento de Matemáticas, University of A Coruña, CITIC
german.aneiros@udc.es

³ Grupo MODES, Departamento de Matemáticas, University of A Coruña, CITIC
manuel.oviedo@udc.es

Abstract

This work introduces a new approach in time-series analysis field for automatic covariates selection in dynamic regression models. Based on [1] and [2] previous study, a forward-selection method is proposed for adding new significant covariates from a given set. This algorithm has been implemented and optimized in R as a package, and it has been applied to multiple simulations to validate its performance. Finally, the obtained results from the IRAS database of Catalonia are presented to analyze the COVID-19 evolution.

1 Introduction

In time-series analysis, the well-known dynamic regression models allow formally modelling the dependence between a set of covariates and a dependent variable considering the intrinsic temporal component of all participant variables. Thus, this type of regression models are of widespread application in diverse scenarios where it is desired to analyze the effect of recollected data in a time series of interest.

Formally, dynamic linear regression models define the linear dependence between a stochastic process Y_t (the dependent variable) and a set of processes $\mathcal{X} = \{X_t^{(1)}, X_t^{(2)}, \dots, X_t^{(m)}\}$ (candidates for regressor variables) in times non-greater than t :

$$Y_t = \beta_0 + \beta_1 X_{t-r_1}^{(1)} + \beta_2 X_{t-r_2}^{(2)} + \dots + \beta_m X_{t-r_m}^{(m)} + \eta_t \quad (1)$$

where $r_i \geq 0$, for $i = 1, \dots, m$, and $\eta_t \sim \text{ARMA}(p, q)$.

In this work we formally introduce a new algorithm to select covariates which significantly influence the behavior of a dependent variable. Due to the impact of COVID-19 around the world, we use this method to formalize and study the relation of the COVID-19 evolution in Catalonia (Spain) with the flu syndrome, COVID-19 vaccination and other recollected variables from the IRAS database.

2 Methodology

Following the definition in 1, [1] proposed a method named *prewhitening* for removing spurious correlation (false linear correlation) between two processes X_t and Y_t (where one of them is not white noise and/or the other is not stationary) by analyzing the cross correlation function

$$\rho_k(\ddot{X}_t, \ddot{Y}_t) = \frac{\text{Cov}(\ddot{X}_t, \ddot{Y}_{t-k})}{\sigma_{\ddot{X}_t} \sigma_{\ddot{Y}_t}} \text{ where } \sigma_{Z_t} \text{ denotes the standard deviation of a stochastic process } Z_t$$

and \ddot{X}_t and \ddot{Y}_t are obtained via some linear filter application to X_t and Y_t ensuring one of

them is white noise and the other is a stationary process. Specifically, [1] proposes a real linear correlation between X_t and Y_t if exists some k where $\rho_k(\ddot{X}_t, \ddot{Y}_t)$ is statistically significant. This method is applied to obtain the optimal lags of each regressor in 1, considering the condition of k being less or equal than 0.

Our approach iteratively adds dependent processes to a model by checking if a significant correlation (as in [1]) exists between a new process (candidate for regressor variable) and the residuals η_t of a simpler model.

Let Y_t be the stochastic dependent process and \mathcal{X} be the set of processes that might act as regressor variables in the model (candidates), and an information criterion (IC) for model evaluation. Our method proceeds as follows:

1. Initialization. Consider the process $\tilde{Y}_t = Y_t$ that will be used to check the existence of linear correlation between Y_t and each $X_t \in \mathcal{X}$ with [1] method, $\nu = \infty$ the value of the IC corresponding to the best model with 1 form, $\mathcal{X}^{(s,r)}$ the set of selected covariates paired with their respective optimal lags and $\mathcal{X}^{(s)}$ the set of selected covariates (with no lag information). Let $\mathcal{M}(\mathcal{Z})$ be the fitted dynamic regression model regarding Y_t where \mathcal{Z} is the set of covariates paired with their optimal lags:

$$\mathcal{M}(\mathcal{Z}) := Y_t = \beta_0 + \sum_{(Z_t, r) \in \mathcal{Z}} \beta^{(Z_t, r)} Z_{t-r} + \eta_t$$

where $\beta^{(Z_t, r)}$ is obtained via some estimation.

2. Iterative selection. For each $X_t \in \mathcal{X} - \mathcal{X}^{(s)}$, obtain the optimal lag where the maximum linear cross correlation between X_t and \tilde{Y}_t occurs (via [1] method). Consider the process $X_t^{\text{best}} \in \mathcal{X} - \mathcal{X}^{(s)}$ that minimizes and improves ν value, based on the selected IC, by including it in the model with its optimal lag ($r_{X_t^{\text{best}}}$):

$$X_t^{\text{best}} = \arg \min_{X_t \in \mathcal{X} - \mathcal{X}^{(s)}} \left\{ \text{criteria} \left(\mathcal{M} \left(\mathcal{X}^{(s,r)} \cup \{(X_t, r_{X_t})\} \right) \right) \right\} \quad (2)$$

conditioned to $\text{criteria}(\cdot) < \nu^1$. If X_t^{best} exists, consider $\mathcal{X}^{(s,r)} = \mathcal{X}^{(s,r)} \cup \{(X_t^{\text{best}}, r_{X_t^{\text{best}}})\}$, $\tilde{Y}_t = \eta_t$ and $\nu = \text{criteria}(\mathcal{X}^{(s,r)})^2$. Repeat this step until no process $X_t \in \mathcal{X} - \mathcal{X}^{(s)}$ can be added to the model, i.e. X_t^{best} does not exist.

3. Finalization. If the errors η_t of $\mathcal{M}(\mathcal{X}^{(s,r)})$ are not stationary and no model with $\eta_t \sim \text{ARMA}(p, q)$ and $\mathcal{X}^{(s,r)}$ covariates can be adjusted, consider the regular differentiation of all data (dependent variable and regressor candidates) and return to (1). Otherwise, it is proven that $\mathcal{M}(\mathcal{X}^{(s,r)})$ with stationary errors defines the significant correlation between the set of $\mathcal{X}^{(s)}$ regressor variables and the dependent process Y_t .

This algorithm was implemented in R programming language. The step 2 was optimized by parallelizing the fit of independent models of each candidate in \mathcal{X} . Dickey-Fuller test is used for checking processes stationary, Ljung-Box to check the independence, Shapiro-Wilks and Jarque-Bera tests for normality and t-test for zero mean of ARIMA residuals.

¹for simplicity, we denote the expression in $\text{criteria}()$ in 2 as \cdot

²once X_t^{best} has been added to the model

Table 1: Information about the dynamic regression model constructed via selection of multiple vaccination variables to model COVID19 evolution

Covariate	Lag	Coefficient est. (s.e)
vac4565	-3	-0.0410 (0.0057)
vac6580	-2	-0.0468 (0.0120)
vac1845	-6	-0.0901 (0.0047)
vac1218	Not included in the model	
vac80	Not included in the model	
residuals	ARIMA(4, 0, 0)	$\phi_1 = 2.0816(0.0810)$ $\phi_2 = -1.2837(0.1152)$ $\phi_4 = 0.1919(0.0432)$

3 Simulation results

In order to validate the performance of our selection method, we simulate multiple scenarios where a time series Y_t was artificially constructed with other variables (introduced with their respective coefficients and lags as in 1), which were added to a set of candidates along with more variables which do not influence in the construction of Y_t . The algorithm was tested when the residuals of the model η_t were stationary and non-stationary.

4 COVID-19 application

Our approach was tested in the IRAS (*acute respiratory infections*) database of Catalonia (Spain) in order to analyze the evolution of COVID-19 and the impact of other variables, such as the vaccination progress and influenza confirmed cases. In addition, individual data was aggregated by age ranges and Health Areas to study the correlation between groups and their influence in the global evolution.

Table 1³ resumes the algorithm trace and the order of covariates addition to the model.

References

- [1] Jonathan D Cryer and Kung-Sik Chan. *Time series analysis: with applications in R*, chapter 11. 2, 2008.
- [2] Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*, chapter 9. OTexts, 2018.

³where **vac1218**, **vac1845**, **vac4565**, **vac6580** denote the vaccination in population from 12, 18, 45 and 65 up to 18, 45, 65 and 80 years (exclusive), respectively, and **vac80** denotes the vaccination in population from 80 years.