

Distributed-lag linear structural equation models in R: the **dlsem** package

Alessandro Magrini
Dep. Statistics, Computer Science, Applications
University of Florence, Italy
<magrini@disia.unifi.it>

dlsem version 2.2 – 05 March 2018

Contents

1	Introduction	1
2	Theory	2
3	Installation	5
4	Illustrative example	5
4.1	Specification of the model code	6
4.2	Specification of control options	6
4.3	Parameter estimation	7
4.4	Assessment of causal effects	10
4.5	Model comparison	12
5	Final remarks	14

1 Introduction

Structural causal models (SCMs, Pearl, 2000) are a mathematical framework describing the behaviour of a multivariate system, and represent one of the prevalent methodologies for causal inference in contemporary applied sciences. Markovian SCMs are a special case where the joint probability distribution of the considered variables can be factored according to a directed acyclic graph. *Distributed-lag linear structural equation models* (DLSEMs) are Markovian SCMs, where each factor of the joint probability distribution is a distributed-lag linear regression model (see, for example, Judge *et al.*, 1985, Chapters 9-10). DLSEMs account for temporal delays in the dependence relationships among the considered variables and allow to perform dynamic causal inference by assessing causal effects at different time lags. As such, they are suitable to investigate the effect of an external impulse on a multivariate system through time. Econometrics and epidemiology are two of the main fields of application for DLSEMs.

Package **dlsem** implements inference functionalities for DLSEMs with several types of constrained lag shapes. This vignette is structured as follows. In Section 2, theory on the DLSEM is presented. In Section 3, instructions for the installation of the **dlsem** package are provided. In Section 4, the practical use of **dlsem** is illustrated through a simple impact assessment problem. Section 5 includes final remarks and considerations on future development of the package.

2 Theory

Distributed-lag linear regression Lagged instances of one or more covariates may be included in the linear regression model to account for temporal delays in their influence on the response:

$$y_t = \beta_0 + \sum_{j=1}^J \sum_{l=0}^{L_j} \beta_{j,l} x_{j,t-l} + \epsilon_t \quad (1)$$

where y_t is the value of the response variable at time t , $x_{j,t-l}$ is the value of the j -th covariate at l time lags before t , and ϵ_t is the random error at time t uncorrelated with the covariates and with ϵ_{t-1} . The set $(\beta_{j,0}, \beta_{j,1}, \dots, \beta_{j,L_j})$ is denoted as the *lag shape* of the j -th covariate and represents its regression coefficient (in the remainder, simply ‘coefficient’) at different time lags.

The model in Formula 1 has the disadvantage that a parameter is required for each lagged instance of a covariate, and lagged instances of the same covariate tend to be highly correlated. The Almon’s polynomial lag shape (Almon, 1965) overcomes these limitations by forcing the coefficients for lagged instances of the same covariate to follow a polynomial of order Q :

$$\beta_{j,l} = \begin{cases} \phi_0 & l = 0 \\ \sum_{q=0}^Q \phi_q l^q & \text{otherwise} \end{cases} \quad (2)$$

For instance, for $q = 2$ we have that $\beta_{j,l} = \phi_0 + \phi_1 l + \phi_2 l^2$. Unfortunately, the Almon’s polynomial lag shape may show multiple modes and coefficients with different signs, thus entailing problems of interpretation. Constrained lag shapes (Judge *et al.*, 1985, Chapters 9-10) overcome this issue. Package `dlsem` includes the *endpoint-constrained quadratic* lag shape:

$$\beta_{j,l} = \begin{cases} \theta_j \left[-\frac{4}{(b_j - a_j + 2)^2} l^2 + \frac{4(a_j + b_j)}{(b_j - a_j + 2)^2} l - \frac{4(a_j - 1)(b_j + 1)}{(b_j - a_j + 2)^2} \right] & a_j \leq l \leq b_j \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

the *quadratic decreasing* lag shape:

$$\beta_{j,l} = \begin{cases} \theta_j \frac{l^2 - 2(b_j + 1)l + (b_j + 1)^2}{(b_j - a_j + 1)^2} & a_j \leq l \leq b_j \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

and the *gamma* lag shape:

$$\beta_{j,l} = \theta_j (l + 1)^{\frac{\delta_j}{1-\delta_j}} \lambda_j^l \left[\left(\frac{\delta_j}{(\delta_j - 1) \log(\lambda_j)} \right)^{\frac{\delta_j}{1-\delta_j}} \lambda_j^{\frac{\delta_j}{(\delta_j - 1) \log(\lambda_j)} - 1} \right]^{-1} \quad (5)$$

$$0 < \delta_j < 1 \quad 0 < \lambda_j < 1$$

The endpoint-constrained quadratic lag shape is zero for a time lag $l < a_j$ or $l > b_j$, and symmetric with mode equal to θ_j at lag $(a_j + b_j)/2$. The quadratic decreasing lag shape decreases from value θ_j at lag a_j to value 0 at lag $b_j + 1$ according to a quadratic function. The gamma lag shape is positively skewed with mode equal to θ_j at lag $\frac{\delta_j}{(\delta_j - 1) \log(\lambda_j)}$. Value a_j is denoted as the *gestation lag*, value b_j as the *lead lag*, and value $b_j - a_j$ as the *lag width*. A static coefficient is obtained if $a_j = b_j = 0$. Since it is not expressed as a function of a_j and b_j , the gamma lag shape cannot reduce to a static coefficient, but the corresponding values of a_j and b_j may be computed through numerical approximation. For these three lag shapes it holds:

$$\begin{aligned} \beta_{j,l} > 0 &\iff \theta_j > 0 \\ \beta_{j,l} < 0 &\iff \theta_j < 0 \end{aligned} \quad \forall a_j \leq l \leq b_j \quad (6)$$

and we refer to the *lag sign* as the sign of parameter θ_j .

A linear regression model with constrained lag shapes is linear in parameters $\beta_0, \theta_1, \dots, \theta_J$, provided that the values of $a_1, \dots, a_J, b_1, \dots, b_J$ are known. Thus, one may use ordinary least squares to estimate parameters $\beta_0, \theta_1, \dots, \theta_J$ for several models with different values of $a_1, \dots, a_J, b_1, \dots, b_J$, and then select the one with the minimum value of the Akaike Information Criterion (AIC, Akaike, 1974) or the Bayesian Information Criterion (BIC, Schwarz, 1978)¹.

¹ Note that neither the response variable nor the covariates must contain a trend in order to obtain unbiased estimates with ordinary least squares (Granger and Newbold, 1974). A reasonable procedure is to sequentially apply differentiation to all variables until the Augmented Dickey-Fuller test (Dickey and Fuller, 1981) rejects the hypothesis of unit root for all of them.

Structural causal models Structural causal models (SCMs) were developed by Pearl (2000) in the context of causal inference. They are rooted to path analysis (Wright, 1934) and simultaneous equation models (Haavelmo, 1943; Koopmans *et al.*, 1950). A SCM consists of a tuple $\{\mathbf{V}, \mathbf{U}, \Omega_{\mathbf{V}}, \Omega_{\mathbf{U}}, \mathbf{f}, \mathbb{P}_{\mathbf{U}}\}$, where:

- $\mathbf{V} = \{V_1, \dots, V_J\}$ is a set of endogenous variables;
- $\Omega_{\mathbf{V}} = \Omega_{V_1} \times \dots \times \Omega_{V_J}$ is the cartesian product of the domains of variables in \mathbf{V} ;
- $\mathbf{U} = \{U_1, \dots, U_K\}$ is a set of unobserved variables;
- $\Omega_{\mathbf{U}} = \Omega_{U_1} \times \dots \times \Omega_{U_K}$ is the cartesian product of the domains of variables in \mathbf{U} ;
- $\mathbf{f} : \Omega_{\mathbf{V}} \times \Omega_{\mathbf{U}} \rightarrow \Omega_{\mathbf{V}}$ is a measurable function;
- $\mathbb{P}_{\mathbf{U}}$ is a probability measure on $\Omega_{\mathbf{U}}$.

Markovian SCMs (Pearl, 2000, Chapter 3) are a special case where \mathbf{f} is acyclic and variables in \mathbf{U} are each other independent. In a Markovian SCM, the following factorization of the joint probability distribution of variables in \mathbf{V} holds:

$$p(v_1, \dots, v_J) = \prod_{j=1}^J p(v_j \mid \Pi_j = \pi_j) \quad (7)$$

where Π_j is the set of variables in \mathbf{V} such that, for $j > 1$, V_j is independent of variables in $\{V_1, \dots, V_{j-1}\} \setminus \Pi_j$, given variables in Π_j . This means that the joint probability distribution of variables in \mathbf{V} can be factored according to conditional independence relationships holding among them disregarding variables in \mathbf{U} . Pearl (2000, pages 12 and following) shows that these conditional independence relationships are encoded into a directed acyclic graph (DAG) such that Π_j is the parent set of V_j , $\forall j = 1, \dots, J$. For example, in the Markovian SCM associated to the DAG in Figure 1, it holds:

$$p(v_1, v_2, v_3, v_4) = p(v_1) p(v_2 \mid v_1) p(v_3 \mid v_1) p(v_4 \mid v_2, v_3) \quad (8)$$

and, for example, V_4 is independent of V_1 given V_2 and V_3 .

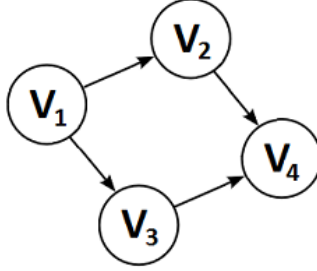


Figure 1: An example of directed acyclic graph.

Let $\text{do}(V_i = v_i)$ denote an intervention setting the value of V_i to v_i . Then, in a Markovian SCM it holds:

$$p(v_1, \dots, v_J \mid \text{do}(V_i = v_i)) = \prod_{j \neq i} p(v_j \mid \pi_j) \mid_{V_i = v_i} \quad (9)$$

where $\mid_{V_i = v_i}$ indicates that $p(v_i \mid \pi_i)$ is replaced by value v_i . This formula, called *truncated factorization* (Pearl, 2000, Section 3.2), allows to compute the effect of an intervention from the (pre-intervention) distribution in Formula 7, that is to predict such effect from non-experimental (observational) data. In a Markovian SCM, the effect of $\text{do}(V_i = v_i)$ on V_j , called *causal effect* of V_i on V_j , is given by the following expression (see Pearl, 2000, page 70 and following):

$$p(V_j = v_j \mid \text{do}(V_i = v_i)) = \sum_{\pi_i} p(V_j = v_j \mid V_i = v_i, \Pi_i = \pi_i) p(\Pi_i = \pi_i) \quad (10)$$

where Π_i is the parent set of V_i .

In a linear parametric formulation of SCMs (linear Markovian SCMs), each factor $p(v_j | \pi_j)$ of the joint probability distribution in Formula 7 is the linear regression model where V_j is the response variable and variables in Π_j are the covariates. For example, in the linear Markovian SCM associated to the DAG in Figure 1, $p(v_4 | v_2, v_3)$ is the linear regression model where V_4 is the response variable and V_2 and V_3 are the covariates.

In a linear Markovian SCM, the computation of causal effects involves the coefficients of the regression models only, without the need of Formula 10. Let $\text{do}(\Delta V_i = 1)$ be an intervention changing the value of V_i by a unit. Under such intervention:

- if V_i is parent of V_j , the *direct* causal effect of V_i on V_j is equal to the coefficient of V_i in the regression model of V_j ;
- the causal effect of V_i on V_j through a multi-edge directed path $\langle V_i, \dots, V_j \rangle$ connecting V_i to V_j , called *indirect* causal effect of V_i on V_j through $\langle V_i, \dots, V_j \rangle$, is equal to (see, for example, Wright, 1934):

$$e(\langle V_i, \dots, V_j \rangle) = \prod_{k: V_k \in \langle V_i, \dots, V_j \rangle \wedge k \neq i} b_{k|k-1} \quad (11)$$

where $b_{k|k-1}$ is the coefficient of V_{k-1} in the regression model of V_k . The direct causal effect of V_i on V_j and all the indirect causal effects of V_i on V_j are denoted as *pathwise* causal effects of V_i on V_j ;

- the *overall* causal effect of V_i on V_j is equal to the sum of all the pathwise causal effects of V_i on V_j .

For example, in the linear Markovian SCM associated to the DAG in Figure 1, there are two directed paths connecting V_1 to V_4 : $\langle V_1, V_2, V_4 \rangle$ with pathwise causal effect $b_{2|1} \cdot b_{4|2}$, and $\langle V_1, V_3, V_4 \rangle$ with pathwise causal effect $b_{3|1} \cdot b_{4|3}$. Thus, the overall causal effect of V_1 on V_4 is equal to $b_{2|1} \cdot b_{4|2} + b_{3|1} \cdot b_{4|3}$.

Distributed-lag linear structural equation models Distributed-lag linear structural equation models (DLSEMs) are Markovian SCMs where each factor of the joint probability distribution in Formula 7 is a distributed-lag linear regression model (see, for example, Judge *et al.*, 1985, Chapters 9-10). The DAG of a DLSEM would involve all the possible temporal instances of each variable in \mathbf{V} . Here, for simplicity, a static DAG is still used for a DLSEM, where the edge $\langle V_i, V_j \rangle$ exists if and only if there exists at least one time lag where the coefficient of variable V_i in the regression model of variable V_j is non-zero. Causal effects at different time lags in a DLSEM are defined as follows:

- if V_i is parent of V_j , the *direct* causal effect of V_i on V_j at lag l is equal to the coefficient of V_i at lag l in the regression model of V_j ;
- Let $\langle V_{d_0}, \dots, V_{d_m} \rangle$, $d_0 = i$ and $d_m = j$, be a directed path composed of m edges connecting V_i to V_j , and \mathcal{L}_m be the set of all the possible ordered m -uples of time lags such that their sum is equal to l . The *indirect* causal effect of V_i on V_j through such path at lag l is equal to:

$$e(\langle V_{d_0}, \dots, V_{d_m} \rangle; d_0 = i, d_m = j) = \sum_{(l_1, \dots, l_m) \in \mathcal{L}_m} \prod_{k=1}^m b_{d_k|d_{k-1}}^{(l_k)} \quad (12)$$

where $b_{d_k|d_{k-1}}^{(l_k)}$ is the coefficient of $V_{d_{k-1}}$ at lag l_k in the regression model of V_{d_k} ;

- the *overall* causal effect of V_i on V_j at lag l is equal to the sum of all the pathwise causal effects of V_i on V_j at lag l .

The causal effects just defined are denoted as *instantaneous* causal effects, as they are evaluated at a single time lag. The *cumulative* causal effect at a pre-specified time lag, say l , is obtained by summing all the instantaneous causal effects at each time lag up to l . A *pathwise causal lag shape* is the set of causal effects associated to a path at different time lags. An *overall causal lag shape* is the set of the overall causal effects of a variable on another one at different time lags.

3 Installation

Before installing `dlsem`, you must have installed R version 2.1.0 or higher, which is freely available at <http://www.r-project.org/>.

To install the `dlsem` package, type the following in the R command prompt:

```
> install.packages("dlsem")
```

and R will automatically install the package to your system from CRAN. In order to keep your copy of `dlsem` up to date, use the command:

```
> update.packages("dlsem")
```

The latest version of `dlsem` is 2.2.

4 Illustrative example

The practical use of package `dlsem` is illustrated through a simple impact assessment problem denoted as “industrial development problem”. The objective is to test whether the influence through time of the number job positions in industry (proxy of the industrial development) on the amount of greenhouse gas emissions (proxy of pollution) is direct and/or mediated by the amount of private consumption. The DAG for the industrial development problem is shown in Figure 2. The analysis will be conducted on the dataset `industry`, containing simulated data for 10 imaginary regions in the period 1983-2015.

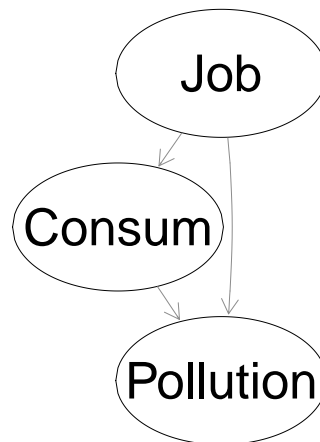


Figure 2: The DAG for the industrial development problem. ‘Job’: number of job positions in industry. ‘Consum’: private consumption index. ‘Pollution’: amount of greenhouse gas emissions.

```
> data(industry)
> summary(industry)
```

	Region	Year	Population	GDP
1	: 32	Min. :1983	Min. : 4771649	Min. : 97119
2	: 32	1st Qu.:1991	1st Qu.: 8310737	1st Qu.: 186783
3	: 32	Median :1998	Median :25381874	Median : 463942
4	: 32	Mean :1998	Mean :32368547	Mean : 727735
5	: 32	3rd Qu.:2006	3rd Qu.:56273337	3rd Qu.:1307044
6	: 32	Max. :2014	Max. :78308254	Max. :1883702
(Other):	128			
	Job	Consum	Pollution	
Min.	: 34.77	Min. : 37.35	Min. : 3161	

1st Qu.:105.07	1st Qu.: 87.88	1st Qu.: 7536
Median :137.03	Median :108.47	Median : 25320
Mean :127.61	Mean :108.17	Mean : 32202
3rd Qu.:152.68	3rd Qu.:124.85	3rd Qu.: 47109
Max. :200.83	Max. :211.16	Max. :101441

4.1 Specification of the model code

The first step to build a DLSEM with the `dlsem` package is the definition of the model code, which includes the formal specification of the regression models. The variables for which a regression model is specified are called *endogenous* variables. The other variables are referred as *exogenous* variables (not to be confused with the unobserved disturbances).

The model code must be a list of formulas, one for each regression model. In each formula, the response and the covariates must be quantitative variables², and operators `quec.lag()`, `qdec.lag()` and `gamma.lag()` may be employed to specify, respectively, an endpoint-constrained quadratic, a quadratic decreasing or a gamma lag shape. Operators `quec.lag()` and `qdec.lag()` have three arguments: the name of the covariate to which the lag shape is applied, the gestation lag (a_j) and the lead lag (b_j). Operator `gamma.lag()` has three arguments: the name of the covariate to which the lag shape is applied, parameter δ_j and parameter λ_j . If none of these two operators is applied to a covariate, it is assumed that its coefficient is equal to 0 for time lags greater than 0 (no lag shape). The group factor and exogenous variables must not appear in the model code (see Subsection 4.3 for the way to include them). The specification of regression models with no endogenous covariates may be omitted from the model code (for example, one could avoid to specify the regression model for the number of job positions). In this problem, all lag shapes are assumed to be endpoint-constrained quadratic lag shapes between 0 and 15 time lags:

```
> indus.code <- list(
+   Job ~ 1,
+   Consum~quec.lag(Job,0,15),
+   Pollution~quec.lag(Job,0,15)+quec.lag(Consum,0,15)
+ )
```

4.2 Specification of control options

The second step to build a DLSEM with the `dlsem` package is the specification of control options. Control options are distinguished into global (applied to all regression models) and local (model-specific) options. Global control options must be a named list with one or more of the following components:

- **adapt**: a logical value indicating if adaptation of lag shapes must be performed, that is parameters of lag shapes must be chosen on the basis of fit to data. Default is **FALSE**, meaning no adaptation;
- **max.gestation**: the maximum gestation lag for all lag shapes. If not provided, it is taken as equal to **max.lead** (see below);
- **max.lead**: the maximum lead lag for all lag shapes. If not provided, it is computed accordingly to the sample size;
- **min.width**: the minimum lag width for all lag shapes. It cannot be greater than **max.lead**. If not provided, it is taken as 0;
- **sign**: the lag sign for all lag shapes, that may be either '+' for positive or '-' for negative. If not provided, adaptation will disregard the lag sign.

Local control options must be a named list containing one or more among the following components:

- **adapt**: a named vector of logical values, where each component must have the name of one endogenous variable and indicate if adaptation of lag shapes must be performed for the regression model of that variable;

² Qualitative variables may be included only as exogenous variables, as described in Subsection 4.3.

- **max.gestation**: a named list. Each component of the list must have the name of one endogenous variable and be a named vector. Each component of the named vector must have the name of one covariate in the regression model of the endogenous variable above and include the maximum gestation lag for its lag shape;
- **max.lead**: the same as **max.gestation**, with the exception that the named vector must include the maximum lead lag;
- **min.width**: the same as **max.gestation**, with the exception that the named vector must include the minimum lag width;
- **sign**: the same as **max.gestation**, with the exception that the named vector must include the lag sign (either '+' for positive or '-' for negative).

Local control options have no default values, and global ones are applied in their absence. If some local control options conflict with global ones, only the former are applied.

Suppose that one wants to perform adaptation with the following constraints for all lag shapes: (i) maximum gestation lag of 3 years, (ii) maximum lead lag of 15 years, (iii) minimum lag width of 5 years, (iv) positive lag sign. Control options for these constraints may be expressed in several ways. The most simple solution is to specify only global control options, as the constraints hold for all regression models:

```
> indus.global <- list(adapt=T,max.gestation=3,max.lead=15,min.width=5,sign="+")
> indus.local <- list()
```

In alternative, one may specify only local control options, by repeating them for each regression model:

```
> indus.global <- list()
> indus.local <- list(
+   adapt=c(Consum=T,Pollution=T),
+   max.gestation=list(Consum=c(Job=3),Pollution=c(Job=3,Consum=3)),
+   max.lead=list(Consum=c(Job=15),Pollution=c(Job=15,Consum=15)),
+   min.width=list(Consum=c(Job=5),Pollution=c(Job=5,Consum=5)),
+   sign=list(Consum=c(Job="+"),Pollution=c(Job="+",Consum="+"))
+ )
```

or both local and global control options:

```
> indus.global <- list(adapt=T,min.width=5)
> indus.local <- list(
+   max.gestation=list(Consum=c(Job=3),Pollution=c(Job=3,Consum=3)),
+   max.lead=list(Consum=c(Job=15),Pollution=c(Job=15,Consum=15)),
+   sign=list(Consum=c(Job="+"),Pollution=c(Job="+",Consum="+"))
+ )
```

4.3 Parameter estimation

Once the model code and control options are specified, parameter estimation can be performed using the command `dlsem(·)`. The user may indicate a single group factor (just one) to argument **group** and one or more exogenous variables to argument **exogenous**. By indicating the group factor, one intercept for each level of the group factor will be estimated in each regression model, in order to explain the variability due to differences between groups. By indicating exogenous variables, they will be included as non-lagged covariates in each regression model, in order to eliminate cross-sectional spurious effects. Each exogenous variable may be either qualitative or quantitative and its coefficient in each regression model is 0 for time lags greater than 0 (no lag). The user may decide to apply the logarithmic transformation to all strictly positive quantitative variables by setting argument **log** to **TRUE**, in order to interpret each coefficient as an elasticity (percentage increase in the value of the response variable for 1% increase in the value of a covariate). Before parameter estimation, differentiation is performed until the hypothesis of unit root is rejected by the Augmented Dickey-Fuller test for all quantitative variables³, and each

³ If the group factor is specified, the panel version of the Augmented Dickey-Fuller test proposed by Levin *et al.* (2002) is used instead.

missing value is replaced by its conditional mean computed through the Expectation-Maximization algorithm (Dempster *et al.*, 1977)⁴. In this problem, the region is indicated as the group factor, while population and gross domestic product are indicated as exogenous variables. Also, the logarithmic transformation is requested, and global and local control options are provided to arguments `global.control` and `local.control`, respectively:

```
> indus.mod <- dlsem(indus.code,group="Region",exogenous=c("Population","GDP"),
+   data=industry,global.control=indus.global,local.control=indus.local,log=T)
```

Checking stationarity...

Order 1 differentiation performed

Starting estimation...

Estimating regression model 1/3 (Job)

Estimating regression model 2/3 (Consum)

Estimating regression model 3/3 (Pollution)

Estimation completed

The results of command `dlsem(.)` is an object of class `dlsem`. Among the components of `dlsem` objects, we found:

- `estimate`: a list of objects of class `lm`, one for each regression model;
- `call`: a list containing the call for each regression model after eventual adaptation of lag shapes;
- `data.used`: data after eventual logarithmic transformation and differentiation.

The `summary` method for class `dlsem` returns the summary of the estimation:

```
> summary(indus.mod)
```

ENDOGENOUS PART

Response: Job

-

Response: Consum

	Estimate	Std. Error	t value	Pr(> t)
quec.lag(Job, 0, 5, Region)	0.1006394	0.01783725	5.642089	4.589874e-08 ***

Response: Pollution

	Estimate	Std. Error	t value	Pr(> t)
quec.lag(Job, 1, 8, Region)	0.1048006	0.03008457	3.483532	5.989626e-04 ***
quec.lag(Consum, 1, 6, Region)	0.2320105	0.03660783	6.337729	1.339514e-09 ***

EXOGENOUS PART

Response: Job

	Estimate	Std. Error	t value	Pr(> t)
Population	-2.015755	0.36919466	-5.45987	1.004944e-07 ***
GDP	-1.274005	0.03253314	-39.16023	1.591909e-119 ***

Response: Consum

	Estimate	Std. Error	t value	Pr(> t)
Population	0.8397265	0.30729012	2.732683	6.735972e-03 **
GDP	-0.8165645	0.02710312	-30.128064	1.096637e-84 ***

⁴ The computation is performed after eventual logarithmic transformation and differentiation, assuming group-specific means and time-invariant covariance matrix. Qualitative variables cannot contain missing values. Each quantitative variable must have at least 3 observed values if the group factor is not specified, otherwise it must have at least 3 observed values per group.

Response: Pollution

	Estimate	Std. Error	t value	Pr(> t)
Population	-0.5335639	0.32247211	-1.654605	9.945701e-02 .
GDP	0.1342472	0.02965881	4.526384	9.908715e-06 ***

INTERCEPTS

Response: Job

	Estimate	Std. Error	t value	Pr(> t)
Region1	-0.027108664	0.002403134	-11.280545	8.189303e-25 ***
Region2	-0.014868387	0.002401561	-6.191135	1.975106e-09 ***
Region3	-0.014228172	0.002401629	-5.924383	8.639991e-09 ***
Region4	-0.005320298	0.002403060	-2.213968	2.758788e-02 *
Region5	-0.008833821	0.002401537	-3.678402	2.784066e-04 ***
Region6	-0.015622725	0.002401342	-6.505831	3.260886e-10 ***
Region7	-0.005154175	0.002401605	-2.146138	3.266936e-02 *
Region8	-0.027052095	0.002401793	-11.263293	9.395308e-25 ***
Region9	-0.046951445	0.002402163	-19.545484	2.514703e-55 ***
Region10	-0.023440072	0.002402647	-9.755938	1.077582e-19 ***

Response: Consum

	Estimate	Std. Error	t value	Pr(> t)
Region1	0.013228135	0.003105034	4.260222	2.905842e-05 ***
Region2	-0.009181367	0.002452433	-3.743779	2.255585e-04 ***
Region3	0.014910423	0.002369592	6.292400	1.413274e-09 ***
Region4	0.012261936	0.002143643	5.720139	3.065699e-08 ***
Region5	0.012591239	0.002189363	5.751097	2.609354e-08 ***
Region6	0.027006345	0.002425256	11.135464	1.319250e-23 ***
Region7	0.023946916	0.002133839	11.222454	6.881615e-24 ***
Region8	-0.014297098	0.003061892	-4.669367	4.962066e-06 ***
Region9	0.019452657	0.004455213	4.366268	1.860065e-05 ***
Region10	0.003490765	0.002834166	1.231673	2.192426e-01

Response: Pollution

	Estimate	Std. Error	t value	Pr(> t)
Region1	0.0181034164	0.005671781	3.1918397	1.624344e-03 **
Region2	0.0166945282	0.002993763	5.5764369	7.290975e-08 ***
Region3	0.0008710423	0.004745009	0.1835702	8.545228e-01
Region4	0.0038743955	0.003341414	1.1595078	2.475294e-01
Region5	-0.0047651007	0.003653564	-1.3042336	1.935420e-01
Region6	-0.0138551604	0.006254540	-2.2152164	2.778958e-02 *
Region7	-0.0133904268	0.004809974	-2.7838873	5.847590e-03 **
Region8	0.0294218841	0.004102569	7.1715751	1.164801e-11 ***
Region9	0.0029735558	0.008691559	0.3421200	7.325933e-01
Region10	0.0171095625	0.004253094	4.0228508	7.951079e-05 ***

GOODNESS OF FIT

R-squared: 0.8609

AIC: -4786.373

BIC: -4636.377

We see that the number of job positions in industry (Job) significantly influences, on one hand, the amount of private consumption (Consum) from 0 to 4 time lags and, on the other hand, the amount of greenhouse gas emissions (Pollution) from 2 to 6 time lags, while the amount of private consumption (Consum) significantly influences the amount of greenhouse gas emissions (Pollution) from 1 to 5 time lags. This result provides evidence that the influence of industrial development on pollution is both direct and mediated by private consumption.

The `plot` method for class `dlsem` displays the DAG of the model where each edge is coloured with respect to the sign of the estimated causal effect (green: positive, red: negative, light gray: not statistically significant):

```
> plot(indus.mod)
```

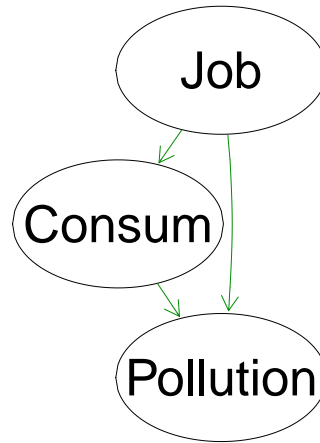


Figure 3: The DAG where each edge is coloured with respect to the sign of the estimated causal effect. Green: positive causal effect. Red: negative causal effect. Grey: not statistically significant causal effect (no such edges here).

The result is shown in Figure 3. Note that the DAG includes only the endogenous variables.

4.4 Assessment of causal effects

After parameter estimation is performed by means of command `dlsem()`, the command `causalEff()` can be used on the resulting object of class `dlsem` to compute all the pathwise causal lag shapes and the overall one connecting two variables. The main arguments of command `causalEff()` include the name of one or more variables generating the causal effect (argument `from`), and the name of the variable receiving the causal effect (argument `to`). Optionally, specific time to which computation should be focused may be provided to argument `lag`, otherwise the whole lag shapes will be considered. Also, the user may choose whether the instantaneous (argument `cumul` set to `FALSE`, the default) or the cumulative (argument `cumul` set to `TRUE`) causal effect must be returned. Only exogenous variables can be indicated as starting or ending variables. Note that, due to the properties of the multiple linear regression model, causal effects are net of the influence of the group factor and exogenous variables.

The cumulative causal effect of the number of job positions on the amount of greenhouse gas emissions may be obtained by means of the following code:

```
> causalEff(indus.mod,from="Job",to="Pollution",cumul=T)
```

```

$`Job*Consum*Pollution`
      estimate  lower 95%  upper 95%
0  0.000000000  0.000000000  0.000000000
1  0.005601519  0.002978257  0.008224781
2  0.024273250  0.017556697  0.030989803
3  0.062239103  0.049859194  0.074619011
4  0.121988641  0.102982419  0.140994863
5  0.200409911  0.174670963  0.226148858
6  0.287544654  0.255879838  0.319209470
7  0.365965924  0.329856126  0.402075722
8  0.425715462  0.386832409  0.464598516
9  0.463681315  0.423431571  0.503931059
10 0.482353046  0.441631154  0.523074937

```

```
11 0.487954565 0.447148267 0.528760863
12 0.487954565 0.447148267 0.528760863
```

```
$`Job*Pollution`
      estimate lower 95% upper 95%
0 0.00000000 0.00000000 0.00000000
1 0.0414027 0.01810801 0.06469739
2 0.1138574 0.06690547 0.16080937
3 0.2070135 0.13664579 0.27738119
4 0.3105202 0.21917950 0.40186096
5 0.4140270 0.30570041 0.52235354
6 0.5071830 0.38684280 0.62752328
7 0.5796378 0.45258023 0.70669529
8 0.6210405 0.49186516 0.75021576
9 0.6210405 0.49186516 0.75021576
10 0.6210405 0.49186516 0.75021576
11 0.6210405 0.49186516 0.75021576
12 0.6210405 0.49186516 0.75021576
```

```
$overall
      estimate lower 95% upper 95%
0 0.00000000 0.00000000 0.00000000
1 0.04700422 0.02108627 0.07292217
2 0.13813067 0.08450297 0.19175836
3 0.26925259 0.18666113 0.35184405
4 0.43250887 0.32250642 0.54251132
5 0.61443689 0.48096437 0.74790940
6 0.79472770 0.64361323 0.94584216
7 0.94560369 0.78369692 1.10751046
8 1.04675592 0.88051424 1.21299761
9 1.08472178 0.91815513 1.25128843
10 1.10339351 0.93671214 1.27007488
11 1.10899503 0.94229301 1.27569704
12 1.10899503 0.94229301 1.27569704
```

The output of command `causalEff(.)` is a list of matrices including point estimates and asymptotic confidence intervals for all the pathwise causal lag shapes and the overall one connecting the starting variables to the ending variable. Since the logarithmic transformation was applied to all quantitative variables, the resulting causal effects are interpreted as elasticities, that is, for a 1% of job positions more, greenhouse gas emissions are expected to grow by 0.61% after 5 years and by 1.11% after 10 years. The influence ends after 11 years, as the cumulative causal effects at 11 and 12 years are equal.

A pathwise or an overall causal lag shape can be displayed using the command `lagPlot(.)`. For instance, one may display the causal lag shape associated to each path connecting the number of job positions to the amount of greenhouse gas emissions:

```
> lagPlot(indus.mod,path="Job*Pollution")
> lagPlot(indus.mod,path="Job*Consum*Pollution")
```

or the overall causal lag shape of the number of job positions on the amount of greenhouse gas emissions:

```
> lagPlot(indus.mod,from="Job",to="Pollution")
```

The resulting graphics are shown in Figure 4. Note that a multi-edge pathwise causal lag shape is a mixture of different lag shapes, thus it may show an irregular aspect, like it is the case of the overall causal lag shape displayed in the lower panel of Figure 4.

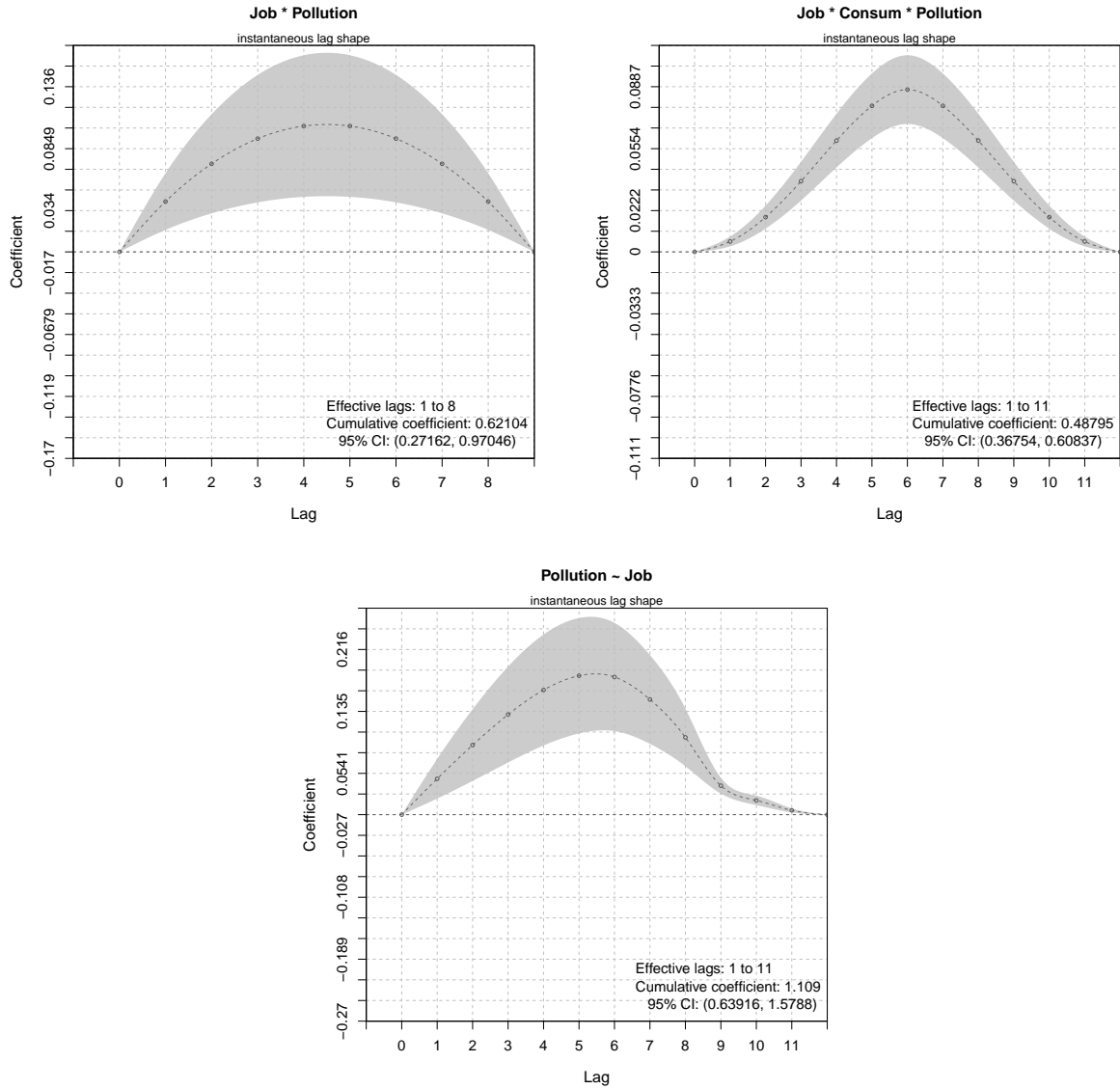


Figure 4: Pathwise causal lag shapes (upper panels) and the overall one (lower panel) connecting the number of job positions to the amount of greenhouse gas emissions. 95% asymptotic confidence intervals are shown in grey.

4.5 Model comparison

We now fit two alternative models for the industrial development problem, such that all lag shapes are quadratic decreasing and gamma lag shapes, respectively:

```
> # model 2: quadratic decreasing lag shapes
> indus.code_2 <- list(
+   Job ~ 1,
+   Consum~qdec.lag(Job,0,15),
+   Pollution~qdec.lag(Job,0,15)+qdec.lag(Consum,0,15)
+ )
> indus.mod_2 <- dlsem(indus.code_2,group="Region",exogenous=c("Population","GDP"),
+   data=industry,global.control=indus.global,local.control=indus.local,log=T)
```

```
Checking stationarity...
Order 1 differentiation performed
Starting estimation...
Estimating regression model 1/3 (Job)
```

```

Estimating regression model 2/3 (Consum)
Estimating regression model 3/3 (Pollution)
Estimation completed

> summary(indus.mod_2)$endogenous

$Job
NULL

$Consum
              Estimate Std. Error  t value    Pr(>|t|)
qdec.lag(Job, 0, 5, Region) 0.1108326 0.02698702 4.106885 5.455783e-05 ***

$Pollution
              Estimate Std. Error  t value    Pr(>|t|)
qdec.lag(Job, 2, 15, Region) 0.22047664 0.03171312 6.952221 1.113042e-10 ***
qdec.lag(Consum, 0, 5, Region) 0.09883959 0.05389121 1.834058 6.868144e-02 .

> # model 3: gamma lag shapes
> indus.code_3 <- list(
+   Job ~ 1,
+   Consum~gamma.lag(Job,0.5,0.5),
+   Pollution~gamma.lag(Job,0.5,0.5)+gamma.lag(Consum,0.5,0.5)
+ )
> indus.mod_3 <- dlsem(indus.code_3,group="Region",exogenous=c("Population","GDP"),
+   data=industry,global.control=indus.global,local.control=indus.local,log=T)

```

Checking stationarity...

Order 1 differentiation performed

Starting estimation...

Estimating regression model 1/3 (Job)

Estimating regression model 2/3 (Consum)

Estimating regression model 3/3 (Pollution)

Estimation completed

```

> summary(indus.mod_3)$endogenous

$Job
NULL

$Consum
              Estimate Std. Error  t value    Pr(>|t|)
gamma.lag(Job, 0.16, 0.36, Region) 0.213074 0.0620565 3.433548 0.0006942689

gamma.lag(Job, 0.16, 0.36, Region) ***

$Pollution
              Estimate Std. Error  t value    Pr(>|t|)
gamma.lag(Job, 0.92, 0.14, Region) 0.3322248 0.02637051 12.598346
gamma.lag(Consum, 0.88, 0.03, Region) 0.0931422 0.03770555 2.470251
              Pr(>|t|)
gamma.lag(Job, 0.92, 0.14, Region) 1.030076e-26 ***
gamma.lag(Consum, 0.88, 0.03, Region) 1.440266e-02 *

```

We see that the three models provide different results. Methods AIC and BIC for class `dlsem` can be used to compare them according to AIC and BIC, respectively:

```

> lapply(list(QUEC=indus.mod,QDEC=indus.mod_2,GAMMA=indus.mod_3),AIC)

$QUEC
      Job      Consum Pollution (overall)
-1781.636 -1603.661 -1401.076 -4786.373

$QDEC

```

```

      Job      Consum Pollution (overall)
-1781.6357 -1589.3165 -912.6081 -4283.5603

$GAMMA
      Job      Consum Pollution (overall)
-1781.636 -1655.876 -1182.643 -4620.154

> lapply(list(QUEC=indus.mod,QDEC=indus.mod_2,GAMMA=indus.mod_3),BIC)

$QUEC
      Job      Consum Pollution (overall)
-1733.060 -1553.811 -1349.505 -4636.377

$QDEC
      Job      Consum Pollution (overall)
-1733.0602 -1539.4669 -866.4805 -4139.0077

$GAMMA
      Job      Consum Pollution (overall)
-1733.060 -1605.498 -1133.168 -4471.726

```

The model with endpoint-constrained quadratic lag shapes has the best fit according to both AIC and BIC. Note that the fit for variable `Job` is constant in each model because it has no endogenous covariates.

5 Final remarks

Lag shapes included in the package may represent a large number of real-world lag structures: unimodal symmetric (with the endpoint-constrained quadratic lag shape), unimodal asymmetric (with the gamma lag shape) and skewed ones (with the quadratic decreasing lag shape). Nevertheless, additional lag shapes with further specific features may be added in future.

Parameter estimation in DLSEM cannot be performed in a single step unless gestation and lead lags are all known. Since complete search over all the possible models is infeasible for most real-world applications, a heuristic search based on forward model selection is currently implemented. Further development of the package may be directed towards the improvement of the search strategy.

Grouped data are currently managed through fixed effects estimation. In the future, random effects estimation may be implemented to enhance inference whenever the considered groups are a subset of the possible ones, or covariates with values constant within groups (second-level covariates) are involved.

Please, do not hesitate to contact me for questions, feedbacks or bug reports.

References

- H. Akaike (1974). A New Look at the Statistical Identification Model. *IEEE Transactions on Automatic Control*, 19: 716-723.
- S. Almon (1965). The Distributed Lag between Capital Appropriations and Net Expenditures. *Econometrica*, 33, 178-196.
- A. P. Dempster, N. M. Laird, and D. B. Rubin (1977). Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1): 1-38.
- D. A. Dickey, and W. A. Fuller (1981). Likelihood Ratio Statistics for Autoregressive Time Series with a Unit Root. *Econometrica*, 49: 1057-1072.
- C. W. J. Granger, and P. Newbold (1974). Spurious Regressions in Econometrics. *Journal of Econometrics*, 2(2), 111-120.
- G. G. Judge, W. E. Griffiths, R. C. Hill, H. Lutkepohl, and T. C. Lee (1985). The Theory and Practice of Econometrics. John Wiley & Sons, 2nd ed., New York, US-NY.

- T. Haavelmo (1943). The Statistical Implications of a System of Simultaneous Equations. *Econometrica*, 11(1): 1-12.
- T. C. Koopmans, H. Rubin, and R. B. Leipnik (1950). Measuring the Equation Systems of Dynamic Economics. In: T. C. Koopmans (ed.), *Statistical Inference in Dynamic Economic Models*, pages 53-237. John Wiley & Sons, New York, US-NY.
- A. Levin, C. Lin, and C. J. Chub (2002). Unit Root Tests in Panel Data: Asymptotic and Finite-Sample Properties. *Journal of Econometrics*, 108: 1-24.
- J. Pearl (2012). The Causal Foundations of Structural Equation Modelling. In: R. H. Hoyle (ed.), *Handbook of Structural Equation Modelling*, pages 68-91. Guilford Press, New York, US-NY.
- J. Pearl (2000). *Causality: Models, Reasoning, and Inference*. Cambridge University Press. Cambridge, UK.
- G. Schwarz (1978). Estimating the Dimension of a Model. *Annals of Statistics*, 6, 461-464.
- S. Wright (1934). The Method of Path Coefficients. *Annals of Mathematical Statistics*, 5(3): 161-215.