# EXAMINING CAUSE-EFFECT RELATIONS IN THE SOCIAL SCIENCES
## A Structural Causal Modelling Approach

The objective of this paper is to present a short overview of the *Structural Causal Modelling* (SCM) framework developed by the present authors in a series of articles spanning the last decade or so. The text is based on a presentation given at Statistics Netherlands in Heerlen on December 4, 2018 (Russo, Wunsch, Mouchart, 2019). The purpose is to explain how the SCM framework provides the tools to hypothesize, model, and test explanatory mechanisms. Our framework proves particularly useful in social science contexts, since it allows us to adopt an explicit causal perspective even when analyzing observational data. Social science experiments are notoriously difficult to carry out for ethical or practical reasons, and our approach allows social scientists to go beyond mere description and to propose a causal explanation even in the absence of experiments and interventions.

Guillaume Wunsch, Michel Mouchart, Federica Russo

As indicated by the name itself, structural models aim for an analysis of the structural relationships among variables and are based on field knowledge and on theoretical contributions. In econometrics or in social sciences, structural models have typically the form of a set of equations. In general, these approaches specify a statistical model without developing a detailed analysis in terms of recursivity and therefore do not end up with an explicit view of the underlying mechanism and sub-mechanisms. The SCM framework is not based on a system of equations, but on an analysis of multivariate distributions. Adopting an SCM approach means endorsing a particular view on modelling in general (the hypothetico-deductive methodology), and a specific stance on exogeneity, namely as a condition of separability of inference, on the one hand, and in interpreting marginal-conditional decompositions as sub-mechanisms, on the other hand. The construction of the statistical model is then deduced from the above approach.

In this paper, we focus on SCM as one possible perspective in quantitative social science research. There are of course other ways to analyze social phenomena, such as, for example, a systemic approach (Loriaux, 1994), agent-based modelling (Billari et al., 2007), or qualitative designs such as the case study. These other ways will not be discussed here.

### A Structural Causal Modelling (SCM) Framework

The *Structural Causal Modelling* approach may be viewed as a chapter in the domain of statistical modelling, where a statistical model is considered as a set of "reasonable" hypotheses concerning the *data generating process* (DGP) represented as a probability distribution. A probabilistic representation of the DGP is used to explain a phenomenon of interest. Such an explanation involves two compo-

nents: (i) a stochastic element that embodies what is not explained in the working of the DGP (*i.e.* of the mechanism), and (ii) a non-stochastic element, the characteristics – or parameters – of the probability distribution, that provides the nature of what is explained by the statistical model. Said differently, the statistical model provides only a partial explanation of the mechanism of the DGP. For more on the relationship between statistics, causality, and explanation, see e.g. the interesting paper published in *STAtOR* by Richard Starmans (2018).

Compared to structural models in econometrics or in social sciences, the present framework takes distance from the latter in several aspects. To begin with, our structural approach is based on a *hypothetico-deductive* (H-D) methodology. This means that a hypothesis is first formulated, a model developed and tested, and the results interpreted in order to confirm or disconfirm the initial hypothesis. H-D methodologies are widely used in science and are often associated with the falsificationist view of Karl Popper (1934, English translation 1959). However, in philosophy of science, hypothetico-deductivism has been developed much beyond the original Popperian approach detailing, among others, the role of background knowledge at the hypothesis formulation stage or the fact that we learn also from disconfirmed hypotheses—so models can be iteratively improved on, and we do not start each time from scratch. Other important methodological features are the following: Causal and structural; Recursive decomposition and DAG; Exogeneity and causation; Distributions rather than equations; Explanation and parametrization; Stability and invariance.

*Causal and structural*
Focusing on causal analysis, SCM depends upon reliable background information and evidence for assessing puta-

tive causes of outcomes and evaluating effects of causes, and more generally on the structure of relations between causes and outcomes. Background knowledge plays a crucial role at each stage of the H-D methodology. Firstly, causal attribution is often quite a difficult issue once a system becomes complex. Secondly, background knowledge typically involves theories concerning the domain of analysis, but also embraces a much wider scope, in particular involving previous results and preliminary analysis of data. It is on this basis that a preliminary hypothesis is formulated. Background knowledge is likewise involved in the process of developing a specific statistical model, where one makes important choices about parametrization, testing methods, etc. Finally, the results of tests are interpreted against available background knowledge.

*Recursive decomposition and DAG*
'Explaining' essentially means representing and decomposing a complex and global mechanism in terms of a set of simpler sub-mechanisms. The explanation is based here on a recursive decomposition of the joint distribution of the variables entering the statistical analysis. This recursive decomposition is equivalent to a systematic marginal-conditional decomposition according to a specific ordering of the variables. For example, if one considers a vector of variables, the joint distribution can be written as: $P(X_1, \ldots X_p) = P(X_1)P(X_2 \mid X_1) \ldots P(X_p \mid X_1 \ldots X_{p-1})$. Thus, the joint distribution is written as a product of conditional distributions where the conditioning variables form an increasing sequence and where each factor of this product represents a sub-mechanism. For this reason, *directed acyclic graphs* (DAGs) provide a privileged tool of representation (Pearl, 2000), though a DAG does not allow representing all particularities of a multivariate distribution nor of a recursive decomposition.

*Exogeneity and causation*
Associations among variables are not necessarily causal. They can be due to the presence of one (or more) *confounders*, the latter being common factors of the putative causes and outcomes. One should therefore control for confounders, in order to avoid making false causal claims. Under a suitable *exogeneity* condition of non-confounding, one can then view the conditioning variables as causing variables in the sub-mechanism where they appear. This is the reason why the structural model is called a causal model, because causation is relative to a particular model built with the purpose of eliciting causes.
*Distributions rather than equations*
The basic objects of analysis are sets (in product form)

of distributions rather than sets of equations. Equations are related, at best, to conditional expectations, although effects of causes may take other ways. For instance, in actuarial applications, the effect of some contracts may be more in the tails of the distributions than in the expectations. To give another example, the analysis of the determinants of fertility should not only focus on the average number of children per woman but also e.g. on women having no children and on those having large families.

*Explanation and parametrization*
In SCM, explanation is based on a recursive decomposition. As mentioned before, representing a DGP by a probability distribution implies that this representation leaves unexplained some part of the DGP, namely the stochastic component of the model. Therefore, the statistical explanation concerns the characteristics, or parameters, of the probability distributions. This fact raises the issue of the specification of the parametrization. Once a conditional distribution is deemed to represent a specific sub-mechanism, the role of the parametrization is to identify the operation of the sub-mechanism (more information in Mouchart and Orsi, 2016).

*Stability and invariance*
Considering as structural a mechanism underlying the workings of a DGP requires that the model enjoys suitable properties of stability, or invariance, under a specific class of interventions and of modifications of the environment. Indeed, a model that would be different for, say, each observation should not be considered as structural. Said differently, the issue here is to look for a proper separation between the incidental and the structural aspects of the DGP. From a statistical point of view, this issue is also that of properly defining the *population of reference*. One reason for this is that no model in the social sciences can pretend to be universal in time and in space. There are no *laws* here. It should be stressed that this stability, or invariance, regards both the ordering of the variables and the value of the parameters.

**An example**

Consider a study on the recourse to contraception in urban Africa, the example being taken from Gourbin et al. (2017). The cities are characterized by different levels of contraceptive prevalence, but also by the different effectiveness of the methods used. Several questions may be raised, the following two amongst others. Firstly, what is
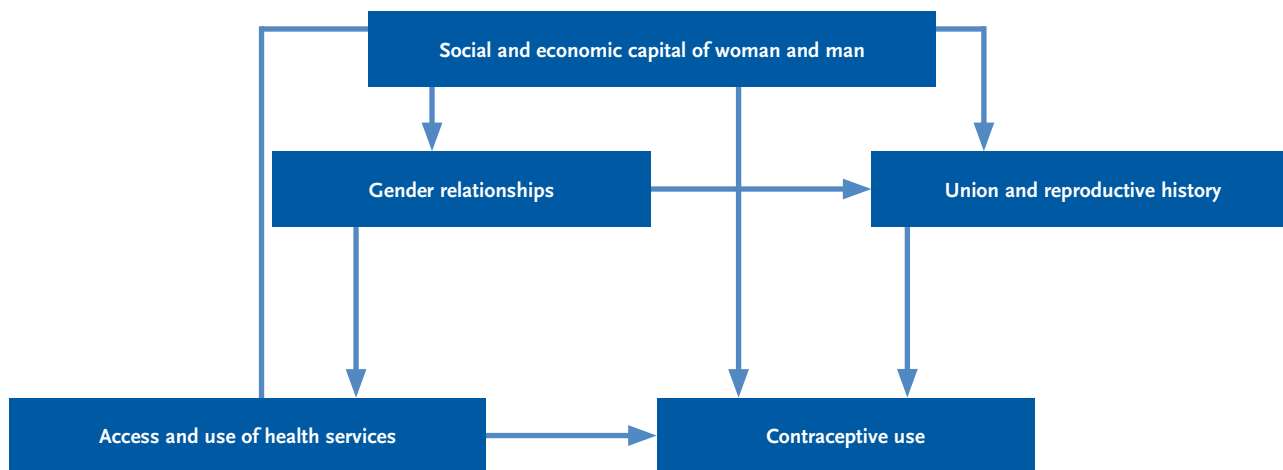
Figure 1. Conceptual model of contraceptive use in Africa (reproduced from Gourbin et al., 2017)

the hierarchical ordering of causal relationships among the individual factors involved in the use of contraception in the urban populations considered? Secondly, as education is a major factor of fertility transition, are two main indirect pathways that have been proposed in the literature (a union-reproductive path and a socio-cultural one) – leading from women's education to contraceptive use – confirmed by the data?

Most analyses of contraceptive use have had resort to statistical methods that do not take into account a possible causal ordering among the variables, implicitly assuming that all the putative determinants just have a direct effect on the dependent variable. However, the impact of these various factors on the use of contraception

can be direct or indirect, meaning in the latter case that the effect of some putative causes can be *mediated* by one or more intermediate factors. To answer the questions raised above, the SCM approach allows researchers to propose an *explanatory mechanism* for the outcome of interest, composed of various sub-mechanisms, and subsuming in particular the distinction between mediators, moderators, and confounding variables.

Based on background knowledge relating to contraceptive use and fertility, the following conceptual framework can be proposed (see Figure 1).

To test this conceptual model, one needs to obtain measurable indicators for each of the concepts in Figure 1. Using existing survey data* for the cities concerned,
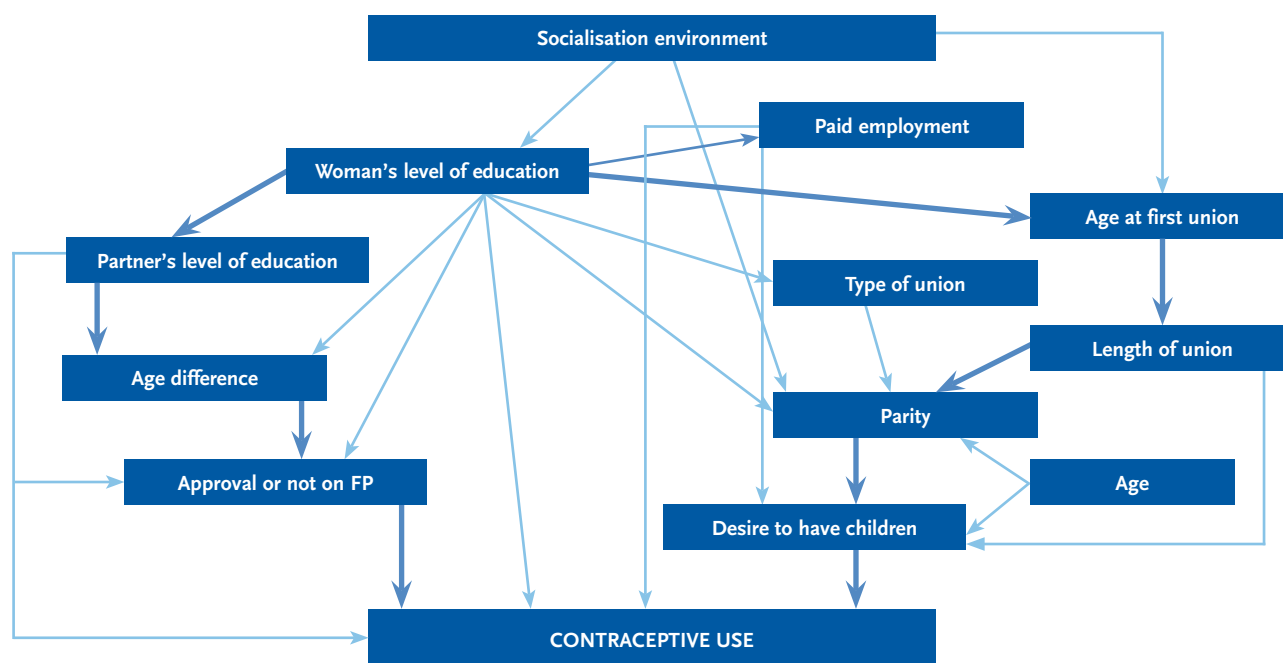


Figure 2. Operational model of contraceptive use in Africa (reproduced from Gourbin et al., 2017)

the authors have proposed the operational model presented in Figure 2. This figure is actually a *directed acyclic graph* (DAG) where the variables are ordered according to the various putative sub-mechanisms. All variables being in categorical format, the operational model has then been tested with the survey data available, using sequential logistic regressions. Concerning the two main indirect paths (in addition to the direct path) leading from education to contraceptive use proposed in the literature and drawn in **bold** in Figure 2, the data have confirmed the union-reproductive indirect path (on the right in the graph) but not the socio-cultural one (on the left). One should point out, however, that the individuals concerned were interviewed at a same moment in time and not over their lifetime. Results refer therefore to inter-individual differences obtained at the time of the surveys providing the data, and not to life-course differentials as could be derived from retrospective or prospective data.

### Some remarks on data, to conclude

As the structural model presented here aims at representing causal relations among variables, the latter should be ordered according to their causal priority, which implies *inter alia* a temporal ordering of causes and effects. This can be based, for example, on retrospective or prospective population surveys, on specific registries and other forms of permanent registration of individual events, such as a national register, etc. Many sources of data actually refer to the same individuals. If each individual receives at birth a personal identification number (PIN), data from multiple sources can be linked together. One can thus examine for an individual e.g. the move from good health to ill health, then to chronic disease, disability, and finally to death, possibly also taking into account various characteristics of the individual (such as education and employment) and their change over time. If individual longitudinal data are available, the causal model can also take into account reverse causation and feedback effects, by time-ordering the variables. One must however consider the fact that an event is often the result of a temporally prior decision-making process, based on the preferences, values, beliefs, emotions, of the agents in possible interaction with others. Data on the decision-making process are unfortunately most often unavailable. Contrary to the time-ordering of events, that of the various decision-making processes is thus difficult to specify.

Due to insufficient background knowledge or to a lack of information on the temporal sequence of events, it happens that the variables cannot be causally ordered. In this case, an exploratory analysis of the data and especially of the so-called *Big Data* (in the sense of very large structured and unstructured data sets) can possibly be helpful in revealing changing characteristics over time and suggesting the temporal sequence of events. An exploratory data approach is never a substitute for sound causal modelling, such as the framework presented in this paper, but it can usefully inform it, especially when background knowledge on the topic of interest is scant.

\* No indicators were however available for the use of health services.

REFERENCES
Billari, F. C., Prskawetz, A., Aparicio Diaz, B., & Fent T. (2007). The "Wedding-Ring": An agent-based marriage model based on social interaction. *Demographic Research*, 17/3, 59–82, doi: 10.4054/DemRes.2007.17.3.
Gourbin, C., Wunsch, G., et al. (2017). Direct and indirect paths leading to contraceptive use in urban Africa. *Revue Quetelet / Quetelet Journal*, 5(1), 33–71.
Loriaux, M. (1994). Des causes aux systèmes: la causalité en question, chap. 2 in R. Franck (Ed.), *Faut-il chercher aux causes une raison?*, Vrin, Paris, 41–86.
Mouchart, M., & Orsi R. (2016). Building a Structural Model: Parameterization and structurality. *Econometrics*, 4(2), doi: 10.3390/econometrics4020023.
Pearl, J. (2000). *Causality. Models, Reasoning, and Inference*. Cambridge: Cambridge University Press (revised and enlarged in 2009).
Popper, K. (1934, English translation 1959). *The Logic of Scientific Discovery*. London: Hutchinson.
Russo, F., Wunsch, G., & Mouchart, M. (2019). Causality in the Social Sciences. A structural causal modelling framework. *Quality & Quantity*, 53(5), 2575–2588.
Starmans R. (2018). Statistiek en causaliteit: voortschrijdende liaison of moeizame samenspraak, *STAtOR, 20*(4), 28–34.

GUILLAUME WUNSCH is Emeritus Professor of Demography at the University of Louvain (UCLouvain) and a member of the Royal Academy of Belgium.
E-mail: g.wunsch@uclouvain.be

MICHEL MOUCHART is Emeritus Professor of Statistics at the University of Louvain (UCLouvain).
E-mail: michel.mouchart@uclouvain.be

FEDERICA RUSSO is Professor of Philosophy of Science at the University of Amsterdam.
E-mail: federica.russo@gmail.com