



RedNHE

Red Nacional de
Investigadores
en Economía

Causality by Vote: Aggregating Evidence on Causal Relations in Economic Growth Processes

Manuel de Mier (UNS)

Fernando Delbianco (UNS/INMABB-CONICET)

Fernando Tohmé (UNS/INMABB-CONICET)

Luisina Patrizio (UNS)

Facundo Rodriguez (UNS)

Mauro Romero Stéfani (UNS)

DOCUMENTO DE TRABAJO N° 260

Julio de 2023

Los documentos de trabajo de la RedNIE se difunden con el propósito de generar comentarios y debate, no habiendo estado sujetos a revisión de pares. Las opiniones expresadas en este trabajo son de los autores y no necesariamente representan las opiniones de la RedNIE o su Comisión Directiva.

The RedNIE working papers are disseminated for the purpose of generating comments and debate, and have not been subjected to peer review. The opinions expressed in this paper are exclusively those of the authors and do not necessarily represent the opinions of the RedNIE or its Board of Directors.

Citar como:

de Mier, Manuel, Fernando Delbianco, Fernando Tohmé, Luisina Patrizio, Facundo Rodriguez y Mauro Romero Stéfani (2023). Causality by Vote: Aggregating Evidence on Causal Relations in Economic Growth Processes. Documento de trabajo RedNIE N°260.

Causality by Vote: Aggregating Evidence on Causal Relations in Economic Growth Processes

Manuel de Mier¹, Fernando Delbianco^{1,2,3}, Fernando Tohmé^{1,2}, Luisina Patrizio¹,
Facundo Rodriguez¹, and Mauro Romero Stéfani¹

¹Department of Economics, Universidad Nacional del Sur

²INMABB-CONICET, Bahía Blanca, Argentina

*Corresponding author: fernando.delbianco@uns.edu.ar, ORCID 0000-0002-1560-2587

July 2023

Abstract

In this paper we investigate the performance of five causality-detection methods and how their results can be aggregated when multiple units are considered in a panel data setting. The aggregation procedure employs voting rules for determining which causal paths are identified for the sample population. Using simulated and real-world panel data, we show the performance of these methods in detecting the correct causal paths in comparison to a benchmark that comprises a standard representation of growth processes as *ground truth* model. We find that the results may be better when only simulated, instead of real-world, data are analyzed. While this may suggest that the methods presented here are currently incapable of detecting causal links, it is plausible that the ground “truth” may incorporate false relations.

Keywords: Granger causality, Transfer Entropy, Stochastic Causality, LiNGAM, Ground Truth, Economic Growth.

JEL Classification: C18, C43, O47.

1 Introduction

The notion of *causal relation* has been, until recently, only used in heuristic terms in empirical studies. That is, it has been applied solely to offer intuitions behind mathematical or statistical relationships within models and datasets. However, in the last decades, there has been a shift as researchers seek to provide precise and formal characterizations of the previously vague conception. The literature now presents a plethora of approaches to both the treatment of causality in

already available networks of relations among variables, as well as the detection of causal relations in data (Peters et al., 2017; Schölkopf et al., 2021). But this diversity comes at the price of a lack of consensus among the different methods when no agreed-on *ground truth* model exists.

In the context of a panel data structure, when the same causal detection method is applied on each individual time series, it could produce dissimilar results. This is particularly worrisome when the phenomenon under study is assumed to be common to all the units in the panel, despite the differences that might arise when each corresponding time series is analyzed in isolation. In more precise terms, the question is, given variables X and Y , to detect whether a causal relation between them exists for all the panel, by checking its validity for each individual i . But a problem arises if this is true for *some* individuals and not for others.

In Econometrics, the standard concept of causality was introduced by Granger (1969), defined in terms of a time series X and another Y . Holtz-Eakin et al. (1988) presented the first version of Granger causality for panel data. But the assertion of the existence of a causal relation between X and Y was predicated under the assumption of the *homogeneity* of the relation, i.e., that the specifics of the causal model was the same for all the individuals in the panel. Dumitrescu and Hurlin (2012) lifted this assumption by allowing heterogeneity among the units and the causal relations. The test statistic they developed is obtained by averaging the Wald statistics of the individual Granger estimates. Juodis et al. (2021) notice that the asymptotic properties of this estimator make it prone to the emergence of a correlation between the estimates and the errors when the number of periods is small while the number of individuals is large. They address it by carrying out the averaging of the Wald statistics based on bias-free estimators.

The key idea behind these extensions of Granger causality to panel data is that the information drawn from running the causality detection method on each individual is then aggregated. In this paper we adopt this idea and extend it to other techniques of causality detection. The difference resides in the method in which the individual estimations is aggregated to find a consensus among them.

A possible way of achieving some degree of consensus among these results is by *voting*. That is, by including a causal relation among variables only when it is detected in a majority (or supermajority) of time series. In this paper, we explore the aggregation of the outcomes of causality detection procedures on economic growth data. The idea is to obtain a network of causally related variables reflecting the main features of growth processes.

An important proviso is that we apply this only on *pairs* of variables. The reason is that, as shown by Papana et al. (2021), causality detection methods like those to be considered here are hard to apply to multivariate time series. If the relations involve more than two variables, these methods are hampered by the existence of *mediators*, *confounders* and *colliders* (Glymour et al., 2016). Thus, it is not convenient to consider more than two variables at a time.

The voting approach to the aggregation of causal evidence is scalable. That is, once a method agrees on some causal link between two variables, we can see this as casting a vote for it in an election towards the construction of an ensemble model. In this way, the agreement obtained in the specific methods can be extended to a consensus among them when a plurality of the methods endorse a link. The causal links in the ensemble model can be deemed as more robust than those just agreed on by the specific methods.

Another advantage of this aggregation procedure can be seen in the light of the meta-analysis conducted by Bruns and Stern (2019) that shows that, in the case of Granger causality, the lags in the specification of the causal relations are obtained by overfitting the model. This indicates that the specification is usually cherry-picked. We address this issue by considering a number of different lags in a relation. Then, the existence of a causal relationship obtains by a weighted voting process where each unit casts votes for the different lags and the total number of its votes it gets is weighed by $\frac{1}{L}$, where L is the number of lags under consideration. This gives us a number for the *relation*, independently of the number of lags. Then, aggregating these numbers, if their sum passes the criterion of acceptance under voting, the relation among the variables is included in the causal graph.

One of the main reasons justifying the adoption of this methodology is the possibility of obtaining empirical evidence upon which future theoretical models can be developed, inverting the usual top-down approach to modeling in Economics. Nevertheless, this, as all bottom-up approaches, has a downside: its sensitivity to the quality of the data. That is, if the datasets are scarce or somehow imprecise, the relations detected may be erroneous while the correct ones may be missed.

The main contribution of this paper resides in the introduction of voting as a method of aggregation of causal information. This methodology facilitates assessing both the positive and negative aspects of the different methods. We also discuss the convenience of using ensemble methods obtained by aggregating the results of the individual procedures of causality detection.

The remainder of the paper is organized as follows. Section 2 describes the causality detection techniques, the vote aggregation procedure, and the data on which they are applied. Section 3 presents the results over synthetic and real-world data. Finally, Section 4 concludes. Additional information is contained in Appendix A.

2 Data and methodology

2.1 Data

The dataset used in this study is extracted from the Penn World Table (PWT) database (Feenstra et al., 2015). Our sample covers 62 developed and developing countries for the period 1961–2019. We work with 11 time series of annual frequency. The variables with their corresponding description are detailed in Table 1. The list of countries from which they were drawn can be found in Appendix A.

Table 1: Variables used in this study (drawn from PWT)

Variable	Description
GDP per capita growth rate (GROWTH)	Expenditure-side real GDP at chained PPPs growth rate
Population growth rate (POP)	Population (in millions) growth rate
Human capital (HK)	Human capital index
Total factor productivity (TFP)	TFP level at current PPPs (USA=1)
Real interest rate (IRR)	Real internal rate of return
Depreciation (DELTA)	Average depreciation rate of the capital stock
Exchange rate (XR)	National currency/USD (market+estimated)
Investment (INV)	Gross capital formation as a percentage of GDP
Government spending (GOV)	Government consumption as a percentage of GDP
Openness (OPEN)	Total merchandise trade as a percentage of GDP
Terms of trade (TOT)	Price level of exports/Price level of imports (USA 2017=1)

The choice of this panel for our analysis is justified by the need to consider a phenomenon like economic growth, which is usually conceived country-wise, with reduced across-border influences. Accordingly, we can analyze the causal relations among variables for the individual countries in the panel, which we later aggregate via voting, constituting an appropriate arena for the evaluation of the performance of our proposal.¹

¹The result of this exercise can be depicted as a *causal graph* As shown by Raunig (2023), this is particularly useful for the analysis of economic policies.

2.2 Causality detection techniques

In this paper we focus on two classes of methods of causality detection.² One is that of methods based on autoregressive models. The other is that of those that rely on structural causal learning techniques. The former methods are specific for time series analyses where a *causal* relation between variables X and Y is understood as indicating that past values of X provide unique information to predict or explain future values of variable Y . Conversely, any structural causal learning method makes assumptions about the functional form, breaking the symmetry between the variables, and enabling the identification of causal effects.

In this paper, we adopt two autoregressive-based methods, namely Granger Causality and Transfer Entropy. In turn, we apply three structural causal learning methods. One is Stochastic Causality and the other two are based on the linear non-Gaussian acyclic model (LiNGAM): ICA-LiNGAM and Direct-LiNGAM.

Granger causality is based on the premise that the causes precede the effects (Granger, 1969). To claim that one variable causes another, past values of an independent variable in a regression must have a significant effect on the dependent variable, given the past values of the latter. This estimation method not only allows us to detect causality but also its direction, which can be unidirectional or bidirectional. It should be emphasized that Granger causality is only applicable to linear structures. Therefore, if there are nonlinearities in the data, the true causal relationships between variables might not be accurately captured.

Since the Granger estimation requires that the time series are stationary, we proceed by differentiating the entire dataset if the Dickey-Fuller unit root test rejects the hypothesis of stationarity of the original series. In the case of simulated data, we run the Granger test with one lag for each pair of variables, while for the real-world growth dataset, we obtain estimations for one to four-time lags. We use for this task the R package *lmtest*. In the vote aggregation procedure, each of the lagged specifications casts a vote if the p-value is less than 0.05. Each of these votes, four in total, counts as 1/4 of a single vote.

Transfer Entropy is a non-parametric method based on Information Theory that uses the concept of *information transfer* to detect causal relationships (Schreiber, 2000). The entropy of a variable measures its degree of uncertainty or randomness. Thus, Transfer Entropy from X to Y

²We do not consider methods based on the independences among variables, like PC or PCMCI (Runge et al., 2019; Spirtes and Glymour, 1991).

measures the reduction in uncertainty (entropy) about the values of Y gained from the knowledge of the values of X , given the past values of Y . It is based on the same general principles of Granger Causality, but Transfer Entropy does not require any assumptions about the underlying structure of the data (Barnett et al., 2009). Therefore, we use Transfer Entropy to capture possible nonlinear causal relations that could not be detected with Granger Causality. Specifically, we apply the Shannon version of Transfer Entropy, using the R package *RTransferEntropy*. In the aggregation stage, a causal relation gets one vote if the p-value of the estimation is less than 0.05.

The Stochastic Causality approach is based on the information obtained by running a kernel regression of Y against some variables, including X to obtain a model $Y = f(X)$ and a kernel regression to get $X = g(Y)$. The results are evaluated according to three empirical criteria and are summarized in a unanimity index that determines whether a causal path exists among X and Y (Vinod, 2017, 2019). The first criterion is based on the application of the Wu-Hausman endogeneity test to choose between $Y = f(X)$ and $X = g(X)$. The second criterion selects the model with the smallest absolute value of the residuals. Finally, the third criterion chooses the model with the higher prediction accuracy, using the general measure of correlation for asymmetry of Zheng et al. (2012). Then, a unanimity index, ranging from -100 to 100, measures the degree of consensus among the three criteria regarding the direction of causality. The index equals 100 if all the criteria agree on the same causal direction. We adopt here a conservative threshold of 90 in absolute value and a significant p-value to indicate the presence of a causal relationship. A shortcoming of this method is that it does not detect bidirectional causal paths. Therefore, the directed acyclic graph (DAG) representing the outcome of the method tends to have less directed edges than those generated by other techniques. In our analyses, we use the R package *generalCorr* to produce the estimations in this approach.

LiNGAM (Linear Non-Gaussian Model) is a statistical model that assumes a linear relationship between variables with a non-normally distributed error term. In this setting, the asymmetry between the residuals in the regressions of X on Y and Y on X allows to distinguish cause from effect. The reason is that the regression residuals are independent of the predictor only for the correct causal direction. This approach was first introduced by Shimizu et al. (2006) with the ICA-LiNGAM algorithm, which uses independent component analysis (ICA) for the estimation of the model. Afterward, Shimizu et al. (2011) developed the Direct-LiNGAM algorithm, which improves upon ICA-LiNGAM by ensuring the convergence of the procedure of inference in a specific number of steps and with a known level of complexity. The R package *rlingam* was used

to implement both versions of the method.

2.3 Vote aggregation

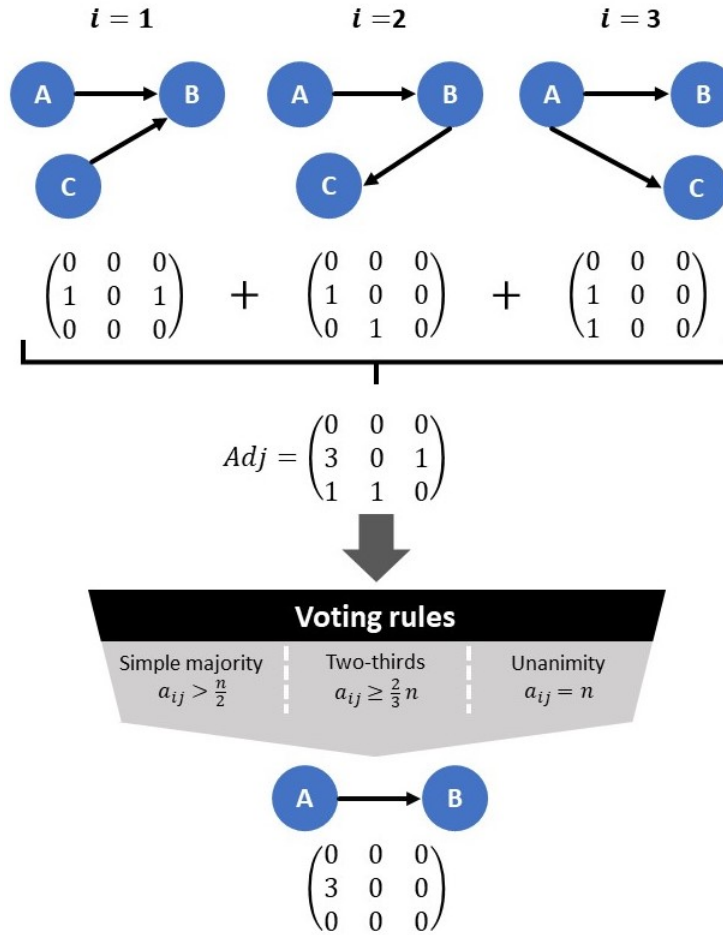
We run each method at the country level, disregarding the panel structure of the dataset. This is, on one hand, to detect patterns that correspond to single economies. But more importantly, with the exception of Granger causality, the remaining causality detection techniques to be applied are not defined for panel data and thus force us to run our analysis country-wise. Consequently, we obtain a large class of results for each method that we seek to aggregate into a unified causal structure valid for a relevant subset of the sample of countries.

A useful way to do this is through a voting system, which involves considering each country as a voter who casts a vote on the presence or absence of a causal relationship between a pair of variables. Then, the existence of a causal path between these variables is determined according to a voting rule that indicates how the votes are aggregated. Among the large class of voting rules, the most commonly used are *unanimity*, the *simple majority* rule, and the *two-thirds* rule. Each of these procedures implies a certain degree of robustness in the final results obtained. The unanimity rule yields a causal path only if that relation is detected in all the countries of the sample. The two-thirds rule, instead, requires that at least two-thirds of the sample share a causal link to accept it, while the simple majority rule requires an agreement of more than half of the sample.

Thus, the application of the simple majority rule mandates that a causal path is included only if it is detected in at least 32 countries among the 62 in the sample. The two-thirds rule, in turn, accepts a causal relationship only if it receives at least 42 votes.

Figure 1 exemplifies the aggregation procedure for three individuals and three variables. Each individual runs the causal detection procedure, obtaining a causal graph, represented as an adjacency matrix, where the rows represent the effects and the columns indicate the causes. In this matrix, a causal relationship from variable A to variable B corresponds to an entry 1 in the (B, A) position, implying that the individual will cast a vote in favor of this relationship. The matrices of the three individuals are then added, yielding an adjacency matrix with the total number of votes for each causal path. A voting rule is then applied on each cell of this matrix. Only the directed edges that get approved by the voting rule are kept in the aggregate causal DAG.

Figure 1: Diagram of the voting aggregation



Let us note that some variables have a constant value in some countries of the sample. It is thus not feasible to apply any of the causality detection methods discussed above to find relations among those variables. This is the case of the exchange rate in Ecuador, the United States, and Zimbabwe, as well as the total factor productivity in the United States. This reduces the size of the sample that can be used to study the relations among these variables, lowering the number of votes required to accept a causal effect between them.

In the next section, we present the results obtained by applying this methodology to simulated and real-world data.

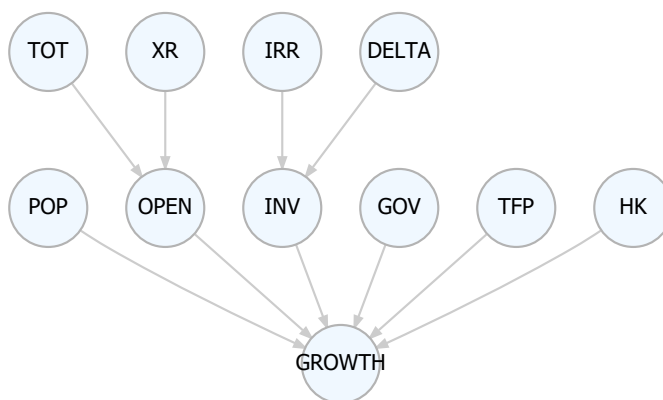
3 Results

3.1 Synthetic data

As a first step in our study, we start by evaluating the quality of the results obtained by each method. We proceed by simulating a growth dataset with the same variables and length as the one obtained from the Penn World Table. This involves creating a causal direct acyclical graph (DAG) in which the same variables are the nodes while the arrows among them represent the causal relations assumed in economic growth theory. We intend this graph to represent the *ground truth* against which we can compare the evidence found by the different methods of causal detection.

Building *the* ground truth model poses a challenge since there exist different models of economic growth, in which the assumptions, functional forms and variables involved differ widely.³ For our exercise, we disregard as much the specific details of the different models and include all the variables in Table 1, establishing pairwise links whenever such relation is favored in the literature, checking for the consistency of the resulting graph, depicted in Figure 2.

Figure 2: Ground truth DAG



We proceed by generating time series based on the ground truth model. We do this for varying lengths ($T = 59, 200, \text{ and } 2000$). For simplicity, we assume that the independent variables are defined by a stationary AR(1) process with a white-noise error term:

$$X_t = \alpha + \beta X_{t-1} + \varepsilon_t ; \quad \beta < 1 , \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2)$$

³To have a taste of this variety of models, see textbooks like Barro and Sala-i-Martin (2004) or surveys like Acemoglu (2012).

The three dependent variables, in turn, are defined by the following linear specifications:

$$\begin{aligned} OPEN_t &= f(TOT_t, XR_t) + \eta_t \\ INV_t &= f(IRR_t, DELTA_t) + v_t \\ GROWTH_t &= f(OPEN_t, INV_t, POP_t, GOV_t, TFP_t, HK_t) + \zeta_t \end{aligned}$$

Each method generates a causal graph, which can be compared to the graph of the ground truth. We also apply two ensembling techniques, combining the predictions of the different methods. On one hand, $ensemble_{\cap}$ incorporates a causal path only if *all* the methods agree on its inclusion, while $ensemble_{\cup}$ includes a causal path if at least one of the methods detects it.

We evaluate the results according to four performance measures: *recall*, *precision*, *specificity*, and the F_1 score. We define TP , TN , FP , and FN as, respectively, the quantities of true positive, true negative, false positive, and false negative edges obtained by a method, compared to the ground truth model. That is, TP is the number of causal links in the ground truth graph that are correctly detected while TN is the number of potential links *absent* in the ground truth that are rightly discarded by the method. In turn, FP is the number of non-causal links wrongly deemed to be causal while FN represents the correct links that are missed by the method. Then:

$$\text{recall} = \frac{TP}{TP+FN}, \quad \text{precision} = \frac{TP}{TP+FP}, \quad \text{specificity} = \frac{TN}{TN+FP}, \quad F_1\text{-score} = \frac{2TP}{2TP+FP+FN}.$$

A method that yields a high recall recovers most of the correct causal links of the ground truth, while higher levels of precision indicate that the method returns more correct results than wrong ones. Specificity measures how well the method can identify links *absent* in the ground truth. The F_1 measure is the harmonic mean of recall and precision and indicates the accuracy of the method.

The results are presented in Table 2, where values are expressed in percentage terms. We can see that Direct-LiNGAM yields, in general, the best results, although no method is perfectly accurate.

Table 2: Performance metrics on ground truth

Method	T	Recall	Precision	Specificity	F1-score
Granger (lags=1)	59	10.0	7.7	88.0	8.7
Transfer Entropy	59	0.0	0.0	94.0	0.0
Stochastic Causality	59	0.0	0.0	95.0	0.0
ICA-LiNGAM	59	20.0	15.4	89.0	17.4
Direct-LiNGAM	59	40.0	28.6	90.0	33.3
$ensemble_{\cap}$	59	0.0	-	100.0	0.0
$ensemble_{\cup}$	59	50.0	13.2	67.0	20.8
Granger (lags=1)	200	20.0	14.3	88.0	16.7
Transfer Entropy	200	30.0	23.1	90.0	26.1
Stochastic Causality	200	10.0	33.3	98.0	15.4
ICA-LiNGAM	200	30.0	25.0	91.0	27.3
Direct-LiNGAM	200	80.0	72.7	97.0	76.2
$ensemble_{\cap}$	200	0.0	-	100.0	0.0
$ensemble_{\cup}$	200	100.0	25.6	71.0	40.8
Granger (lags=1)	2000	10.0	5.6	83.0	7.1
Transfer Entropy	2000	30.0	30.0	93.0	30.0
Stochastic Causality	2000	0.0	0.0	99.0	0.0
ICA-LiNGAM	2000	30.0	23.1	90.0	26.1
Direct-LiNGAM	2000	30.0	18.8	87.0	23.1
$ensemble_{\cap}$	2000	0.0	-	100.0	0.0
$ensemble_{\cup}$	2000	70.0	16.3	64.0	26.4

3.2 Real-world data

When we turn to the real-world data drawn from the Penn World Table dataset, the first result is that no method obtains causal relations unanimously detected in each country. Therefore, all the DAGs presented in this section are the ones that pass either the simple majority or the two-thirds rule. In the DAG corresponding to any method, causal paths are tagged with the number of votes that it gets under that method.

The results obtained according to each method are the following:

- **Granger:** no DAG can be generated by even the majority rule.
- **Transfer Entropy:** two different DAGs are generated under the majority and the two-thirds rules. In both human capital is the source of all the causal links.

Figure 3: Simple majority rule (Transfer Entropy)

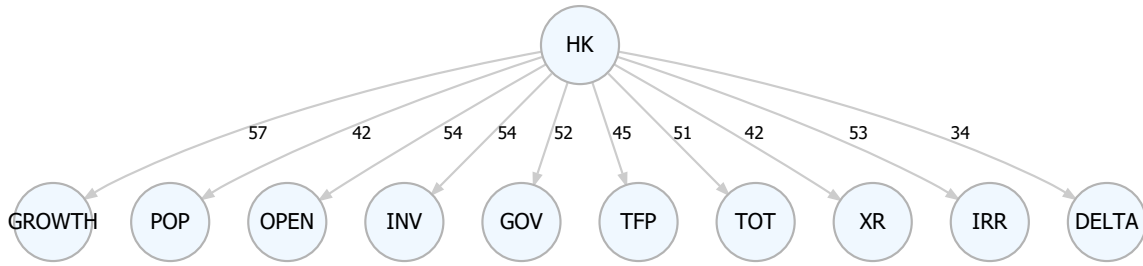
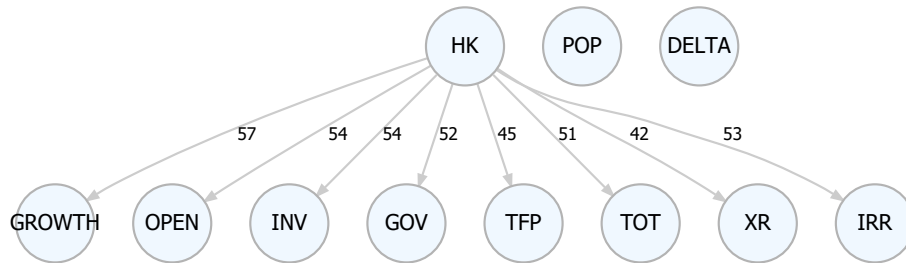
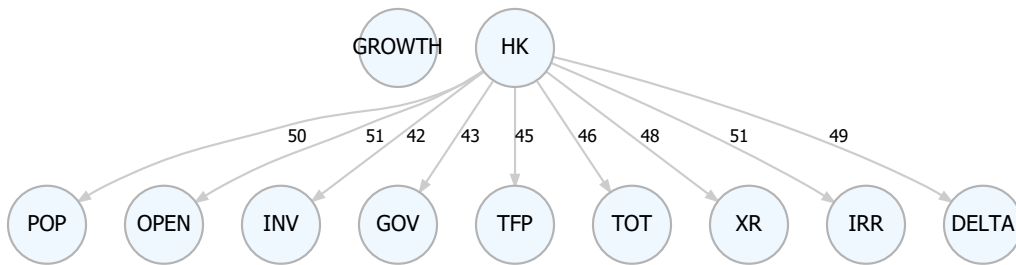


Figure 4: Two-thirds rule (Transfer Entropy)



- **Stochastic Causality:** we find similar results to those obtained by the Transfer Entropy approach, with the human capital as the cause of most of the variables. Unlike with Transfer Entropy, we do not find a causal relationship between human capital and economic growth.

Figure 5: Simple majority and two-thirds rule (Stochastic Causality)



- The **LiNGAM** methods: economic growth is causally independent of the rest of the variables in all the DAGs generated under the two versions of this approach.
 - **ICA-LiNGAM:** the growth rate of the population plays a key role, by initiating all the causal paths in the simple majority DAG, while in the case of the two-thirds DAG, it cedes part of its role. In that case, all the causal paths are initiated either by the growth rate of the population or the rate of depreciation of capital.

Figure 6: Simple majority rule (ICA-LiNGAM)

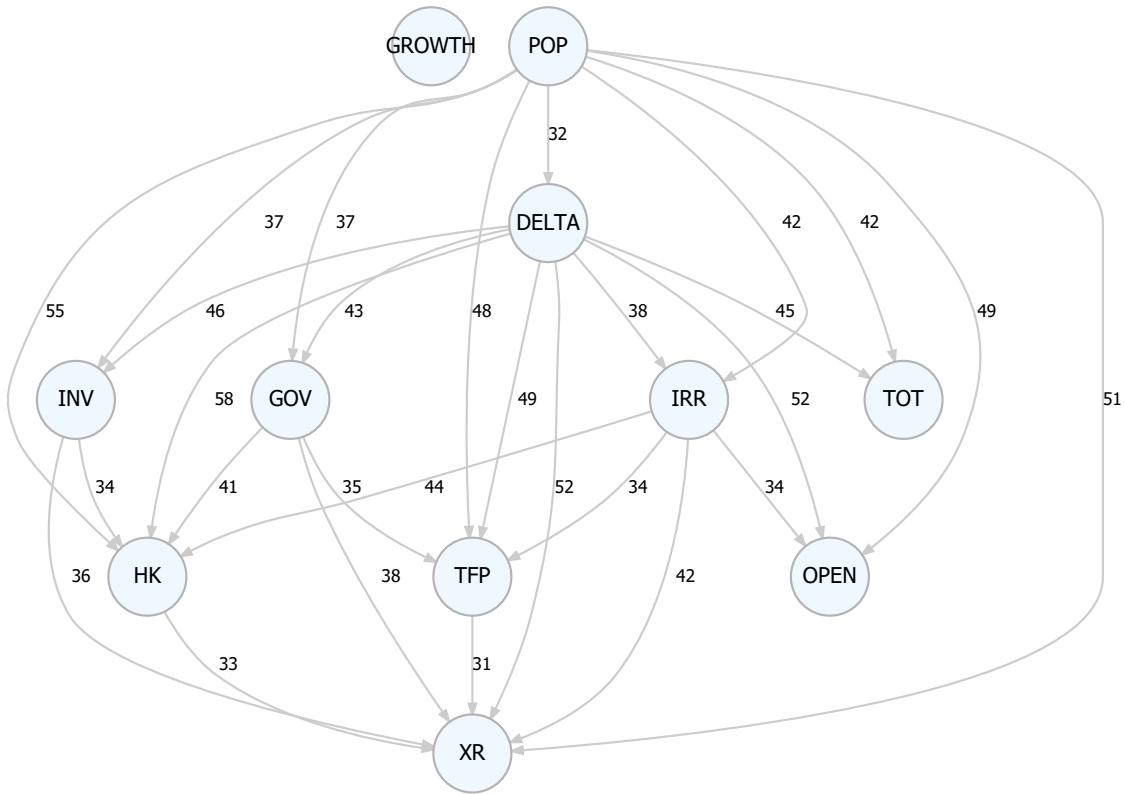
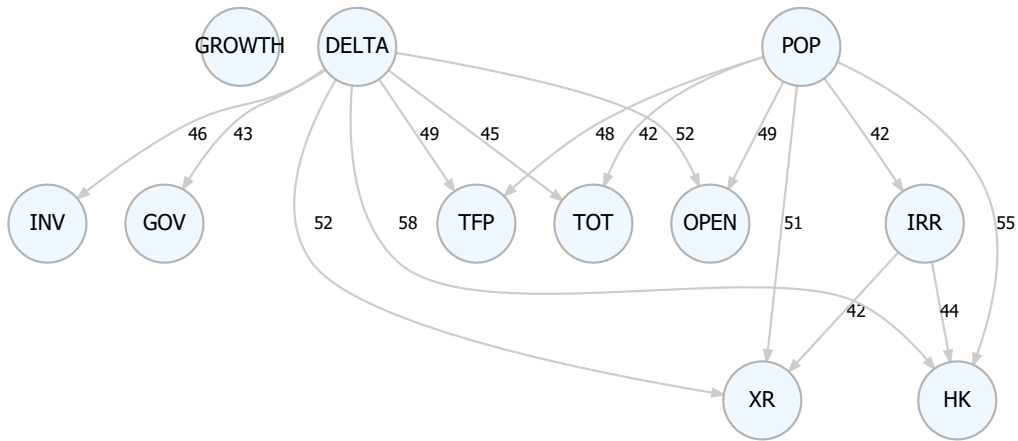
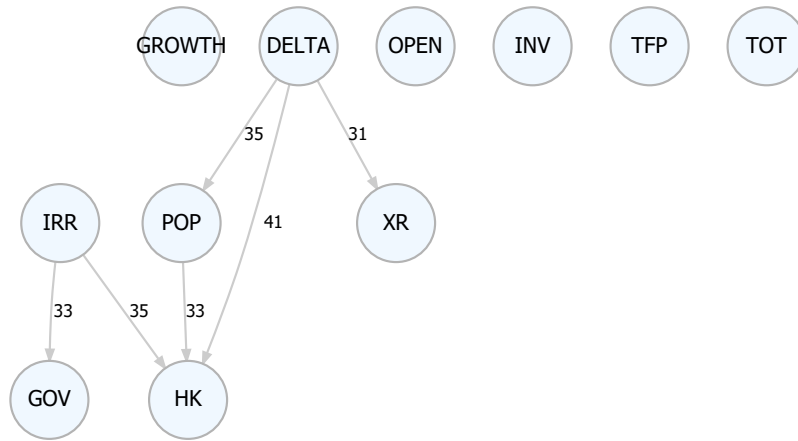


Figure 7: Two-thirds rule (ICA-LiNGAM)



- **Direct-LiNGAM:** the two-thirds rule does not yield a DAG, while many variables are causally independent of the others in the simple majority DAG. The rate of depreciation of capital is one of the sources of causal paths in the graph. The other source is the real interest rate.

Figure 8: Simple majority rule (Direct-LiNGAM)



One aspect worth highlighting in the results is that, except for Stochastic Causality, which lacks the ability to detect it, no method maintained a bidirectional relationship between the variables post-voting. In other words, in this case the methods did not find it for a sufficient majority of countries.

- **ensemble_∪**⁴: the two graphs inherit from *Transfer Entropy* the only causal link involving the variable Growth, namely as caused by Human Capital. The graphs are the following:

⁴*ensemble_∩* is empty both under the majority and the two-thirds rules, even disregarding votes for *Granger* (that yields an empty graph in both cases).

Figure 9: Simple majority rule ($ensemble_{\cup}$)

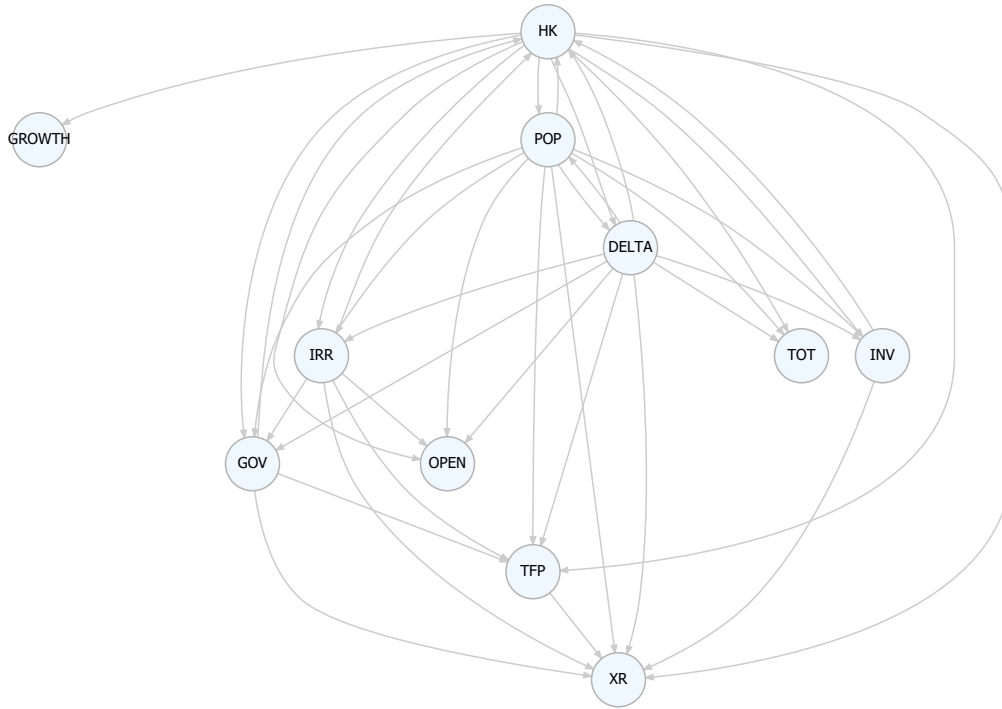
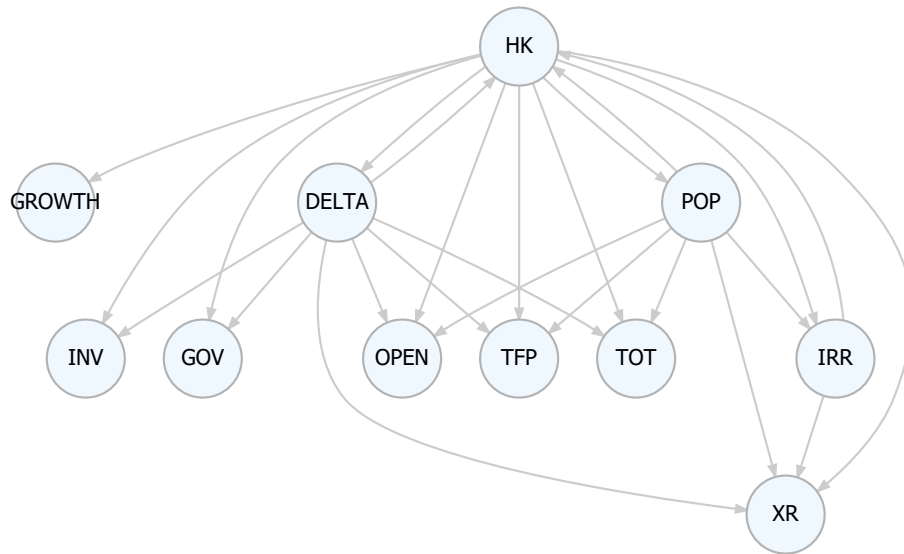


Figure 10: Two-thirds rule ($ensemble_{\cup}$)



The comparison to the ground truth model of subsection 3.1 is summarized in the following table (notice that *Granger* and $ensemble_{\cap}$ are not included since these two methods do not yield graphs):

Table 3: Performance metrics with real-world data (majority rule)

Method	Recall	Precision	Specificity	F1-score
Transfer Entropy	10.0	10.0	91.0	10.0
Stochastic Causality	0.0	0.0	91.0	0.0
ICA-LiNGAM	10.0	3.6	73.0	5.3
Direct-LiNGAM	0.0	0.0	94.0	0.0
<i>ensemble_∪</i>	20.0	5.1	63.0	8.2

Table 4: Performance metrics with real-world data (two-thirds rule)

Method	Recall	Precision	Specificity	F1-score
Transfer Entropy	10.0	12.5	93.0	11.1
Stochastic Causality	0.0	0.0	91.0	0.0
ICA-LiNGAM	10.0	6.7	86.0	8.0
Direct-LiNGAM	0.0	-	100.0	0.0
<i>ensemble_∪</i>	20.0	8.0	77.0	11.4

We can see that scores are consistently equal or higher under the two-thirds rule, although they are (unsurprisingly) worse than with simulated data. We also can see that Transfer Entropy is the method that, under the majority rule, seems to yield better scores than Direct-LiNGAM. In the case of the two-thirds rule, the best results are obtained by *ensemble_∪*.

4 Discussion

In this paper, we have proposed a novel approach based on voting rules for determining causal relationships in panel data settings. This method allows to obtain a causal structure that is supported by several independent sources. Because of the voting structure used, this method could also perform a more comprehensive analysis in terms of country subgroups, in cases where it is relevant. This is because it is possible to disaggregate the votes obtained for a causal relationship into the subgroups identified in the data.

The results reported in Tables 3 and 4 could be seen as indicating that the methods applied in this study may not be too efficient in detecting the causal links present (and absent) in the ground truth. This idea may get reinforced by the results obtained with synthetic data (Table 2).

But while it may be the case that the methods presented here are still too crude to provide causal inferences, an alternative explanation can be put forward. Namely, that the ground truth itself is inaccurate. While the former argument cannot be discarded, the latter shows why the

exploration of causal detection methods matters.

The causal links in the ground truth should not be taken as having indisputable validity. On the contrary, the ground truth model should only be taken provisionally and fully accepted if it gets validated by data. Therefore, the real power of the methods explored here is the possibility of detecting the strongest causal relations in the datasets.

In our analysis, we found that real-world data under the *ensemble_U* method supports, as the only relevant cause of growth, the human capital in the economy. Notice that this result can be understood in terms of the crucial difference between *correlation* and *causation*. So, while the ground truth model is consistent with the theory and evidence of economic growth, most of the links in it may just represent strong correlations, with human capital as the main underlying common cause.

To examine if this is the case, we may need to apply the approach of Pearl (2000) to detect *confounders*. Natural experiments in which some variables are intervened may provide the information needed to implement such exploration. In the meanwhile, the exploration of alternative Machine Learning methods seems worthwhile.

Statements and Declarations

Conflict of interest: Authors declare that they have no conflict of interest.

References

- Acemoglu, D. (2012). Introduction to economic growth. *Journal of economic theory*, 147(2), 545–550.
- Barnett, L., Barrett, A. B., & Seth, A. K. (2009). Granger causality and transfer entropy are equivalent for gaussian variables. *Phys. Rev. Lett.*, 103(23), 238701.
- Barro, R. J., & Sala-i-Martin, X. (2004). *Economic growth*. MIT Press.
- Bruns, S. B., & Stern, D. I. (2019). Lag length selection and p-hacking in granger causality testing: Prevalence and performance of meta-regression models. *Empirical Economics*, 56, 797–830.
- Dumitrescu, E.-I., & Hurlin, C. (2012). Testing for granger non-causality in heterogeneous panels. *Economic Modelling*, 29(4), 1450–1460.

- Feenstra, R. C., Inklaar, R., & Timmer, M. P. (2015). The next generation of the Penn World Table. *American economic review*, 105(10), 3150–3182.
- Glymour, M., Pearl, J., & Jewell, N. P. (2016). *Causal inference in statistics: A primer*. John Wiley & Sons.
- Granger, C. (1969). Investigating causal relations by econometric models and cross-spectral methods. *Econometrica*, 37(3), 424–438.
- Holtz-Eakin, D., Newey, W., & Rosen, H. S. (1988). Estimating vector autoregressions with panel data. *Econometrica*, 56(6), 1371–1395.
- Juodis, A., Karavias, Y., & Sarafidis, V. (2021). A homogeneous approach to testing for granger non-causality in heterogeneous panels. *Empirical Economics*, 60(1), 93–112.
- Papana, A., Siggiridou, E., & Kugiumtzis, D. (2021). Detecting direct causality in multivariate time series: A comparative study. *Communications in Nonlinear Science and Numerical Simulation*, 99, 105797.
- Pearl, J. (2000). *Causality: Models, reasoning and inference*. Cambridge University Press.
- Peters, J., Janzing, D., & Schölkopf, B. (2017). *Elements of causal inference: Foundations and learning algorithms*. The MIT Press.
- Raunig, B. (2023). Using causal graphs to test for the direction of instantaneous causality between economic policy uncertainty and stock market volatility. *Empirical Economics*, 1–20.
- Runge, J., Nowack, P., Kretschmer, M., Flaxman, S., & Sejdinovic, D. (2019). Detecting and quantifying causal associations in large nonlinear time series datasets. *Science Advances*, 5(11), eaau4996.
- Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., & Bengio, Y. (2021). Toward causal representation learning. *Proceedings of the IEEE*, 109(5), 612–634.
- Schreiber, T. (2000). Measuring information transfer. *Phys. Rev. Lett.*, 85(2), 461–464.
- Shimizu, S., Hoyer, P. O., Hyvärinen, A., & Kerminen, A. (2006). A linear non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(72), 2003–2030.

- Shimizu, S., Inazumi, T., Sogawa, Y., Hyvärinen, A., Kawahara, Y., Washio, T., Hoyer, P. O., & Bollen, K. (2011). DirectLiNGAM: A direct method for learning a linear non-gaussian structural equation model. *Journal of Machine Learning Research*, 12(33), 1225–1248.
- Spirtes, P., & Glymour, C. (1991). An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1), 62–72.
- Vinod, H. D. (2017). Generalized correlation and kernel causality with applications in development economics. *Communications in Statistics - Simulation and Computation*, 46(6), 4513–4534.
- Vinod, H. D. (2019). New exogeneity tests and causal paths. In H. D. Vinod & C. Rao (Eds.), *Conceptual econometrics using R* (pp. 33–64). Elsevier.
- Zheng, S., Shi, N.-Z., & Zhang, Z. (2012). Generalized measures of correlation for asymmetry, nonlinearity, and beyond. *Journal of the American Statistical Association*, 107(499), 1239–1252.

A Appendix

Table 5: Countries in the sample

Argentina	United Kingdom	Malaysia
Australia	Greece	Nigeria
Austria	Guatemala	Netherlands
Belgium	Indonesia	Norway
Bolivia	India	New Zealand
Brazil	Ireland	Peru
Canada	Iran	Philippines
Switzerland	Iceland	Portugal
Chile	Israel	Paraguay
China	Italy	Sweden
Colombia	Jamaica	Thailand
Costa Rica	Jordan	Trinidad and Tobago
Cyprus	Japan	Turkey
Germany	Kenya	Taiwan
Denmark	Republic of Korea	Uruguay
Dominican Republic	Sri Lanka	United States
Ecuador	Luxembourg	Venezuela
Egypt	Morocco	South Africa
Spain	Mexico	Zambia
Finland	Malta	Zimbabwe
France	Mauritius	