

Assessing the presence of causal effects between climatic variables and crops extension: the case of Indonesia

Benedetta Francesconi

Vrije Universiteit Amsterdam, NL & University of Luxembourg, LU

b.francesconi@student.eur.nl

October 26, 2023

Abstract

Causal analysis can drastically improve our understanding of climate change impacts over crop yields. This knowledge can guide policy makers in the design of more effective adaptation and mitigation policies. We propose an integration of machine learning and econometric causal methods, which have the potential to provide new insights over the cause-effect dynamics in place and drastically improve the forecasting of future crop yields.

1 Motivation

Accurately assessing the causal effects that climate change has on crops growth and extension is critical for the correct design of adaptation and mitigation policies. According to a growing body of literature, crop yields are expected to be damaged by changes in climate [1, 2, 3]. This represents a serious problem especially for developing countries whose economies still heavily rely on the agricultural sector. In this work, we focus on the case of Indonesia and on the consequences that climate change may have over its crop extension. Indonesia is known for being one of the countries where climate change is expected to exert its harshest effects [4, 5, 6, 7]. Its agricultural industry is still heavily based on smallholder farmers, who occupy 89% of farming land [8] and who may not be able to cope with the effects of future changes in climate [9, 10, 11, 12]. The impacts on the overall economy and society have the potential to be catastrophic.

Causal inference is an intricate problem, in this case further complicated by known variations among plant types, twisted unknown causal structures, and possibly the presence of confounders [13, 14, 15]. There is a growing body of literature studying the effects of climate change over crop yields [16, 17, 18, 19, 20, 21, 22]. However, most of this research does not focus on the assessment of causal effects.

When they do, they mostly employ techniques relying on assumptions of linearity; yet, it is known that such assumptions might be too restrictive for the problem at hand [23, 24, 25]. Following [26, 27], we propose to integrate the Granger causality test and the Kolmogorov-Smirnov test with the novel DeepAR machine learning model.

2 Goals and contributions

First, we apply powerful nonlinear causal analysis techniques to satellite data, in the specific field of climate change effects over crops growth. Per our knowledge, such methods have so far only been employed for research topics different than the one pursued in this work. Second, we check the performance of the DeepAR model and of the modified Granger causality test and Kolmogorov-Smirnov (KS) test in the presence of nonstationary variables. Third, we compare the results of the Granger and KS tests as designed in this proposal, which has not been done before. Fourth, as econometric research done on this topic mostly uses year or seasonal level data [28, 29, 30], we contribute to it by employing satellite monthly level data. Last, we study whether the presence of non-climatic variables may also have causal effects. This will be done through the use of data about the spatial extension of other types of vegetation, such as trees or pasture. In the absence of better indicators at this granular spatial level, they may behave as proxies of broader socio-economic forces, such as economic incentives or technological development.

3 Data

The data used are those published by the Coupled Mode Intercomparison Project (CMIP), promoted by the World Climate Research Programme [31]. We refer to the data elaborated and released under phase 6 [32], using the following variables: the percentage of total land covered respectively by crops, grass, trees and pasture; precipitation in $kg/m^2/s$; surface downwelling and upwelling shortwave radiations in W/m^2 ; near-surface air temperature in K. In Figure 1, as an example, we plot near-surface air temperature degrees distribution through time. More can be found in the Appendix.

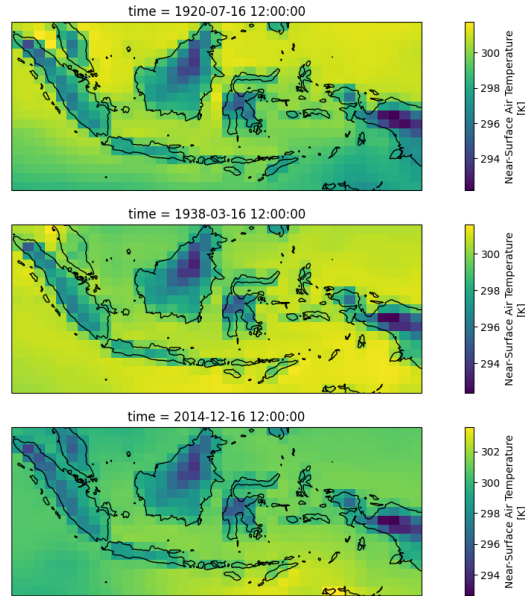


FIGURE 1: Near-surface air temperature concentration through time

The data have a spatial resolution of 1.25° latitude x 2.5° longitude and span from 1850 to 2014. Despite the advantages that such a granular spatial resolution may entail, for the purposes of this proposal we have aggregated the data at the country level.

4 Proposed methodology

The Granger causality test has already been used, as originally envisioned in [33], to investigate the cause-effect relations spanning between climate and crops [34, 35]. In its original form the test is not designed to deal with nonlinear causal effects. For this reason, we apply the DeepAR model in place of the standard Vector Auto Regressive (VAR) model. We further propose to use deep knockoff counterfactuals to robustly check for the presence of causal effects between the independent variables and our target variable. The knockoffs of a variable are values that have the same distribution as the original data, while being independent of the model output. This means that the knockoffs will simulate the behaviour and values of the original variable, while being generated such that they have no causal relation on the target variable. The strategy here is to run the DeepAR model using only the original data and then run again the model intervening each time on one of the independent variables; by intervention we mean that the original values of the variable chosen are substituted with its knockoffs.

The Granger causality test and the KS test will be employed together with the technique of deep knockoff interventions in the following way. The metric used for the Granger causality test is designed according to the Causal Significance Score (CSS)

outlined in [27]: $CSS_{i \rightarrow j} = \ln \frac{MAPE_j^i}{MAPE_j}$. The $MAPE_j$ is the Mean Absolute Percentage Error (MAPE) obtained when running the DeepAR model to forecast the variable j , using only the original data; the $MAPE_j^i$ is the one obtained when running the DeepAR model with intervention on the i independent variable. With respect to the KS test, we also make use of the DeepAR model, run with and without knockoffs intervention, although in a different way. Here we compare the mean of the distributions of the residuals obtained from the model run with original data only and the model run each time with the knockoffs intervention on a separate variable. We then test whether they are statistically significantly different from each other. In case the means of the two distributions are different, it means that the presence of a significant effect over the target variable has been detected. Such effect is spanning from the variables whose values have been substituted with its knockoffs.

In [27], there is the assumption that the data considered is stationary and no test is run for it. When considering climatic variables and crops growth, prior work has often found them to be nonstationary [36, 37, 38, 39]. For this reason, we will start our work performing an Augmented Dickey Fuller (ADF) unit root test. The literature does not always agree about the ability of deep learning models to deal with nonstationary data [40, 41, 42, 43] and the DeepAR model is no exception.

For this reason, we run the DeepAR model using both raw data and the data that has been manipulated and made stationary. Given the big differences in the value scales of the variables, we will also run the DeepAR with standardised data. We will then compare the models performance through means of MAPE scores.

5 Preliminary results

The ADF test results confirmed the presence of a unit root for the percentages of land covered by crops, grass, trees and pasture. They have been differenced taking the difference between the value at time t_1 and the value at time t_{-12} . This is because we believe we should consider a time series as nonstationary the moment in which its statistical properties are not preserved throughout the same period of the year. The spatial extension of crops and of the other types of vegetation considered can highly vary throughout the year naturally, and we can assume that even between two consecutive months there can be drastic differences.

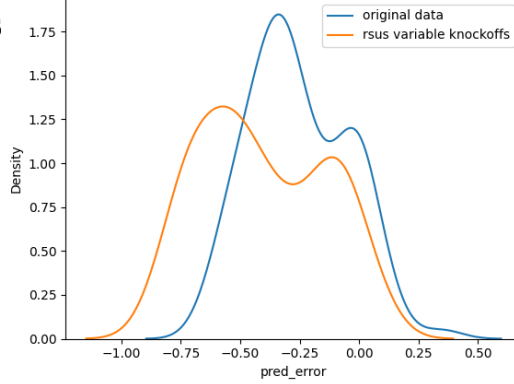


FIGURE 2: Residuals distribution, with and without intervention over surface upwelling shortwave radiation

We ran the DeepAR model with 3 layers, 40 cells, 5 epochs, and employing Long-Short Term Memory (LSTM) model RNNs. We obtain an MAPE of 0.10 when using raw data and an MAPE of 259 when using differenced data. The MAPE is 1.31 when using standardised data. We then performed the Granger and KS tests using raw and standardised data.

For the Granger causality test, when employing raw data we find traces of a causal effect only spanning from the tree fraction of land cover. When using standardised data, we find signs of causal effects also from percentage of pasture cover, surface downwelling and surface upwelling shortwave radiation.

The KS test results do not indicate any presence of causal effect neither in the case of raw data nor with standardised data. As the KS test is based upon the distribution of forecasting errors, we also plotted the forecasting errors obtained under the DeepAR model with and without interventions. Despite the fact that the KS test results do not indicate any causal effect, the forecasting errors of the different model runs are often characterised by different distributions as in Figure 2.

6 Limitations and further work

Further work will be performed to integrate a more granular spatial dimension. As variables involved may interact between each other and lead to new causal factors,

we will also apply the Statistically Enhanced Learning (SEL) techniques proposed in [44], a novel feature engineering methodology able to build new features or co-variates which cannot be directly observed.

References

- [1] Robert Mendelsohn. The impact of climate change on agriculture in developing countries. *Journal of Natural Resources Policy Research*, 1(1):5–19, 2009.
- [2] Oluwole Matthew Akinagbe and Ifeoma Jonathan Irohibe. Agricultural adaptation strategies to climate change impacts in africa: A review. *Bangladesh Journal of Agricultural Research*, 39(3):407–418, 2014.
- [3] Alex C Ruane, David C Major, H Yu Winston, Mozaharul Alam, Sk Ghulam Hussain, Abu Saleh Khan, Ahmadul Hassan, Bhuiya Md Tamim Al Hossain, Richard Goldberg, Radley M Horton, et al. Multi-factor impact analysis of agricultural production in bangladesh with climate change. *Global environmental change*, 23(1):338–350, 2013.
- [4] Andrianto Ansari, Yu-Pin Lin, and Huu-Sheng Lur. Evaluating and adapting climate change impacts on rice production in indonesia: a case study of the keduang subwatershed, central java. *Environments*, 8(11):117, 2021.
- [5] Ria Cahyaningsih, Jade Phillips, Joana Magos Brehm, Hannes Gaisberger, and Nigel Maxted. Climate change impact on medicinal plants in indonesia. *Global Ecology and Conservation*, 30:e01752, 2021.
- [6] Hafidha Asni Akmalia. The impact of climate change on agriculture in indonesia and its strategies: A systematic review. *AGRITEPA: Jurnal Ilmu dan Teknologi Pertanian*, 9(1):145–160, 2022.
- [7] Rizaldi Boer and Yuli Suharnoto. Climate change and its impact on indonesia’s food crop sector. pages 11–13, 2012.
- [8] Silmi Tsurayya¹, Alya Malika, Ardina Latifah Azzahra, Haikal Fadlurrahman¹, and Febriantina Dewi. Determining success criteria for agricultural social start-ups in indonesia. 236:167, 2023.
- [9] Rodel D Lasco, Christine Marie D Habito, Rafaela Jane P Delfino, Florencia B Pulhin, and Rogelio N Concepcion. Climate change adaptation for smallholder farmers in southeast asia. 2011.
- [10] William E Easterling, Pramod K Aggarwal, Punsalma Batima, Keith M Brander, Lin Erda, S Mark Howden, Andrei Kirilenko, John Morton, Jean-François Soussana, Josef Schmidhuber, et al. Food, fibre and forest products. *Climate change*, 2007:273–313, 2007.
- [11] Sixteenth Session. Climate change and the future of smallholder agriculture. 2008.

- [12] D Qin, Z Chen, KB Averyt, HL Miller, S Solomon, M Manning, M Marquis, and M Tignor. Ipcc, 2007: summary for policymakers. 2007.
- [13] A. Bonfante, A. Impagliazzo, N. Fiorentino, G. Langella, M. Mori, and M. Fagnano. Supporting local farming communities and crop production resilience to climate change through giant reed (*arundo donax* L.) cultivation: An italian case study. *Science of The Total Environment*, 601-602:603–613, 2017.
- [14] Samuel Asumadu-Sarkodie and Phebe Asantewaa Owusu. The causal nexus between carbon dioxide emissions and agricultural ecosystem—an econometric approach. *Environmental Science and Pollution Research*, 24(2):1608–1618, October 2016.
- [15] Imran Baig, Farhan Ahmed, Md. Abdus Salam, and Shah Khan. An assessment of climate change and crop productivity in india: A multivariate cointegration framework. *Test Engineering and Management*, 83:3438–52, 08 2020.
- [16] Imran Baig, Farhan Ahmed, Md. Abdus Salam, and Shah Khan. An assessment of climate change and crop productivity in india: A multivariate cointegration framework. *Test Engineering and Management*, 83:3438–52, 08 2020.
- [17] Faiza Ahsan, Abbas Chandio, and Wang Fang. Climate change impacts on cereal crops production in pakistan: Evidence from cointegration analysis. *International Journal of Climate Change Strategies and Management*, ahead-of-print, 02 2020.
- [18] Andrew Crane-Droesch. Machine learning methods for crop yield prediction and climate change impact assessment in agriculture. *Environmental Research Letters*, 13, 11 2018.
- [19] Georgios Giannarakis, Vasileios Sitokonstantinou, Roxanne Suzette Lorilla, and Charalampos Kontoes. Personalizing sustainable agriculture with causal machine learning. 2022.
- [20] X.E. Pantazi, D. Moshou, T. Alexandridis, R.L. Whetton, and A.M. Mouazen. Wheat yield prediction using machine learning and advanced sensing techniques. *Computers and Electronics in Agriculture*, 121:57–65, 2016.
- [21] Evan J. Coopersmith, Barbara S. Minsker, Craig E. Wenzel, and Brian J. Gilmore. Machine learning assessments of soil drying for agricultural planning. *Computers and Electronics in Agriculture*, 104:93–104, 2014.
- [22] S. Veenadhari, Bharat Misra, and CD Singh. Machine learning approach for forecasting crop yield based on climatic parameters. pages 1–5, 2014.
- [23] Wolfram Schlenker and Michael J Roberts. Estimating the impact of climate change on crop yields: The importance of nonlinear temperature effects. 2008.
- [24] Carlo Fezzi and Ian Bateman. The impact of climate change on agriculture: nonlinear effects and aggregation bias in ricardian models of farmland values.

Journal of the Association of Environmental and Resource Economists, 2(1):57–92, 2015.

- [25] Richard W Katz. Assessing the impact of climatic change on food production. *Climatic Change*, 1(1):85–96, 1977.
- [26] Benedetta Francesconi and Ying-Jung C Deweese. Robustly modeling the nonlinear impact of climate change on agriculture by combining econometrics and machine learning. *NeurIPS 2023 Workshop on Tackling Climate Change with Machine Learning*, 2023, 2023.
- [27] Wasim Ahmad, Maha Shadaydeh, and Joachim Denzler. Causal inference in non-linear time-series using deep networks and knockoff counterfactuals. 09 2021.
- [28] Md Abdur Rashid Sarker, Khorshed Alam, and Jeff Gow. Assessing the effects of climate change on rice yields: An econometric investigation using bangladeshi panel data. *Economic Analysis and Policy*, 44(4):405–416, 2014.
- [29] Ebrima K Ceesay and Mohamed Ben Omar Ndiaye. Climate change, food security and economic growth nexus in the gambia: Evidence from an econometrics analysis. *Research in Globalization*, 5:100089, 2022.
- [30] Saeid Satari Yuzbashkandi and Sadegh Khalilian. On projecting climate change impacts on soybean yield in iran: an econometric approach. *Environmental Processes*, 7:73–87, 2020.
- [31] WCRP Coupled Model Intercomparison Project. <https://www.wcrp-climate.org/wgcm-cmip>. Accessed: 2023-08-24.
- [32] Coupled Model Intercomparison Project phase 6. <https://wcrp-cmip.org/cmip-phase-6-cmip6/>. Accessed: 2023-08-24.
- [33] Clive WJ Granger. Some recent development in a concept of causality. *Journal of econometrics*, 39(1-2):199–211, 1988.
- [34] Ebrima K Ceesay, Phillips C Francis, Sama Jawneh, Matarr Njie, Christopher Belford, and Mustapha Momodou Fanneh. Climate change, growth in agriculture value added, food availability and economic growth nexus in the gambia: a granger causality and ardl modeling approach. In *Food Security and Safety Volume 2: African Perspectives*, pages 435–468. Springer, 2022.
- [35] Opeyemi Eyitayo Ayinde, Oluwafemi Olajide Ajewole, Israel Ogunlade, Mathew Olaniyi Adewumi, et al. Empirical analysis of agricultural production and climate change: A case study of nigeria. *Journal of Sustainable Development in Africa*, 12(6):275–283, 2010.
- [36] Gabriele Villarini, James A Smith, and Francesco Napolitano. Nonstationary modeling of a long record of rainfall and temperature over rome. *Advances in Water Resources*, 33(10):1256–1267, 2010.

- [37] Taha BMJ Ouarda and Christian Charron. Nonstationary temperature-duration-frequency curves. *Scientific reports*, 8(1):15493, 2018.
- [38] Yixuan Wang, Jianzhu Li, Ping Feng, and Rong Hu. A time-dependent drought index for non-stationary precipitation series. *Water Resources Management*, 29:5631–5647, 2015.
- [39] Wenbin Wu, Ximing Wu, Yu Yvette Zhang, and David Leatham. Gaussian process modeling of nonstationary crop yield distributions with applications to crop insurance. *Agricultural Finance Review*, 81(5):767–783, 2021.
- [40] Jiachen Zhang, Xingquan Zuo, Mingying Xu, Jing Han, and Baisheng Zhang. Base station network traffic prediction approach based on lma-deepar. pages 473–479, 2021.
- [41] Siqi Liu and Andreas Lehrmann. Dynaconf: Dynamic forecasting of non-stationary time-series. *arXiv preprint arXiv:2209.08411*, 2022.
- [42] Bohdan M Pavlyshenko. Forecasting of non-stationary sales time series using deep learning. *arXiv preprint arXiv:2205.11636*, 2022.
- [43] Sercan O Arik, Nathanael C Yoder, and Tomas Pfister. Self-adaptive forecasting for improved deep learning on non-stationary time-series. *arXiv preprint arXiv:2202.02403*, 2022.
- [44] Florian Felice, Christophe Ley, Andreas Groll, and Stéphane Bordas. Statistically enhanced learning: a feature engineering framework to boost (any) learning algorithms. *arXiv preprint arXiv:2306.17006*, 2023.
- [45] Mohammad Taha Bahadori and Yan Liu. An examination of practical granger causality inference. pages 467–475.
- [46] David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.
- [47] Rina Foygel Barber and Emmanuel J. Candès. Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.

7 Appendix

7.1 Data maps

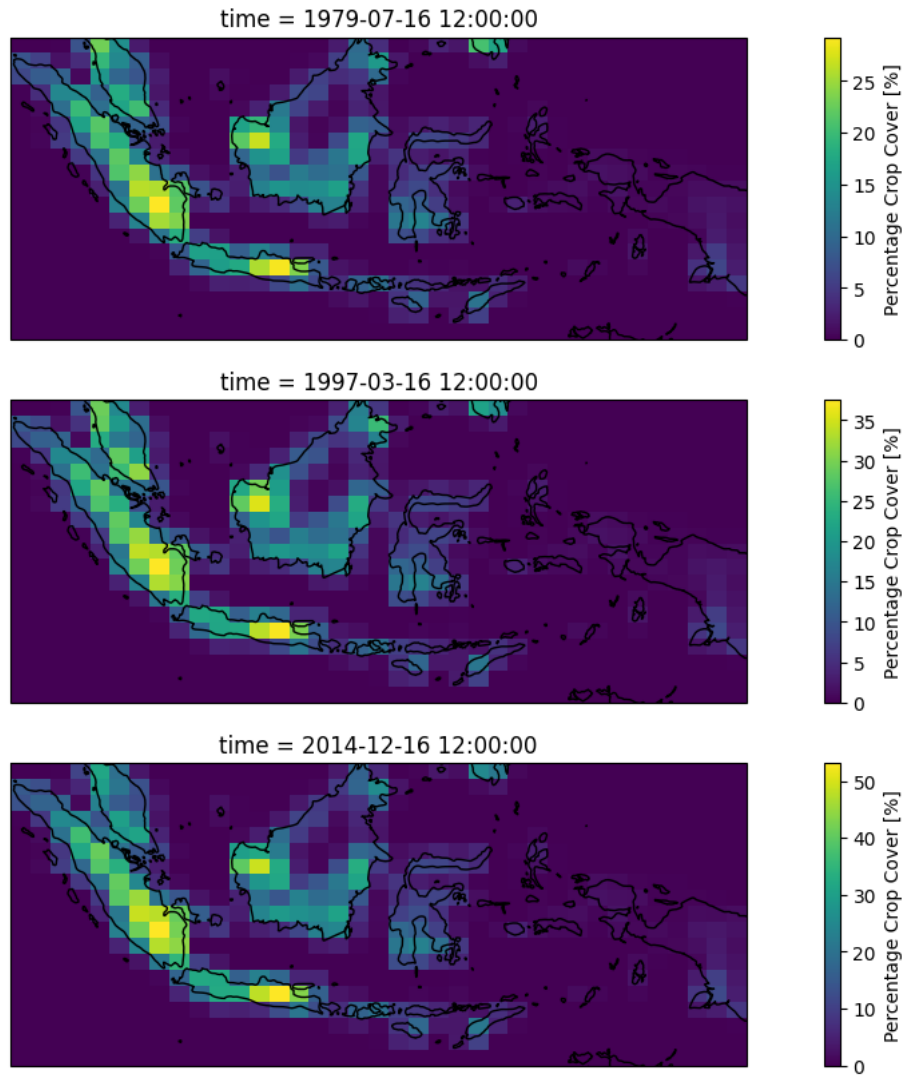


FIGURE 3: Land cover percentage of crops through time

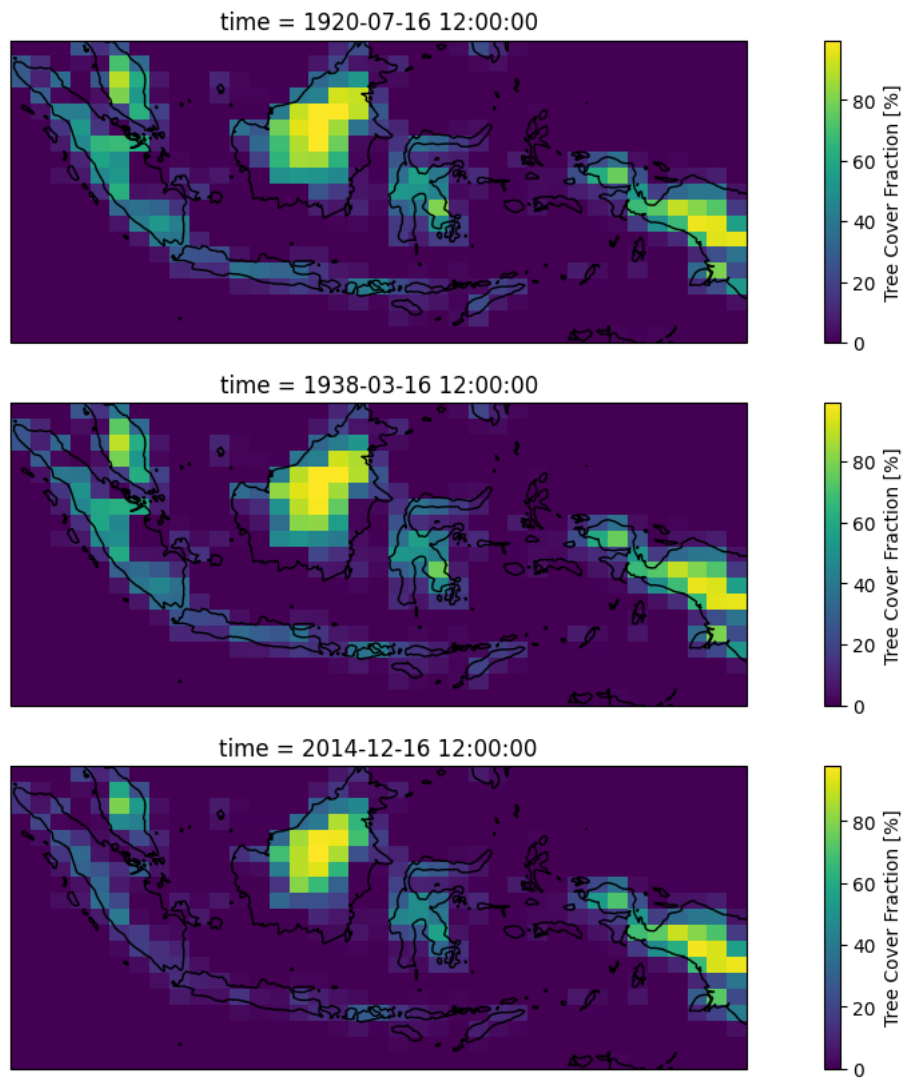


FIGURE 4: Land cover percentage of trees through time

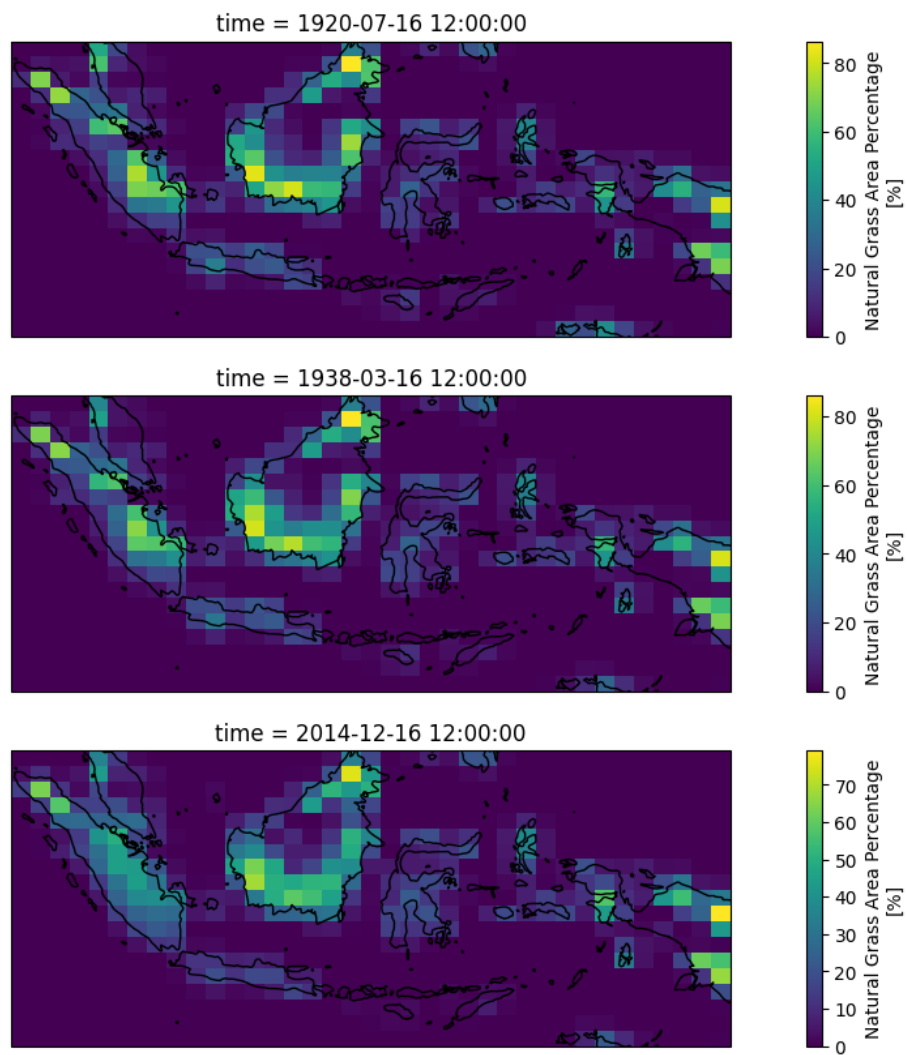


FIGURE 5: Land cover percentage of grass through time

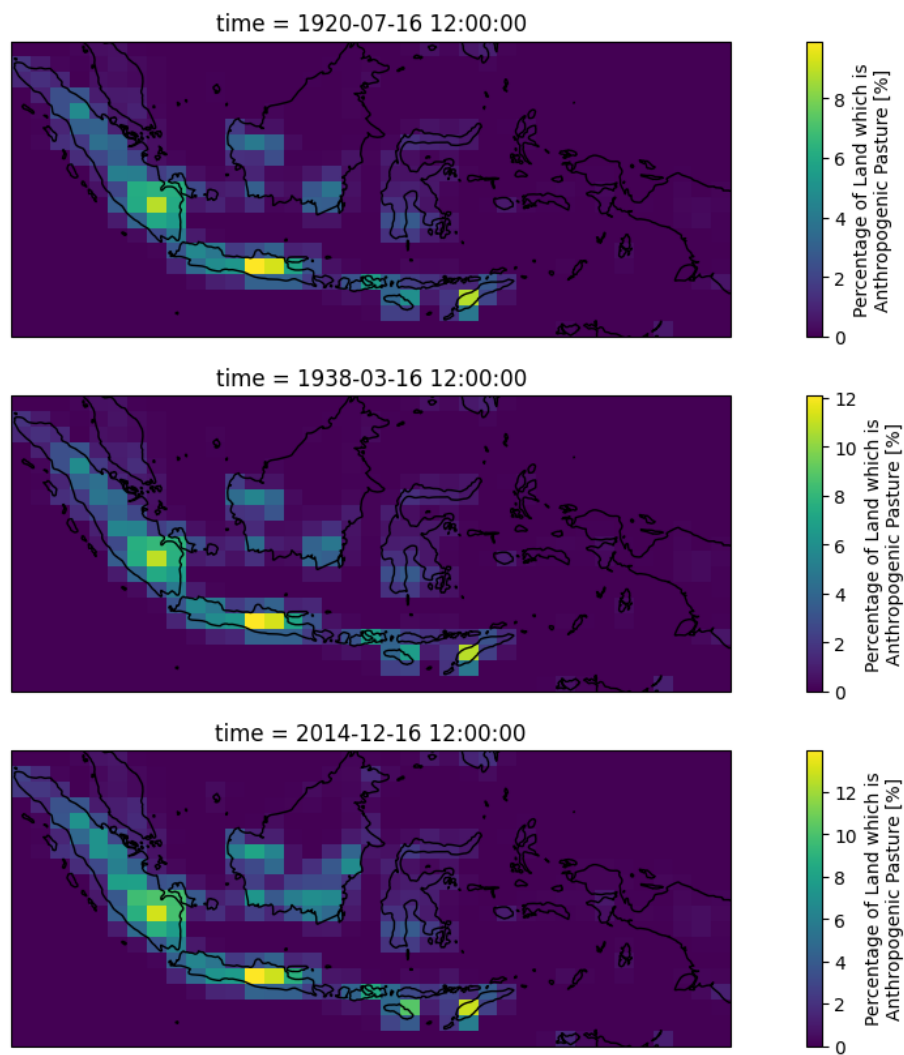


FIGURE 6: Land cover percentage of pasture through time

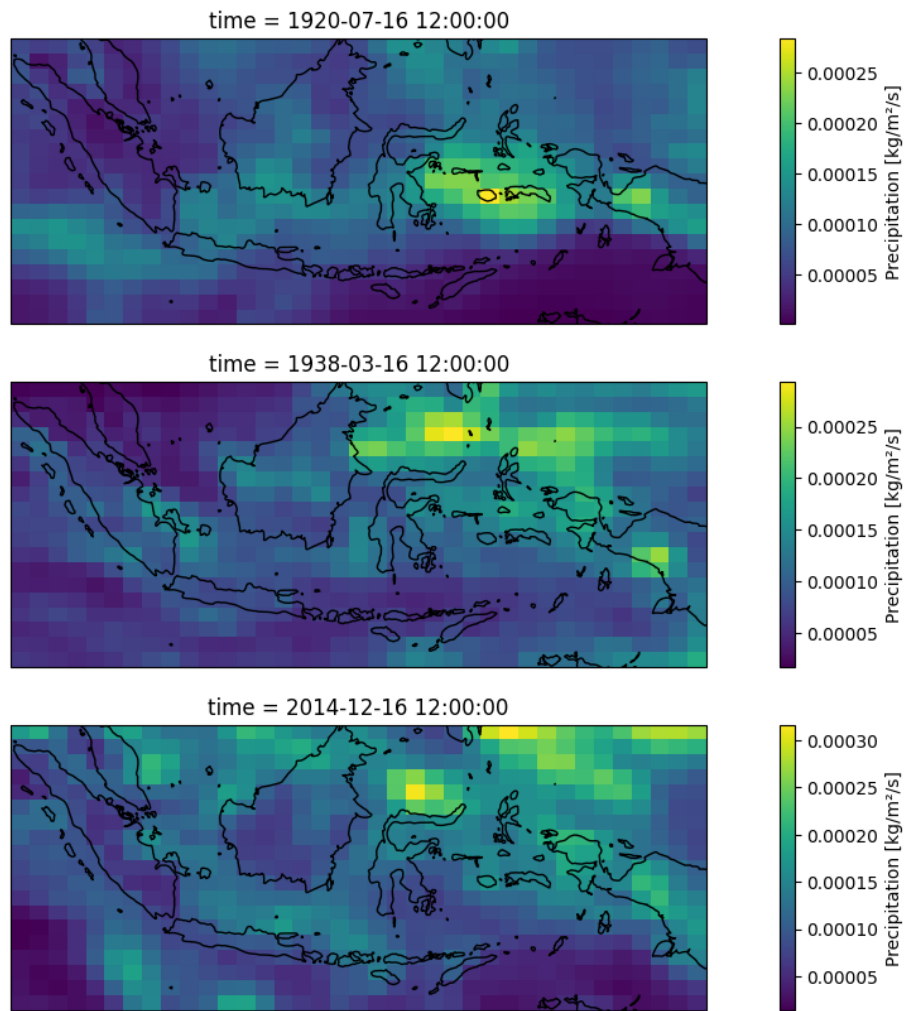


FIGURE 7: Precipitation concentration through time

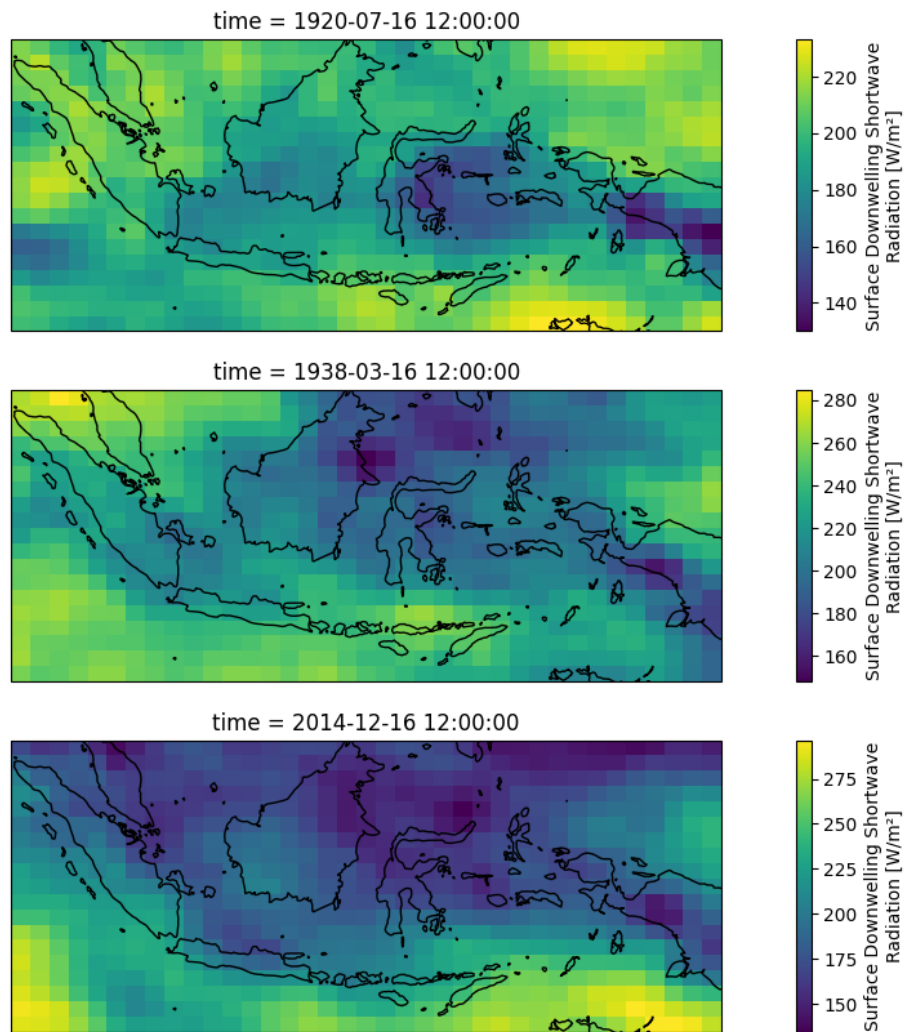


FIGURE 8: Surface downwelling shortwave radiation concentration through time

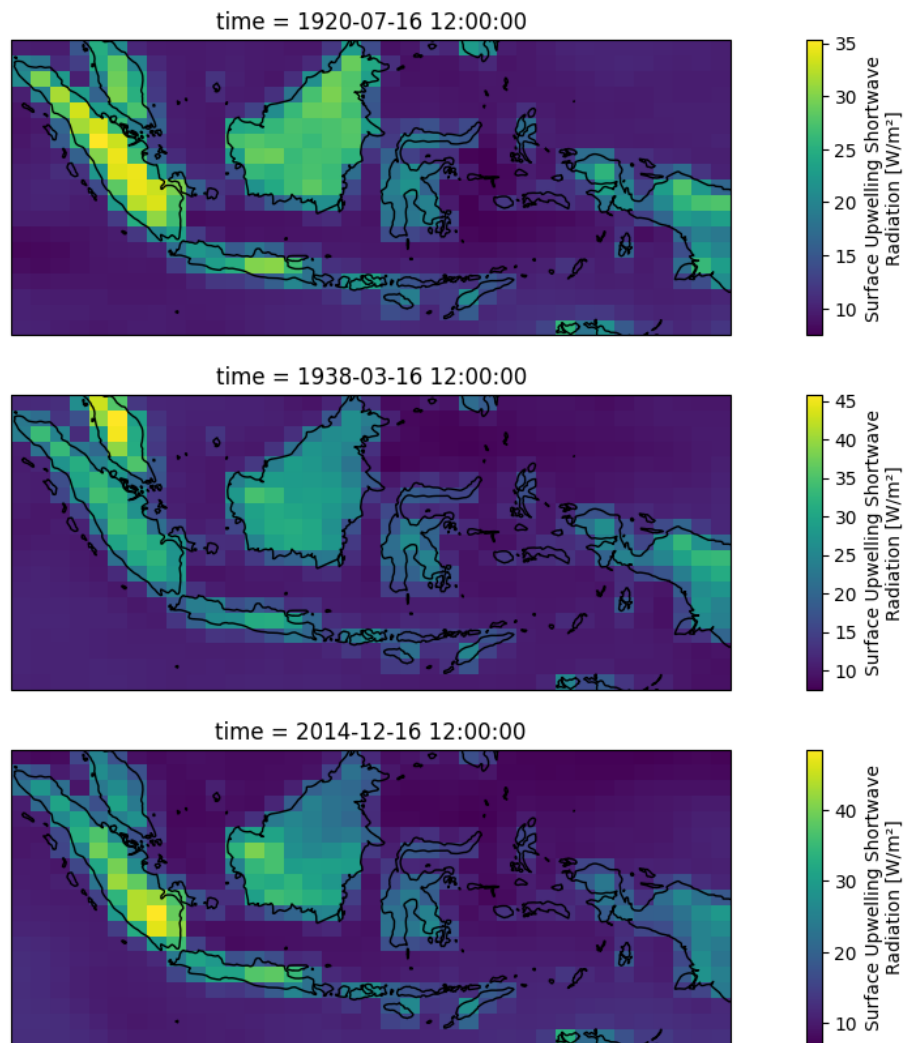


FIGURE 9: Surface upwelling shortwave radiation concentration through time

7.2 Granger and Kolmogorov-Smirnov test results

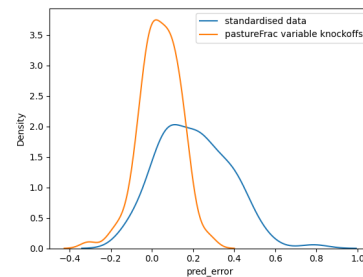
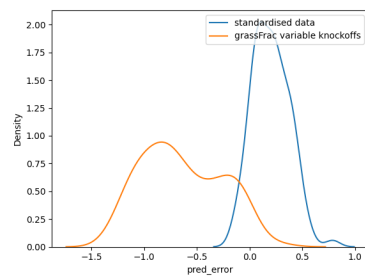
Variable	Granger test	KS test
grassFrac	no causal	no causal
treeFrac	causal effect	no causal
pastureFrac	no causal	no causal
precipitation	no causal	no causal
rsds	no causal	no causal
rsus	no causal	no causal
tas	no causal	no causal

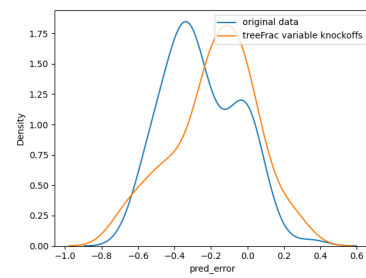
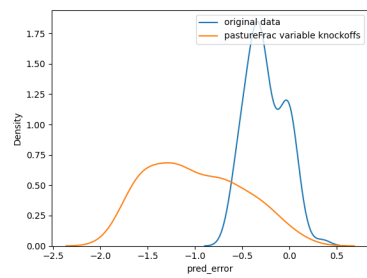
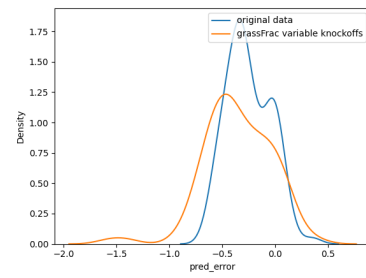
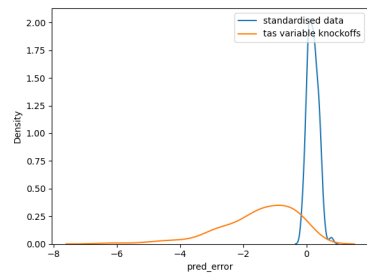
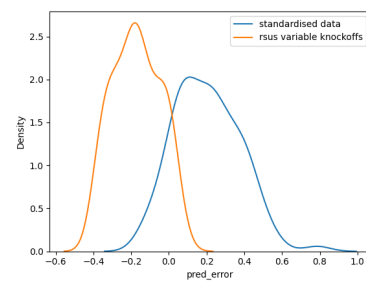
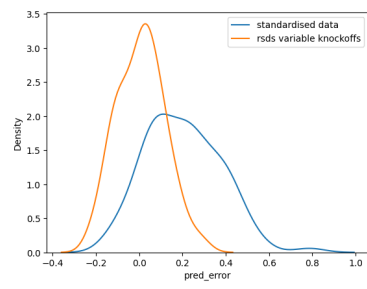
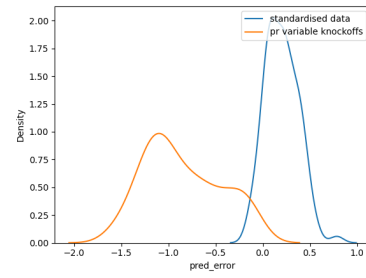
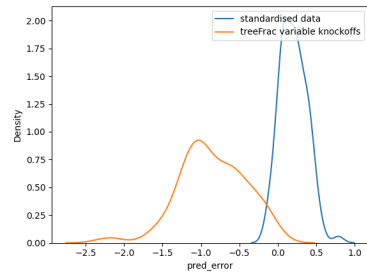
TABLE 1: Raw untransformed data

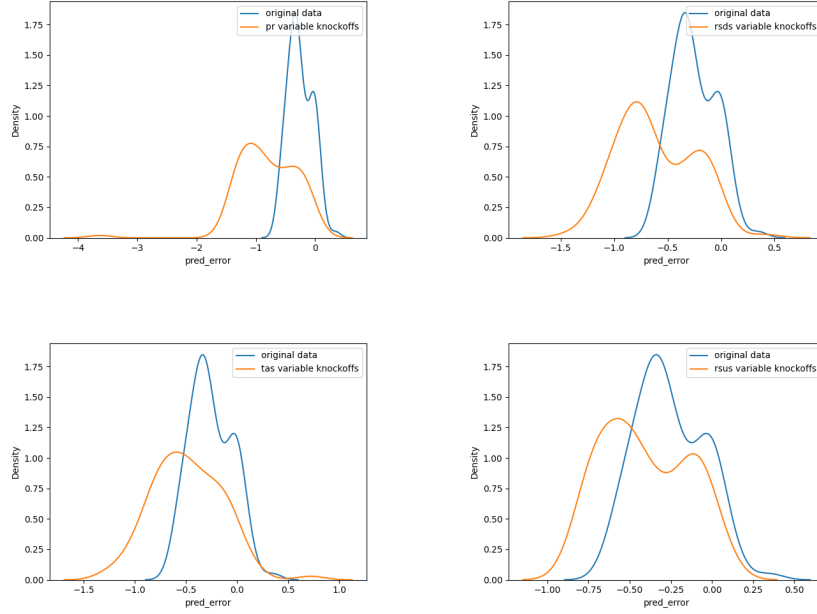
Variable	Granger test	KS test
grassFrac	no causal	no causal
treeFrac	no causal	no causal
pastureFrac	causal	no causal
precipitation	no causal	no causal
rsds	causal	no causal
rsus	causal	no causal
tas	no causal	no causal

TABLE 2: Standardised data

7.3 Forecasting errors distributions







8 Theoretical background

8.1 Augmented Dickey Fuller

The Augmented Dickey Fuller is a test used to check whether a time series is characterised by a unit root, in which case there are grounds to believe the series is non-stationary. The Null hypothesis is that a unit root is present, while the Alternative hypothesis is that no unit root is present.

Given a model of the type:

$$\Delta y_t = \mu + \gamma_t + \alpha y_{t-1} + \sum_{j=1}^{k-1} \beta_j \Delta x_{t-j} + \epsilon_t \quad (1)$$

The Null hypothesis is that $H_0 : \alpha = 0$ against the alternative $H_1 : \alpha < 0$. The test statistic is expressed as $DF_\tau = \frac{\hat{\alpha}}{SE(\hat{\alpha})}$, to compare to the relevant critical values.

8.2 Granger causality

Granger causality is based on two main principles: [i] the cause usually precedes its effects and [ii] the cause makes unique changes in the effects. We say that a variables X Granger causes Y , if the past values of X support in predicting the future values of Y beyond what could have been done with the past values of Y only. The stationarity of the time series under consideration is a fundamental assumption of the Granger causality analysis. If the series involved are not stationary, they have to

be made so through differencing or other techniques, such as taking the log of the series. Specifically, given two stationary time series X and Y we have two different information sets: [i] $I^*(t)$ the set of all available information up to time t and [ii] $I_{-X}^*(t)$ the set of all available information up to time t excluding the information provided by X . If X really aids the prediction of future values of Y , the conditional distribution of the future values of Y should differ under the information set $I^*(t)$ and under $I_{-X}^*(t)$ [45]. Then, X is defined to Granger cause Y if

$$\mathbb{P}[Y(t+1) \in A | I^*(t)] \neq \mathbb{P}[Y(t+1) \in A | I_{-X}^*(t)] \quad (2)$$

for some measurable set $A \subseteq \mathbb{R}$ and $t \in \mathbb{Z}$, with A being the set of future realisations of $Y(t)$.

However, modeling the distribution of multivariate time series can be highly complicated, especially when using functions with non-convex loss landscapes such as deep neural networks. Moreover, the Granger causality definition does not give exact assumptions on the data generating process of the variables involved. For this reason, a usual approach to test for the presence of Granger causality is through the estimation of linear models, which tend to be easy to estimate and yet robust in their estimation. The VAR model is one of such models and one of the most used. The idea is the following. Given several time series X_1, \dots, X_V , we estimate the following VAR for each of the X_j time series:

$$X_j(t) = \sum_{i=1}^V \beta_{j,i}^T \mathbf{X}_i^{t,Lagged} + \epsilon_j(t) \quad (3)$$

where $\mathbf{X}_j^{t,Lagged} = [X_i(t-L), \dots, X_i(t-1)]$ is the history of X_i up to time t , L is the maximal time lag and $\beta_{j,i} = [\beta_{j,i}(1), \dots, \beta_{j,i}(L)]$ is the vector of coefficients modeling the effects of X_i on the target time series. The Granger causality is tested estimating the model with and without all the possible X_i values with $i = 1, \dots, V$. If the conditional probability of the target variable X_j does not change under the different models, then there is no Granger causality as expressed in 3.

8.3 Kolmogorov-Smirnov

The Kolmogorov-Smirnov test for causal inference make use of the distribution of the residuals. Given the forecasting residuals $R \sim e_1, e_2, \dots, e_n$ and $\tilde{R} \sim \tilde{e}_1, \tilde{e}_2, \dots, \tilde{e}_n$, the test statistics uses the supremum distance between the two:

$$D_{n,m} = \sqrt{\frac{nm}{n+m}} \sup |R_n - \tilde{R}_m| \quad (4)$$

The Null hypothesis is that a variable j does not cause a variable i if the distribution of the residuals is approximately identical across environments.

8.4 DeepAR

Given a target time series $\mathbf{z}_{i,1:t_0-1} = [z_{i,1}, \dots, z_{i,t_0-2}, z_{i,t_0-1}]$ and wanting to estimate its future values $\mathbf{z}_{i,t_0:T} = [z_{i,t_0}, z_{i,t_0+1}, \dots, z_{i,T}]$, we need to model the conditional distribution $P(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T})$, where $\mathbf{x}_{i,1:T}$ is a time series of covariates assumed to be known at all time points. Assuming that the model distribution consists of a product of likelihood factors like

$$Q_\theta(\mathbf{z}_{i,t_0:T} | \mathbf{z}_{i,1:t_0-1}, \mathbf{x}_{i,1:T}) = \prod_{t=t_0}^T Q_\theta(z_{i,t} | \mathbf{z}_{i,1:t-1}, \mathbf{x}_{i,1:T}) = \prod_{t=t_0}^T p(z_{i,t} | \theta(\mathbf{h}_{i,t}, \Theta)) \quad (5)$$

where $\mathbf{h}_{i,t}$ is the output of an autoregressive recurrent network $\mathbf{h}_{i,t} = h(\mathbf{h}_{i,t-1}, z_{i,t-1}, \mathbf{x}_{i,t}, \Theta)$. The h is a function implemented by a multi layer Recurrent Neural Network (RNN) estimated with a Long Short Term Memory (LSTM) model and parametrized by Θ [46]. This model can be used in place of the VAR and then proceeding in the testing of Granger causality. Specifically, it's possible to assess differences in the DeepAR model estimation before and after using specific variables, assumed to Granger cause the target variables, by using the following causal significance score (CSS):

$$CSS_{i \rightarrow j} \ln \frac{MAPE_j^i}{MAPE_j} \quad (6)$$

where $MAPE_j^i$ is the mean absolute percentage error between the $\hat{z}_{j,t}$ and the real $z_{j,t}$ using the variable $z_{i,t}$ and $MAPE_j$ without using $z_{i,t}$.

8.5 Knockoff counterfactual

The knockoff counterfactual technique was first proposed in 2015 [47]. The idea of the technique is to swap the original variables with some fake ones and checking if the model estimations change. Given the set of the original variables Z such that $Z = Z_1, Z_2, \dots, Z_n$, with distribution P_z , the knockoffs are created such that they are in-distribution null variables. The knockoffs have the same distribution as the original variables but they do not contain any information about the target variable, and for this reason they can be swapped with the original variables to check how the model estimation change. Moreover, the knockoffs have the same covariance structure and the correlation between the knockoffs is the same as the correlation between the original variables.