



## Effects of spike sorting error on the Granger causality index



Pei-Chiang Shao<sup>a</sup>, Wan-Ting Tseng<sup>b</sup>, Chung-Chih Kuo<sup>c</sup>, Wei-Chang Shann<sup>a</sup>,  
Meng-Li Tsai<sup>d</sup>, Chien-Chang Yen<sup>e,\*</sup>

<sup>a</sup> Department of Mathematics, National Central University, Jhongli 32001, Taiwan

<sup>b</sup> Institute of Zoology, National Taiwan University, Taipei 10617, Taiwan

<sup>c</sup> Department of Physiology, School of Medicine, College of Medicine, Tzu Chi University, Hualien 97004, Taiwan

<sup>d</sup> Department of Biomechatronic Engineering, National Ilan University, Ilan 26047, Taiwan

<sup>e</sup> Department of Mathematics, Fu-Jen Catholic University, Xinzhuang 24205, Taiwan

### ARTICLE INFO

#### Article history:

Received 8 January 2013

Received in revised form 21 May 2013

Accepted 4 June 2013

#### Keywords:

Granger causality index

Spike sorting

Vector autoregressive model

### ABSTRACT

Accurately sorting individual neurons is a technical challenge and plays an important role in identifying information flow among neurons. Spike sorting errors are almost unavoidable and can roughly be divided into two types: false positives (FPs) and false negatives (FNs). This study investigates how FPs and FNs affect results of the Granger causality (GC) analysis, a powerful method for detecting causal interactions between time series signals. We derived an explicit formula based on a first order vector autoregressive model to analytically study the effects of FPs and FNs. The proposed formula was able to reveal the intrinsic properties of the GC, and was verified by simulation studies. The effects of FPs and FNs were further evaluated using real experimental data from the ventroposterior medial nucleus of the thalamus. Some practical suggestions for spike sorting are also provided in this paper.

© 2013 Elsevier Ltd. All rights reserved.

### 1. Introduction

In neuroscience research, it is important to identify information flow among multiple neurons in the brain, according to the recorded neural activity data. A powerful method for achieving this is the Granger causality (GC) which arose in economics after being introduced by Wiener and Granger (Granger, 1969, 1980; Wiener, 1956). The GC is a time series inference (TSI) type of method, proposes that if the prediction of one time series can be improved with the knowledge of a second time series, then there is a causal influence from the second time series to the first. This prediction is made by using the vector autoregressive (VAR) model. In this model, if the variance of the prediction error of one time series at the present time can be reduced by including the past values of another series, then the latter is said to Granger-cause the former. This causality can be quantified by the so-called GC Index (GCI) which can be used to determine whether there is any causal interaction between time series. The GC was shown to be effective and has been widely deployed in recent neuroscience research (Bressler, Richter, Chen, & Ding, 2007; Cadotte, DeMarse, He, & Ding, 2008; Cadotte et al., 2010; Cao, Maran, Dhamala, Jaeger, & Heck, 2012; Zhang et al., 2012). In addition to the time domain GC, other versions of the GC (e.g., frequency, and time–frequency

domain) have been developed as well (Baccala & Sameshima, 2001; Dhamala, Rangarajan, & Ding, 2008). The time domain formulation of GCI is briefly introduced in the next section, and we refer the reader to an article by Bressler and Seth (2011) for more details about the GC.

Neurons emit action potentials (APs) that are known as *spikes* and play an important role in communicating among cells. The temporal sequence of APs produced by a neuron, which shows its own activity, is also known as a *spike train*. In multi-channel recordings (Brown, Kass, & Mitra, 2004), the APs of neurons are detected and differentiated from background electrical noise before single-unit spike trains are used to probe neural behaviors. This technical procedure is called *spike sorting*. However, it is not easy to obtain spike train data that fully agree with the AP because of noise, superimposed APs, and difficulties of differentiating waveforms of APs from different neurons. Spike sorting often introduces unavoidable errors (Deborah, Won, & Patrick, 2003; Lewicki, 1998). These errors can roughly be divided into two types, false positives (FPs) and false negatives (FNs). An FP means an error detection of an event that is not a real spike (just an electrical noise) or is a spike from another neuron. Conversely, an FN means that real spikes were not detected or were classified into groups of other neurons. One may be interested in the question: “How do FPs and FNs affect the estimation of functional connectivity among neurons?” This study answered this question analytically and also via numerical simulations. The change in the GCI due to spike sorting errors was derived analytically to form an explicit formula,

\* Corresponding author. Tel.: +886 2 2905 3547.

E-mail address: [yen@math.fju.edu.tw](mailto:yen@math.fju.edu.tw) (C.-C. Yen).

and a direct discussion of the effects of FPs and FNs is possible. Moreover, numerical simulations were used to verify the analyses. We constructed three types of models for sorting errors: those with uniform, random, and concentrative distributions. That is, errors occur uniformly, randomly, and concentratively in spike trains. Changes in the GCI were computed as these types of spike sorting errors were artificially added to the simulated spike trains, and the effects on the directional interactions were also investigated.

Finally, it is worth noting that spike trains are non-equally spaced data and are regarded as a point process. Interpolation or filtering is usually employed to convert point processes to equally spaced time series. Previous studies on spike trains (Kaminski, Ding, Truccolo, & Bressler, 2001; Zhu, Lai, Hoppensteadt, & He, 2003) proposed several methods to convert a time series from being non-equally spaced to equally spaced. This study adopted the procedure of binning to convert spike trains into time series data, which are suitable for GC analyses. Although the GCI between two point processes has been directly defined in Kim, Putrino, Ghosh, and Brown (2011) recently, we still cannot abandon binning because it reduces the complexity of analysis, and considers also the effect of temporal summation of action potentials in the neuroscience.

This article is organized as follows. Section 2 presents an analytic formula based on a first order autoregression to show how error processes affect the GCI. Section 3 presents some models for sorting errors and probes the proposed formula further via numerical simulations. Section 4 presents a real data evaluation where the effects of sorting error on the GCI are evaluated using real experimental data. Section 5 provides some suggestions for spike sorting and the discussion.

## 2. Modeling and analysis

Based on a first order autoregression, we derived an explicit formula for changes in the GCI in terms of four parameters involving the error process. We also investigated the influences of various types of errors on the GCI indicated by the proposed formula.

### 2.1. A short introduction to the GCI

Let  $x$  and  $y$  be two stationary time series with zero means. The first order linear autoregressive model for  $x$  and  $y$  is given by

$$\begin{bmatrix} x(n) \\ y(n) \end{bmatrix} = \mathbf{A} \begin{bmatrix} x(n-1) \\ y(n-1) \end{bmatrix} + \begin{bmatrix} \epsilon(n) \\ \eta(n) \end{bmatrix}, \quad (1)$$

where  $\mathbf{A}$  is the model coefficient matrix, and the residuals  $\epsilon$  and  $\eta$  are zero-mean uncorrelated white noises with covariance matrix  $\Sigma$ . Here the variances  $\text{Var}(\epsilon)$  and  $\text{Var}(\eta)$  are called prediction errors, which measure the accuracy of the autoregressive prediction. More specifically,  $\text{Var}(\eta)$  measures the accuracy of the prediction of  $y(n)$  based on the previous values  $x(n-1)$  and  $y(n-1)$ .

Now consider the reduced model that excludes the time series variable  $x$

$$y(n) = By(n-1) + \zeta(n), \quad (2)$$

where  $B$  is the corresponding model coefficient. The variance  $\text{Var}(\zeta)$  measures the accuracy of the prediction of  $y(n)$  based only on its previous value  $y(n-1)$ . For  $\eta$  in (1) and  $\zeta$  in (2), if  $\text{Var}(\eta)$  is significantly less than  $\text{Var}(\zeta)$  in some statistical sense, then we say that  $x$  Granger-cause  $y$ . This causality can be quantified by the GCI from  $x$  to  $y$  formulated as:

$$F_{x \rightarrow y} = \ln \frac{\text{Var}(\zeta)}{\text{Var}(\eta)}. \quad (3)$$

It is clear that  $F_{x \rightarrow y} = 0$  when  $\text{Var}(\eta) = \text{Var}(\zeta)$ , i.e.,  $x$  has no causal influence on  $y$ , and  $F_{x \rightarrow y} > 0$  when  $x$  Granger-cause  $y$ . Notice that  $F_{x \rightarrow y}$  is nonnegative, i.e.,  $\text{Var}(\eta)$  is bounded above by  $\text{Var}(\zeta)$ , since the full model defined in (1) should have a better prediction ability than the reduced model defined in (2). Finally, we note that the GCI values should be checked for significance by using hypothesis testing, and more details of the GCI can be found in Ding, Chen, and Bressler (2006), Granger (1969, 1980).

### 2.2. An explicit formula

When inaccurate spike sorting occurs, the sorting errors can be regarded as a perturbed error process. For simplicity, we assume that only the source process  $x$  has a sorting error and the corresponding error process is denoted by  $\delta x$ . We can assume that  $\delta x$  is zero mean and the model in (1) is perturbed as follows when  $\delta x$  is superposed on  $x$ :

$$\begin{bmatrix} x + \delta x(n) \\ y(n) \end{bmatrix} = \tilde{\mathbf{A}} \begin{bmatrix} x + \delta x(n-1) \\ y(n-1) \end{bmatrix} + \begin{bmatrix} \tilde{\epsilon}(n) \\ \tilde{\eta}(n) \end{bmatrix}, \quad (4)$$

where  $\tilde{\mathbf{A}}$  is the corresponding model coefficient matrix, and the residuals  $\tilde{\epsilon}$  and  $\tilde{\eta}$  have the covariance matrix  $\tilde{\Sigma}$ . Let  $S_y := \text{Var}(\zeta)$ ,  $S := \text{Var}(\eta)$ , and  $\tilde{S} := \text{Var}(\tilde{\eta})$ . Since the perturbed quantity  $\delta x$  is superposed only on  $x$ , the reduced models for (1) and (4) are the same as (2). Then the original GCI from  $x$  to  $y$  and the perturbed GCI from  $x + \delta x$  to  $y$  are

$$F = \ln \frac{S_y}{S} \quad \text{and} \quad \tilde{F} = \ln \frac{S_y}{\tilde{S}}, \quad (5)$$

respectively. To investigate the perturbed GCI, we derived an explicit formula for  $\tilde{F}$  in terms of four parameters involving  $\delta x$  which are  $\xi_1 := E(\delta x_1^2)$ ,  $\xi_2 := E(x_1 \delta x_1)$ ,  $\xi_3 := E(y_2 \delta x_1)$ , and  $\xi_4 := E(y_1 \delta x_1)$ . Further denote  $X_0 = E(x_1^2)$ ,  $Y_0 = E(y_1^2)$ ,  $Y_1 = E(y_1 y_2)$ ,  $Z_1 = E(x_1 y_1)$ , and  $Z_2 = E(x_1 y_2)$ . We are now ready to present the formula for  $\tilde{F}$ .

**Proposition 1.** In the situation described above,  $\tilde{F}$  can be presented explicitly by the following formula (for calculation see the Appendix):

$$\tilde{F} = \ln \frac{S_y}{S + \Theta}, \quad \Theta = (S_y - S)I, \quad (6)$$

where

$$I = \frac{1}{Y_0(X_0 + \xi_1 + 2\xi_2) - (\xi_4 + Z_1)^2} \left\{ Y_0(X_0 + \xi_1 + 2\xi_2) - \frac{1}{S_y - S} [Y_0(\xi_3 + Z_2)^2 + (Y_0 - S)(\xi_4 + Z_1)^2 - 2Y_1(\xi_3 + Z_2)(\xi_4 + Z_1)] \right\}. \quad (7)$$

Note that since  $S + \Theta$  in (6) is bounded above by  $S_y$ , we have that  $\Theta$  is upper bounded by  $S_y - S$ , i.e.,  $I$  has an upper bound 1.

We end this subsection by the following two remarks.

**Remark 1.** In the same situation of Proposition 1, the following inequalities hold:

$$Y_0 \geq S_y \geq S \quad \text{and} \quad Y_1 \leq 0. \quad (8)$$

According to (2), we have  $Y_0 = \text{Var}(y_1) \geq \text{Var}(\zeta) = S_y$ . The remainder  $S_y \geq S$  just follows by the reason that the prediction error of the reduced model in (4) is always less than or equal to that of the full model in (1). The latter holds because of the stationary assumption. If  $Y_1 = E(y_1 y_2) > 0$ , then  $y$  will not be stationary. Thus  $Y_1 \leq 0$ .

Download English Version:

<https://daneshyari.com/en/article/6863472>

Download Persian Version:

<https://daneshyari.com/article/6863472>

[Daneshyari.com](https://daneshyari.com)