



Localization and regularization of normalized transfer entropy



Heeyoul Choi

Samsung Advanced Institute of Technology, Samsung Electronics 97, Samsung2-ro, Giheung-gu, Yongin-si, Gyeonggi-do 446-712, South Korea

ARTICLE INFO

Article history:

Received 13 February 2013

Received in revised form

29 December 2013

Accepted 22 February 2014

Communicated by S. Choi

Available online 13 April 2014

Keywords:

Information theory

Local transfer entropy

Dirichlet distribution

Regularization

ABSTRACT

To find hidden structures of a data set, it is important to understand the relationship between variables such as genes or neurons. As a measure of such relationship, causality is to find directed relations between the variables, which can reveal more of the structures than undirected relations. As a quantitative measure of such causal relationship, transfer entropy has been proposed and successfully applied to capture the amount of information flow between events and sequences. In order to analyze the flow locally in time, we propose to localize normalized transfer entropy and regularize it to avoid the unstable result. Experiment results with synthetic and real-world data confirm the usefulness of our algorithm.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

In data analysis, it is important to understand the relationship among variables, events and sequences. Once we measure relationships between the variables, we can represent the data set on a metric space [1–4] or build up a model for inference on them [5]. Such relationship can be obtained by symmetric measurements like correlation and distance, or asymmetric ones like transition probability and causality.

Causality is a kind of asymmetric (or directional) relationship between two variables. Although real causality is not possible to find, there are some methods to measure quantitative amounts of causality-like properties. Granger causality test is one of the popular methods to measure such amounts [6], which is based on a linear regression method. To measure the casual relation which cannot be captured by the covariance, transfer entropy (TE) (also known as directed transinformation [7], and directed information [15]) was proposed several decades ago and recently started to draw much attention [7,8]. Transfer entropy has been applied to many research areas like neuroscience and cognitive science [9,10], where it measures how much information flows between variables over the whole time series. Transfer entropy was normalized by the total information of the variable after removing the bias effect, which leads to *normalized transfer entropy (NTE)* [11].

NTE does not measure the flow locally in time, while localized information flow could be more helpful in understanding the relationship. For example, after finding out that neuron *A* affects

neuron *B* over the whole time, it can be more informative to figure when and how much *A* affects *B*. Recently several localized versions of transfer entropy have been proposed [12–14], but normalized transfer entropy has not been localized. So, in this paper, we localize NTE after removing the bias effect as the previous NTE does, which leads to local NTE. The rationale for the need of localization of NTE over TE is that NTE can make the transfer information comparable in different environments or between different pairs of variables without bias effects, while TE cannot. More discussions about normalization can be found in [11].

Another issue in transfer entropy is that when data set is not dense enough, there is a sampling problem as pointed out by [9], since the entropies are calculated based on the probabilities. That is, when the data samples are sparse in the data space, the probability density estimation is too much biased to a few samples. To overcome this sparseness problem, a priori distribution can be used [9]. In this paper, we propose to regularize the local NTE based on Dirichlet distribution as a prior distribution for the probabilities, which leads to *regularized local NTE*.

The rest of this paper is organized as follows. First, we briefly review some previous works including transfer entropy and its localized version in Section 2. Then in Section 3, we define local NTE and regularize it to avoid sparseness problem. Experiment results with synthetic and real-world data sets show how our algorithm works in Section 4, followed by conclusion in Section 5.

2. Previous works

Transfer entropy (TE) was proposed decades ago by Marko [7], and recently started to draw much attention since Schreiber's

E-mail address: heeyoul@gmail.com

work [8]. It has been successfully applied to understand how much information flows between variables (or channels) in many research areas [15,11,10]. In this section, we briefly review some previous works including all the steps in transfer entropy and its localized version.

2.1. Symbolization and probability mass function

When the variables in the data set are continuous variables, to make the problem simple, symbolization of the variables can be applied so that the data can be considered as a set of discrete sequences. As in the SAX algorithm [16], we use the distribution of the piecewise aggregate approximates (PAAs) and make a uniformly distributed symbol set based on the histogram of the PAAs. If the data set is given in discrete sequences, symbolization is not necessary.

Let $A = \{a_1, a_2, \dots, a_K\}$ be a set of the symbols of the two discrete variables along the time $t = 1, \dots, N$: $\mathbf{x} = [x_1, x_2, \dots, x_N]$ and $\mathbf{y} = [y_1, y_2, \dots, y_N]$, and $x_t \in A$ and $y_t \in A$. $p(x)$ be the probability mass function of \mathbf{x} . Basically, the probability functions are obtained by frequency of the symbols, that is, $p(x) = n_x/N$, where n_x is the number of times that symbol x happens in the sequence.

2.2. Transfer entropy

There are several entropy measures related to TE [7]. Shannon entropy of \mathbf{x} and conditional entropy (or entropy rate) of \mathbf{x} given \mathbf{x}^p are, respectively, defined by

$$H(\mathbf{x}) = - \sum_{x \in A} p(x) \log p(x), \quad (1)$$

$$H(\mathbf{x}|\mathbf{x}^p) = - \sum_{x, \mathbf{x}^p \in A} p(x, \mathbf{x}^p) \log p(x|\mathbf{x}^p), \quad (2)$$

where \mathbf{x}^p is the past values of \mathbf{x} . Conditional entropy means the total information of current \mathbf{x} knowing the previous values. Free information of \mathbf{x} is given by \mathbf{x}^p and \mathbf{y}^p , $H(\mathbf{x}|\mathbf{x}^p, \mathbf{y}^p)$ is defined in a similar way and it means the information of current \mathbf{x} is independent of \mathbf{y} :

$$H(\mathbf{x}|\mathbf{x}^p, \mathbf{y}^p) = - \sum_{x, \mathbf{x}^p, \mathbf{y}^p \in A} p(x, \mathbf{x}^p, \mathbf{y}^p) \log p(x|\mathbf{x}^p, \mathbf{y}^p). \quad (3)$$

Transfer entropy from \mathbf{y} to \mathbf{x} is given by

$$TE(\mathbf{y}, \mathbf{x}) = H(\mathbf{x}|\mathbf{x}^p) - H(\mathbf{x}|\mathbf{x}^p, \mathbf{y}^p), \quad (4)$$

which means the information of current \mathbf{x} is coming from \mathbf{y} but not from the previous values of \mathbf{x} itself.

2.3. Normalized transfer entropy

Transfer entropy might include another kind of information from \mathbf{y} to \mathbf{x} that is not what we want to count [17]. For example, the future of \mathbf{y} can incidentally give some amount of information to \mathbf{x} , especially when the data samples are sparse. To remove such inevitable information from \mathbf{y} to \mathbf{x} in the transfer entropy, we need to keep the distribution of \mathbf{y} as independent of \mathbf{x} as possible over the whole sequence. If \mathbf{y} is independent of \mathbf{x} , $H(\mathbf{x}|\mathbf{x}^p, \mathbf{y}^p) = H(\mathbf{x}|\mathbf{x}^p)$ which makes $TE(\mathbf{y}, \mathbf{x}) = 0$. That is, even after shuffling (or reordering randomly) \mathbf{y} , if $TE(\mathbf{y}^s, \mathbf{x}) > 0$, where \mathbf{y}^s is the shuffled sequence of \mathbf{y} , that is not the information we want to measure. So we remove that from the real $TE(\mathbf{y}, \mathbf{x})$ as in [17]. Instead of shuffling, we can shift \mathbf{y} assuming that if a signal is shifted enough, there is no mutual information between two signals.

Finally, normalized transfer entropy (NTE) from \mathbf{y} to \mathbf{x} [11] is given by

$$NTE(\mathbf{y}, \mathbf{x}) = \frac{TE(\mathbf{y}, \mathbf{x}) - TE(\mathbf{y}^s, \mathbf{x})}{H(\mathbf{x}|\mathbf{x}^p)}, \quad (5)$$

where $TE(\mathbf{y}^s, \mathbf{x})$ is a bias term with a shuffled (or shifted) sequence of \mathbf{y} . That is, the bias term indicates the amount of information flowing even when the sequence is shuffled. NTE represents the ratio of information in the current \mathbf{x} transferring from \mathbf{y} to the total information of \mathbf{x} after removing the bias effect. The rationale for the need of the bias term and normalization was discussed in [11]. As briefly described, NTE makes comparison between transfer entropies possible.

2.4. Local transfer entropy

Given two sequences of symbols, TE or NTE gives us just one value which indicates how much information is flowing between the two sequences over the whole time. To understand the information locally in space or time, local transfer entropy can be used [12,13]. Local TE is almost the same as Eq. (4), except not taking the expectation. The local entropies at time $t \in [1, N]$ are defined as follows.

Local entropy, also known as surprise or information, of \mathbf{x} at a specific time t is defined by

$$H_t(\mathbf{x}) = -\log p(x_t), \quad (6)$$

where $\sum_t H_t(\mathbf{x}) = H(\mathbf{x})$ in Eq. (1). That is, $H_t(\mathbf{x})$ indicates the contribution to the entropy $H(\mathbf{x})$ at a specific time t . The conditional entropy is localized in the same way as follows:

$$H_t(\mathbf{x}|\mathbf{x}^p) = -\log p(x_t|\mathbf{x}_{t-1}), \quad (7)$$

where the range of t is $[2, N]$. Here, we consider just the order 1 Markov process to reduce the computation complexity so that \mathbf{x}^p is replaced by \mathbf{x}_{t-1} . Also, free information of \mathbf{x} is localized by

$$H_t(\mathbf{x}|\mathbf{x}^p, \mathbf{y}^p) = -\log p(x_t|\mathbf{x}_{t-1}, \mathbf{y}_{t-1}). \quad (8)$$

Finally, local transfer entropy from \mathbf{y} to \mathbf{x} [12] is defined by

$$TE_t(\mathbf{y}, \mathbf{x}) = H_t(\mathbf{x}|\mathbf{x}^p) - H_t(\mathbf{x}|\mathbf{x}^p, \mathbf{y}^p), \quad (9)$$

and $\sum_t TE_t(\mathbf{y}, \mathbf{x}) = TE(\mathbf{y}, \mathbf{x})$, where $TE_t(\mathbf{y}, \mathbf{x})$ indicates the contribution to the transfer entropy $TE(\mathbf{y}, \mathbf{x})$ at a specific time t .

3. Proposed algorithm

3.1. Localization of normalized transfer entropy

Although NTE is more appropriate than TE to compare many pairs, there is no localized version of NTE. So, in light of localizing TE, we localize NTE, which leads to *local NTE*, as follows:

$$NTE_t(\mathbf{y}, \mathbf{x}) = \frac{TE_t(\mathbf{y}, \mathbf{x}) - \frac{1}{N-1} \sum_t TE_t(\mathbf{y}^s, \mathbf{x})}{\sum_t H_t(\mathbf{x}|\mathbf{x}^p)}, \quad (10)$$

where

$$\sum_t TE_t(\mathbf{y}^s, \mathbf{x}) = TE(\mathbf{y}^s, \mathbf{x}),$$

$$\sum_t H_t(\mathbf{x}|\mathbf{x}^p) = H(\mathbf{x}|\mathbf{x}^p).$$

To make the shuffled effect consistent over the whole time series, we take the average of $TE_t(\mathbf{y}^s, \mathbf{x})$ over the whole time which is $TE(\mathbf{y}^s, \mathbf{x})/(N-1)$, and we normalize the equation by $\sum_t H_t(\mathbf{x}|\mathbf{x}^p)$ to make sure that $\sum_t NTE_t(\mathbf{y}, \mathbf{x}) = NTE(\mathbf{y}, \mathbf{x})$ as follows:

$$\begin{aligned} \sum_t NTE_t(\mathbf{y}, \mathbf{x}) &= \frac{\sum_t TE_t(\mathbf{y}, \mathbf{x}) - \sum_t \frac{1}{N-1} TE(\mathbf{y}^s, \mathbf{x})}{H(\mathbf{x}|\mathbf{x}^p)} \\ &= \frac{TE(\mathbf{y}, \mathbf{x}) - TE(\mathbf{y}^s, \mathbf{x})}{H(\mathbf{x}|\mathbf{x}^p)} \\ &= NTE(\mathbf{y}, \mathbf{x}). \end{aligned}$$

Download English Version:

<https://daneshyari.com/en/article/407872>

Download Persian Version:

<https://daneshyari.com/article/407872>

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