

Causal DAG for genes obtained via Mappa Mundi algorithm

Robert R. Tucci
tucci@ar-tiste.com

March 14, 2025

Abstract

1 Introduction

This paper can be viewed as an application and further refinement of the Mappa Mundi (MM) algorithm. In this case, we apply it to finding what are called Gene Regulatory Networks (GRN) and Network Motifs in the Genomics and Systems Biology literature (Ref.[1]).

The MM algorithm was first proposed in Ref.[4] for DAG Extraction From Text (DEFT). In Ref.[4], it was used to compare 3 P.G. Wodehouse short stories and the scripts of 3 Pixar movies and extract from those causal DAGs. In general, the MM algorithm can extract causal DAGs from 2 or more text files, as long as each of those text files recounts actions in chronological order. So it will work with time stamped lab notebooks or medical records, but it won't work with textbooks or fiction with time travel or flashback shenanigans.

After Ref.[4], the MM algorithm was later applied in Ref.[3] to extracting causal DAGs from FitBit times series data.

In this paper, we use the MM algorithm to extract DAGs from time series data for concentrations of gene expressions and transcription factors. The DAGs we obtain are called GRN in the Genomics literature. GRN are a special case of autoregulation networks.

Autoregulation Networks (AN) are discussed in the chapter entitled "Autoregulation Networks (Network Motifs)" of my book Ref.[2].

AN are a special case of Dynamical systems (DS). DS are discussed in the chapter entitled "Dynamical Systems" of my book Ref.[2].

2 Comparing 2 TS Records

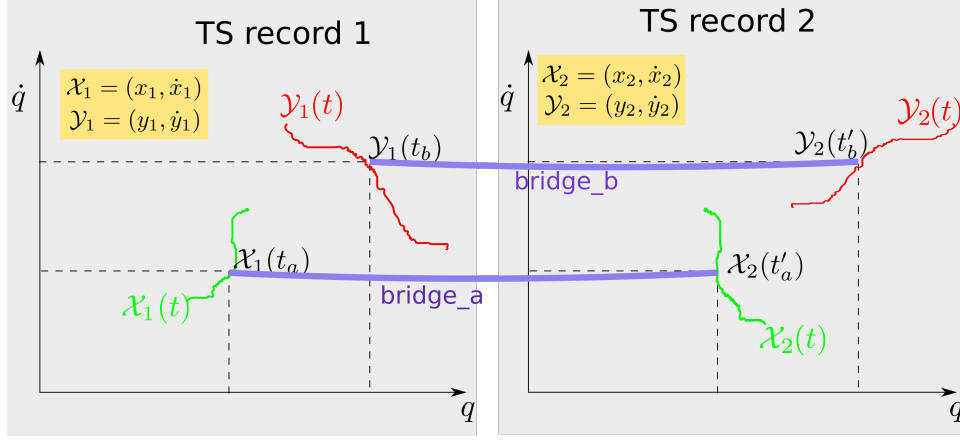


Figure 1: Causal bridges a and b spanning phase planes for TS records 1 and 2.

We will use the term **time series (TS) record** to refer to a data file such as a spreadsheet with the first column giving time increasing downwards and additional column giving the values of q_i (a state variable) and \dot{q}_i (time derivative of q_i) for $i = 1, 2, \dots, N$ at the time indicated by the first column. Any q_i or \dot{q}_i is called a **state variable**. The multidimensional space $(q_i)_{i=0}^N$ is called **configuration space**. The multidimensional space $(q_i, \dot{q}_i)_{i=0}^N$ is called **phase space**. A two dimensional space (q_i, \dot{q}_i) for any i is called a **phase plane**.

The TS record may contain initially only a column for time and columns for configuration space, if the change in time between rows is small. In that case we can subtract two consecutive q_i readings and divide by the difference in times to obtain the \dot{q}_i column and same row.

For this paper, each q_i represents either a **translation factor (TF)** concentration, or a **gene expression (GE)** concentration.¹

Suppose x and y are any two q_i . For $\xi \in \{x \rightarrow y, y \rightarrow x\}$, let

n_{acc}^ξ : number of arrows, initially zero

n_{rej}^ξ : number of arrows rejected, initially zero

$N^\xi = n_{acc}^\xi + n_{rej}^\xi$: number of arrows detected

$p_{acc}^\xi = \frac{n_{acc}^\xi}{n_{acc}^\xi + n_{rej}^\xi}$: probability of causal arrow, initially zero.

N^* : threshold value for N^ξ

p_{acc}^* : threshold value for p_{acc}^ξ

The MM algorithm for finding an AN from 2 TS records, consists of the following steps:

¹The terms “translation factor” and “gene expression” are both defined in the AN chapter of my free book Ref.[2]

1. **Compare 2 TS records and score arrows between any two autoregulations nodes \boxed{x} and \boxed{y}**

Consider Fig.1. In that figure, suppose the two ends of bridge a are equal: $\mathcal{X}_1(t_a) \approx \mathcal{X}_2(t'_a)$ and the two ends of bridge b are equal too: $\mathcal{Y}_1(t_b) \approx \mathcal{Y}_2(t'_b)$.²

At the very least, one must store (unless they are the default value zero) the current values of n_{acc}^ξ and n_{rej}^ξ for $\xi \in \{x \rightarrow y, y \rightarrow x\} = \mathcal{A}$, where x and y are any two q_i .

- if $t_a < t_b$ and $t'_a < t'_b$ (bridges are parallel in time)³

$$\begin{cases} n_{acc}^{x \rightarrow y} + + \\ N^{x \rightarrow y} + + \end{cases} \quad (1)$$

- if $t_a > t_b$ and $t'_a > t'_b$ (bridges are parallel in time)

$$\begin{cases} n_{acc}^{y \rightarrow x} + + \\ N^{y \rightarrow x} + + \end{cases} \quad (2)$$

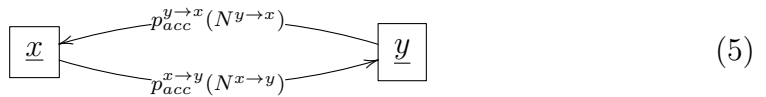
- if $t_a < t_b$ and $t'_a > t'_b$ (bridges are crossing in time)

$$\begin{cases} n_{rej}^{x \rightarrow y} + + \\ N^{x \rightarrow y} + + \end{cases} \quad (3)$$

- if $t_a > t_b$ and $t'_a < t'_b$ (bridges are crossing in time)

$$\begin{cases} n_{rej}^{y \rightarrow x} + + \\ N^{y \rightarrow x} + + \end{cases} \quad (4)$$

2. **Draw DAG**



If x and y are any two q_i , draw an arrow from autoregulation \boxed{x} to autoregulation \boxed{y} iff both $p_{acc}^{x \rightarrow y} > p^*$ and $N^{x \rightarrow y} > N^*$ are true.

Likewise, draw an arrow from autoregulation \boxed{y} to autoregulation \boxed{x} iff both $p_{acc}^{y \rightarrow x} > p^*$ and $N^{y \rightarrow x} > N^*$ are true.

When drawing an arrow, put the values p_{acc}^ξ and N^ξ over the arrow, where $\xi \in \mathcal{A}$. See Eq.(5) where this is done with variables. Do it with the values of those variables instead.

²By $\mathcal{X} \approx \mathcal{Y}$ we mean that both of these vectors are inside the same small bin or open ball of size given by a pre-specified precision.

³ $x++$ means add 1 to x

At first glance, Eq.(5) doesn't look like a DAG, because it has a cycle and DAGs are, by definition, acyclic. But Eq.(5) does indeed represent a DAG because, as explained in the AN chapter of Ref.[2], Eq.(5) represents this net:



which is acyclic.

3 Comparing More Than 2 TS Records

References

- [1] Uri Alon. *An introduction to systems biology: design principles of biological circuits*. Chapman and Hall/CRC, 2019.
- [2] Robert R. Tucci. Bayesuivus (free book). <https://github.com/rrtucci/Bayesuivus>.
- [3] Robert R. Tucci. <https://github.com/rrtucci/causalfitbit> (software and paper). <https://github.com/rrtucci/CausalFitbit>.
- [4] Robert R. Tucci. Mappa Mundi (software and paper). https://github.com/rrtucci/mappa_mundi.