

Unsupervised learning of chemical reaction networks from Time Series Data

Julien Martinelli

December, 19th 2018

François Fages

Annabelle Ballesta

Mechanistic Model Learning

The Machine Learning area provides tools to analyze time series data and yield predictions. Classical examples are Recurrent Neural Networks.

- While these predictions can be accurate, they do not come with an interpretation
- We say that the model is *Black Box*

On the contrary, Mechanistic Model Learning aims at achieving the same predictive results while being explainable

(*XAI : Explainable Artificial Intelligence*)

Some attempts at Mechanistic Model Learning

- DREAM3 (2008) - Network Inference Challenge
- Logic programming combined with prior knowledge on the network's structure allows to learn the boolean function responsible for each species

Boolean Network Identification from Perturbation Time Series Data combining Dynamics Abstraction and Logic Programming. L. Pauleve et al.

- Evolutionary Algorithms : based on the minimization of a fitness criterion measuring the difference between the observed data and the proposed mechanistic models

Inferring Reaction Networks using Perturbation Data. H. Sauro et al.

- TimeDelay-ARACNE: Reverse engineering of gene networks from time-course data by an information theoretic approach

P. Zoppoli et al.

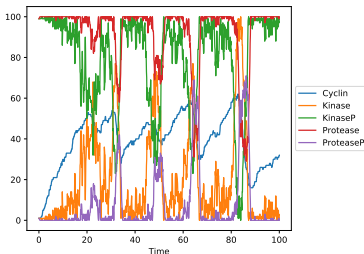
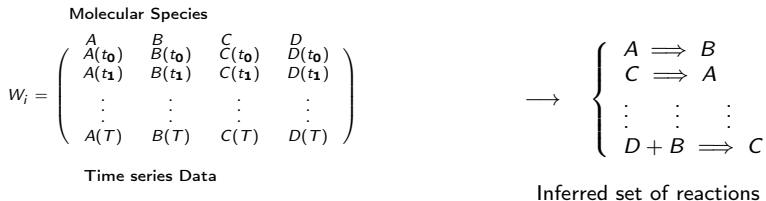
Toward Personalized Chronotherapies at INSERM U935

- In the biomedical case, predictions are required for instance to determine the optimal hour of drug delivery.
- Moreover, in the case of Personalized Medicine, we want to learn a model of the patient
- Learning a mechanistic model would give these predictions consistency through the understanding of the biological processes underneath
- We aim at learning not only parameters but also model structure

The Problem of Reaction Network Inference

Input : observed time-series data from biological experiments such as proteomic data.

Output : a set of reactions defining a model \mathcal{M} reproducing similar time-series data



Chemical Reaction Network (CRN)

Hypothesis : Stoichiometry coefficients are less or equal to 1.

Definition

A reaction j is a triplet (R_j, P_j, h_j)

R_j is the set of reactants

P_j the set of products

h_j is the rate function

A CRN is a set of reactions $\mathcal{M} = (R_j, P_j, h_j)_{1 \leq j \leq J}$

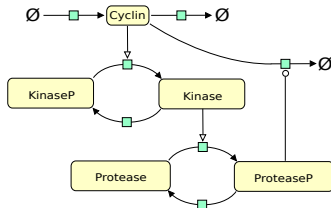
A catalyst is a species $B \in R_j \cap P_j$

Example

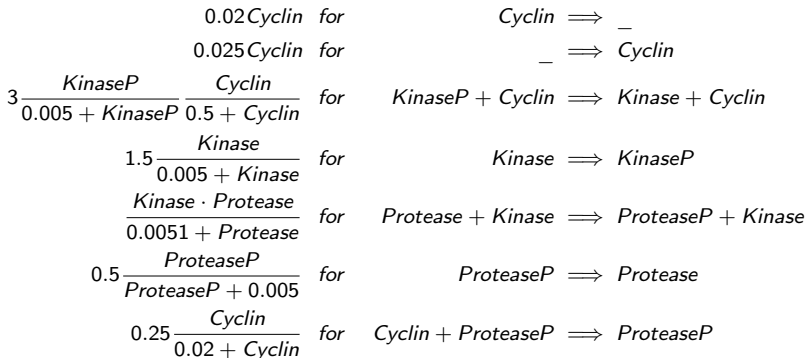
$R = \{A\}$ $P = \{B\}$ $h : x \mapsto k \cdot x$

$k * A \text{ for } A \implies B$

Simulated Data from Minimal Mitotic Oscillator

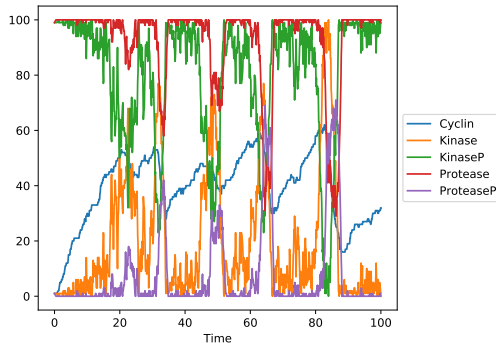


Goldbeter,
1991 -
Biomodels
Repository



Stochastic Simulation Algorithm

We consider stochastic simulation traces from an hidden model
(Continuous time Markov chain)

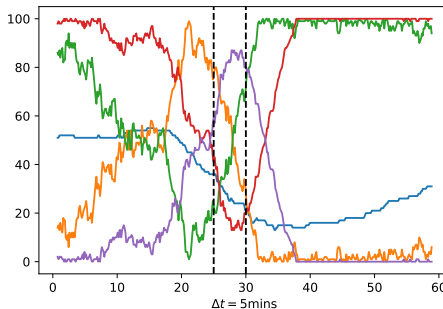


Numerical simulation using the *Gillespie Algorithm* - Minimal Mitotic Oscillator

Adding Subsampling to traces

Subsampling hypothesis : We do not observe every transition from the Markov chain simulation, only a sample of them every $\Delta t = 5mins$

→ Therefore we do not observe reactions one by one but *macro transitions*.



Adding Noise to traces

Multiplicative Gaussian noise is added to the predecessor state and the successor state.

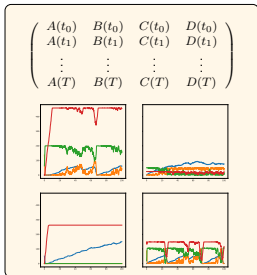
$$X_{meas} = X_{sim} * e^w \text{ where } w \sim \mathcal{N}(0, \sigma) \text{ and } \sigma = 0.003$$

→ A species more present than another will then be more noisy.

Noise is then suppressed by rounding to the closest integer.

Workflow of the learning algorithm

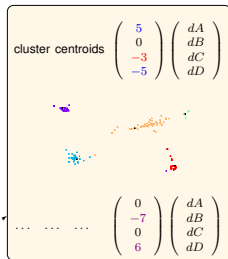
Time series data



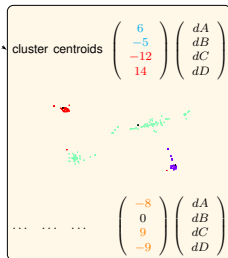
Observed set of macro transitions

$$\{(X_i, X_{i+1}) \mid \forall i \in 0, \dots, T-1\}$$

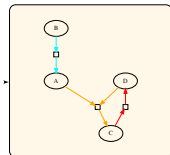
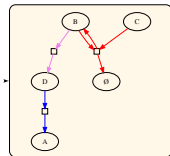
Clustering of observed macro transitions



$$k = 2, 3, \dots, 20$$



Inferred Reactions



Clustering of the Observed Macro Transitions

Finite differences between the successor state and the predecessor state are computed while predecessor and successor state are stored.

Example

$$\begin{pmatrix} 32 \\ 19 \\ 88 \\ 57 \\ 6 \end{pmatrix} \rightarrow \begin{pmatrix} 32 \\ 24 \\ 82 \\ 49 \\ 14 \end{pmatrix} \begin{pmatrix} Cyclin \\ Kinase \\ KinaseP \\ ProteaseP \\ Protease \end{pmatrix} \quad \begin{pmatrix} 0 \\ 5 \\ -6 \\ -8 \\ 8 \end{pmatrix} \begin{pmatrix} Cyclin \\ Kinase \\ KinaseP \\ ProteaseP \\ Protease \end{pmatrix}$$

Macro transition (P_i, S_i) and associated difference vector $\delta_i = S_i - P_i$

Macro Transitions Clustering based on difference vectors

Clustering as a way to extract information from the dataset

→ We choose the *K-medoids* algorithm with the squared euclidean distance

- Start with randomly chosen centroids and update the clusters :

$$C_k = \{\delta_i \text{ s.t. } \operatorname{argmin}_{\delta \in M} \|\delta - \delta_i\|_2^2 = M_k\}$$

- Then the **centroids** are updated

$$M_k = \operatorname{argmin}_{\delta \in C_k} \frac{1}{|C_k|} \sum_{\delta_i \in C_k} \|(\delta - \delta_i)\|_2^2$$

- Repeat until the partitioning reaches a stable state

Centroids are actual members of the dataset

Inferring reactions from the centroids

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \textit{Cyclin} \\ \textit{Kinase} \\ \textit{KinaseP} \\ \textit{ProteaseP} \\ \textit{Protease} \end{pmatrix}$$

$$x^* = \underset{\textit{species}}{\operatorname{argmin}} x^i$$

Inferring reactions from the centroids

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \textit{Cyclin} \\ \textit{Kinase} \\ \textit{KinaseP} \\ \textit{ProteaseP} \\ \textit{Protease} \end{pmatrix}$$

$$x^* = \underset{\textit{species}}{\operatorname{argmin}} x^i$$

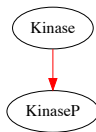
$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \textit{Cyclin} \\ \textit{Kinase} \\ \textit{KinaseP} \\ \textit{ProteaseP} \\ \textit{Protease} \end{pmatrix}$$

$$\text{Collect every } x^i \text{ s.t. } \left| \frac{x^*}{x^i} \right| \leq \alpha$$

Inferring reactions from the centroids

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

$$x^* = \underset{\text{species}}{\operatorname{argmin}} x^i$$



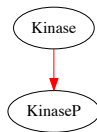
$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Collect every x^i s.t. $|\frac{x^*}{x^i}| \leq \alpha$

Inferring reactions from the centroids

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

$$x^* = \underset{\text{species}}{\operatorname{argmin}} x^i$$



$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Collect every x^i s.t. $|\frac{x^*}{x^i}| \leq \alpha$

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Inferring reactions from the centroids

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

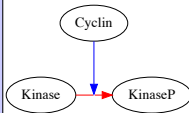
$$x^* = \underset{\text{species}}{\operatorname{argmin}} x^i$$



$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Collect every x^i s.t. $|\frac{x^*}{x^i}| \leq \alpha$

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$



Catalyst Candidate
with 0 variation

Inferring reactions from the centroids

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

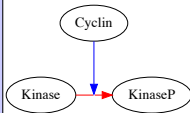
$$x^* = \underset{\text{species}}{\operatorname{argmin}} x^i$$



$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Collect every x^i s.t. $|\frac{x^*}{x^i}| \leq \alpha$

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$



Catalyst Candidate
with 0 variation

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Set species to 0 and start over.

Inferring reactions from the centroids

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

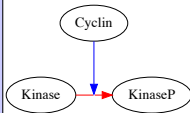
$$x^* = \underset{\text{species}}{\operatorname{argmin}} x^i$$



$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Collect every x^i s.t. $|\frac{x^*}{x^i}| \leq \alpha$

$$\begin{pmatrix} 0 \\ -7 \\ 6 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$



Catalyst Candidate
with 0 variation

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 2 \\ -3 \end{pmatrix} \quad \begin{pmatrix} \text{Cyclin} \\ \text{Kinase} \\ \text{KinaseP} \\ \text{ProteaseP} \\ \text{Protease} \end{pmatrix}$$

Set species to 0 and start over.

Results on Minimal Mitotic Oscillator (Goldbeter, 1991)

$R_1 : \text{ProteaseP} \Rightarrow \text{Protease}$

$R_2 : \text{Kinase} + \text{Protease} \Rightarrow \text{ProteaseP} + \text{Kinase}$

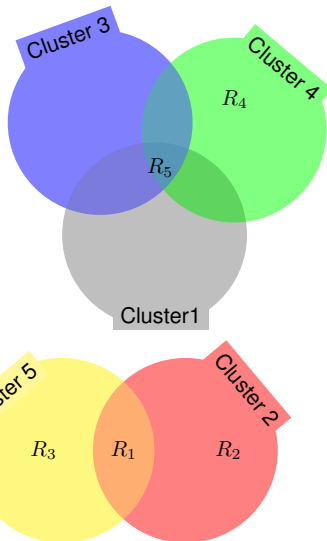
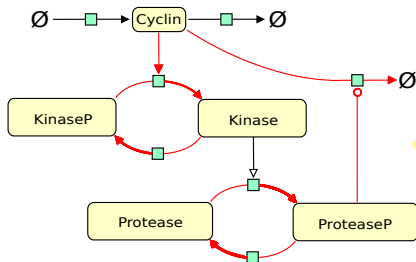
$R_3 : \text{Cyclin} + \text{KinaseP} \Rightarrow \text{Kinase} + \text{Cyclin}$

$R_4 : \text{ProteaseP} + \text{Cyclin} \Rightarrow \text{ProteaseP}$

$R_5 : \text{Kinase} \Rightarrow \text{KinaseP}$

$R_6 : - \Rightarrow \text{Cyclin}$

$R_7 : \text{Cyclin} \Rightarrow -$



Model Selection Step

- Choosing the right aggregated network amounts to choosing the optimal number of clusters k
- The reaction inference algorithm outputs a set of reactions, defining a generative model $\hat{\mathcal{M}}$
- Model quality can be assessed by comparing the distribution of \mathcal{M} to the one described by $\hat{\mathcal{M}}$

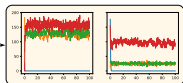
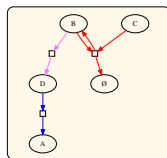
Model selection protocol

Learned
Reactions
Networks

New Simulations

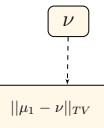
Compute
Distribution of
difference vectors

Comparison
with observed
data's distribution



$$Y_i^1 = X_{i+1}^1 - X_i^1$$

$$\mu_1 = \frac{1}{n-1} \sum_{i=1}^{n-1} \mathbb{1}_{\{Y_i \in \cdot\}}$$



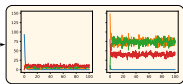
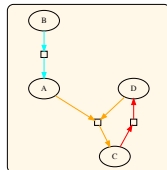
⋮

⋮

⋮
 $k = 2, 3, \dots, 20$
⋮

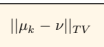
⋮

$$\mu^* = \operatorname{argmin}_{\mu} ||\mu - \nu||_{TV}$$

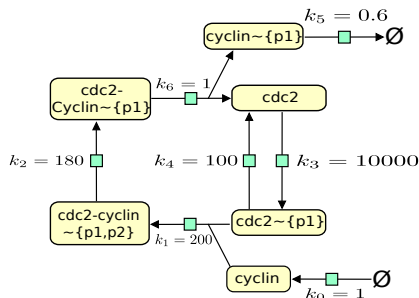


$$Y_i^k = X_{i+1}^k - X_i^k$$

$$\mu_k = \frac{1}{n-1} \sum_{i=1}^{n-1} \mathbb{1}_{\{Y_i \in \cdot\}}$$



6 variables Cell Cycle Model (Tyson, 1991)



$$Cdc2 \Rightarrow Cdc2 \sim \{p1\}$$

$$Cdc2 \sim \{p1\} \Rightarrow Cdc2$$

$$Cdc2-Cyclin \sim \{p1, p2\} \Rightarrow Cdc2-Cyclin \sim \{p1\}$$

$$Cdc2 \sim \{p1\} + Cyclin \Rightarrow Cdc2-Cyclin \sim \{p1, p2\}$$

$$_ \Rightarrow Cyclin$$

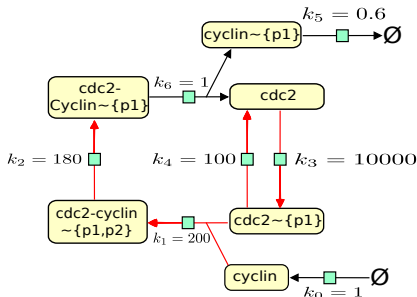
$$Cyclin \sim \{p1\} \Rightarrow _$$

$$Cdc2-Cyclin \sim \{p1\} \Rightarrow Cyclin \sim \{p1\} + Cdc2$$

$$Cdc2-Cyclin \sim \{p1, p2\} + 2 * Cdc2-Cyclin \sim \{p1\} \\ \Rightarrow 3 * Cdc2-Cyclin \sim \{p1\}$$

	mean transition difference vector	max transition difference vector
Cyclin	3.3	45
Cyclin~{p1}	1.02	2
Cdc2	53.54	853
Cdc2~{p1}	50.3	840
Cdc2-Cyclin~{p1}	6.43	123
Cdc2-Cyclin~{p1,p2}	4.7	122

6 variables Cell Cycle Model (Tyson, 1991)



$Cdc2 \Rightarrow Cdc2 \sim \{p1\}$

$Cdc2 \sim \{p1\} \Rightarrow Cdc2$

$Cdc2-Cyclin \sim \{p1, p2\} \Rightarrow Cdc2-Cyclin \sim \{p1\}$

$Cdc2 \sim \{p1\} + Cyclin \Rightarrow Cdc2-Cyclin \sim \{p1, p2\}$

$_ \Rightarrow Cyclin$

$Cyclin \sim \{p1\} \Rightarrow _$

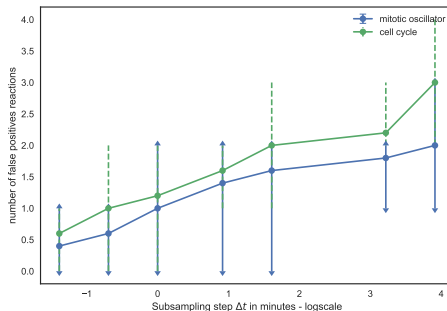
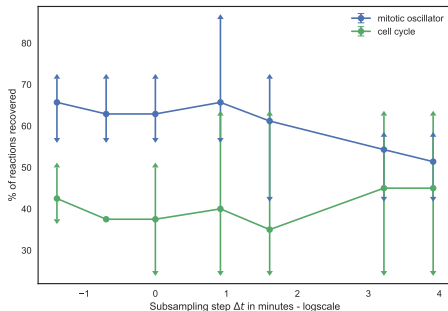
$Cdc2-Cyclin \sim \{p1\} \Rightarrow Cyclin \sim \{p1\} + Cdc2$

$Cdc2-Cyclin \sim \{p1, p2\} + 2 * Cdc2-Cyclin \sim \{p1\}$

$\Rightarrow 3 * Cdc2-Cyclin \sim \{p1\}$

- Reaction recovered are precisely the four fastest ones hence those with the highest probability to occur when possible
(False Positive : 0%, False Negative : 50%)
- The gap between kinetic parameters values results in a slow/fast dynamic, a limit of the stochastic approach.

Subsampling effect on learning



- Rules including catalysts are inferred without the latter : more false positives
- Scarce reactions such as $_ \implies A$ are inferred : less false negatives
- As the subsampling step grows, more false positives appear.

Conclusion and Perspectives

- Unsupervised reaction inference algorithm dealing with subsampled and noisy time-series data
- The algorithm finds original reactions
- But also misses other original reactions (*false negatives*)
→ high precision but low recall

Perspectives:

- Case where not all species are observed (Elisabeth Degrand's Master Thesis - Evolving CRN)