



**HAL**  
open science

# Simulation d'expériences d'intervention biologique dans des cellules cancéreuses à partir de données temporelles d'expression de gènes

Anouk Rago, Nicolas Champagnat, Anne Gégout-Petit, Laurent Vallat

## ► To cite this version:

Anouk Rago, Nicolas Champagnat, Anne Gégout-Petit, Laurent Vallat. Simulation d'expériences d'intervention biologique dans des cellules cancéreuses à partir de données temporelles d'expression de gènes. 54es Journées de Statistique de la SFdS (JdS 2023), Société Française de Statistique, Jul 2023, Bruxelles, Belgique. hal-04408929

**HAL Id: hal-04408929**

**<https://inria.hal.science/hal-04408929>**

Submitted on 22 Jan 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# SIMULATION D'EXPÉRIENCES D'INTERVENTION BIOLOGIQUE DANS DES CELLULES CANCÉREUSES À PARTIR DE DONNÉES TEMPORELLES D'EXPRESSION DE GÈNES

Anouk Rago <sup>1</sup> & Nicolas Champagnat <sup>2</sup> & Anne Gégout-Petit <sup>3</sup> & Laurent Vallat <sup>4</sup>

<sup>1</sup> *Université de Lorraine, CNRS, Inria, IECL, France, anouk.rago@univ-lorraine.fr*

<sup>2</sup> *Université de Lorraine, CNRS, Inria, IECL, France, nicolas.champagnat@univ-lorraine.fr*

<sup>3</sup> *Université de Lorraine, CNRS, Inria, IECL, France, anne.gegout-petit@univ-lorraine.fr*

<sup>4</sup> *Université de Strasbourg, INSERM, CHU de Strasbourg, France, laurent.vallat@inserm.fr*

**Résumé.** En mathématiques comme en biologie, les interactions entre les gènes sont généralement représentées sous la forme d'un graphe orienté où les nœuds représentent les différents gènes et les arêtes une relation de dépendance entre deux gènes. Afin d'inférer ce réseau à partir de données dynamiques d'expression de gènes, de nombreuses techniques ont été développées ces dernières années. On peut citer par exemple l'utilisation de modèles graphiques gaussiens, de modèles linéaires avec inférence pénalisée ou encore des forêts aléatoires. À partir d'un graphe inféré grâce à un modèle et des données temporelles d'expression de gènes, nous nous intéressons à la modélisation d'une expérience biologique dite de *silencing*, consistant à réduire fortement l'expression de certains gènes dans la cellule, et à mesurer l'impact de ce *silencing* sur un ensemble de gènes appelés "cibles". Ces expériences sont un espoir pour réduire la prolifération cellulaire incontrôlée qui survient dans les cellules leucémiques. En prenant en compte les spécificités de notre problème, notamment le faible nombre de données médicales et la structure du graphe inféré, nous proposons de développer et comparer deux méthodes différentes pour simuler mathématiquement ce *silencing*. Celles-ci seront testées numériquement sur des données temporelles simulées dans le cas d'un modèle linéaire standard.

**Mots-clés.** Réseaux de gènes, modélisation d'expérience biologique, simulation numérique, données d'expression de gènes

**Abstract.** In mathematics as in biology, the interactions between genes are usually depicted as an oriented graph, where the vertices represent the different genes and the edges indicate a dependence relationship between two nodes. A lot of methods have been recently developed in order to infer this network from temporal genomic data. Graphical gaussian models, penalized linear models or random forests can be cited as examples.

Given a graph, which was inferred thanks to a model and temporal data, we aim to model a biological experiment known as *silencing*. This involves reducing the level of expression of a group of genes in order to observe the impact on another set of genes, known as "targets". These experiments are a hope for physicians to reduce the cell proliferation which occurs in cancer cells. We would like to develop and compare two methods to mathematically simulate the phenomenon of *silencing*. These methods will take into account the lack of medical data and the structure of the graph. Their performance will then be tested numerically on temporal data simulated by a standard linear model.

**Keywords.** Gene network, silencing modelling, numerical simulations, gene expression data

## 1 Introduction

Lorsqu'elles sont soumises à des modifications d'environnement, comme l'activation d'un récepteur membranaire par exemple, les cellules de notre corps disposent d'un mécanisme de réponse adapté. Celui-ci se caractérise par l'activation de voies de signalisation à l'intérieur de la cellule : les gènes s'activent les uns les autres via l'augmentation de leur niveau d'expression (nombre d'ARN messenger codant pour ce gène), et modulent l'abondance des protéines de fonctions, qui sont les véritables actrices face à l'évènement initial.

En mathématiques comme en biologie, ce processus est généralement représenté sous la forme d'un graphe orienté dont les nœuds représentent les différents gènes, et où les arêtes indiquent des liens d'interactions entre deux gènes, dont la signification dépend du modèle sous-jacent utilisé pour construire le graphe. Ces réseaux, appelés graphes d'interactions, sont cruciaux en biologie : mieux connaître les différents liens entre les gènes permet de mieux comprendre les mécanismes cellulaires et leur dérégulation dans des pathologies tumorales. Ces dernières années, de nombreuses techniques d'inférence de réseaux de gènes basées sur l'utilisation de données d'expression temporelles ont été développées. Le livre de *Sanguinetti et al* (2019) présente certaines de ces techniques qui s'appuient sur différentes méthodes mathématiques : les modèles graphiques gaussiens, les réseaux bayésiens dynamiques ou encore les forêts aléatoires. Au delà de la connaissance biologique apportée par la construction des réseaux de gènes, ces derniers ainsi que les méthodes utilisées pour leur inférence peuvent être ensuite utilisés afin de simuler théoriquement certaines expériences biologiques intéressantes.

L'une d'entre elles, connue sous le nom de *silencing*, *Vallat et al* (2013), ou extinction de gènes en français, est un espoir pour le traitement du cancer et notamment des leucémies. Les cellules sanguines tumorales caractéristiques de cette maladie ont acquis une altération génique qui modifie les voies de signalisation habituelles en favorisant l'expression de certains gènes. Cela se traduit par une prolifération incontrôlée de cellules immatures dans le sang au détriment des cellules sanguines normales qui ne peuvent plus jouer leur rôle habituel. Le *silencing* permet de réduire en continu l'expression de certains gènes dans la cellule. Ainsi, en silençant le bon gène, fortement exprimé dans les premiers temps de l'expérience, les médecins peuvent espérer agir également sur l'expression d'un ou plusieurs gènes "cibles", responsables de la prolifération cellulaire.

Expérimentalement parlant, il est cependant impossible de tester un à un les 20000 gènes contenus dans une cellule pour trouver le ou les plus intéressants à silencer.

Notre objectif est donc de modéliser l'expérience de *silencing* d'un gène et d'en mesurer les effets sur les cellules cibles. Pour cela, nous proposons tout d'abord de choisir un modèle représentant les interactions entre les gènes, de l'inférer ainsi que le graphe correspondant en utilisant des données d'expression de gènes non silencées, avant de l'exploiter pour en

déduire un modèle silencé. Celui-ci sera ensuite utilisé pour en déduire l'effet du *silencing* sur les gènes cibles, comme représenté Figure 1.

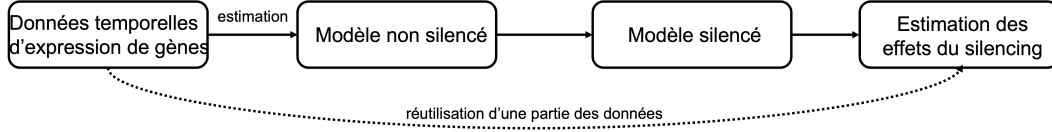


FIGURE 1 – Représentation schématique des différentes étapes de la modélisation mathématique du silencing.

Les données médicales nécessaires pour résoudre ce problème sont peu nombreuses et nous souhaitons construire une méthode d'estimation qui tient compte de cette particularité. Nous pourrions alors estimer théoriquement et numériquement l'effet du *silencing* sur un gène cible afin de sélectionner les candidats les plus prometteurs pour réduire la prolifération cellulaire.

## 2 Méthodes de simulation du *silencing*

Afin de simuler une expérience de *silencing*, nous supposons dorénavant disposer de données temporelles d'expressions de gènes. Ces données se présentent sous la forme de  $N$  matrices de taille  $p \times T$  où  $N$  est le nombre de patients, généralement moins d'une dizaine,  $p$  le nombre de gènes, proche de 20 000, et  $T$  le nombre de temps de mesures, généralement moins d'une dizaine :

$$\mathbf{X}^n = \begin{pmatrix} X_{1,0}^n & X_{1,1}^n & \cdots & X_{1,T}^n \\ X_{2,0}^n & X_{2,1}^n & \cdots & X_{2,T}^n \\ \vdots & \vdots & \ddots & \vdots \\ X_{p,0}^n & X_{p,1}^n & \cdots & X_{p,T}^n \end{pmatrix}$$

A partir de ces données, nous pourrions inférer un modèle ainsi que le graphe d'interaction correspondant. Ce modèle, construit à partir de données non silencées d'expressions de gènes, nous servira ensuite à simuler mathématiquement une expérience de silencing.

### 2.1 Structure du graphe et modèle

Nous inférons donc un modèle et un graphe orienté d'interactions  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  entre les gènes considérés. Nous supposons que l'orientation correspond à une influence avec décalage temporel : l'expression d'un gène à une date donnée influence celle d'un autre gène à la date suivante. Dans ce graphe, nous définissons les parents et les enfants directs d'un gène  $k$  de la manière suivante :

$$\begin{aligned} Pa^1(k) &= \{v \in \mathcal{V} \mid (v, k) \in \mathcal{E}\} \\ Enf^1(k) &= \{v \in \mathcal{V} \mid (k, v) \in \mathcal{E}\} \end{aligned}$$

Par récurrence, nous pouvons étendre cette définition pour obtenir la descendance et l'ascendance de degré  $t$  du gène  $k$  :

$$\begin{aligned} \forall t \geq 1, \quad Pa^t(k) &= \{v \in \mathcal{V} \mid \exists x \in Pa^{t-1}(k) \quad (v, x) \in \mathcal{E}\} \\ Enf^t(k) &= \{v \in \mathcal{V} \mid \exists x \in Enf^{t-1}(k) \quad (x, v) \in \mathcal{E}\} \end{aligned}$$

Nous notons également  $pa^n(k, t) = (X_{j,t}^n, j \in Pa^1(k))$  le vecteur des expressions des parents directs du gène  $k$  au temps  $t$  dans la  $n^{\text{ième}}$  matrice.

Concernant le modèle associé au graphe, nous supposons d'une part que celui-ci admet une forme mécaniste, c'est-à-dire qu'il s'écrit sous la forme d'une famille de fonctions  $(f_j)_{1 \leq j \leq p}$  dépendant de paramètres  $\beta$  qui devront être estimés, et d'autre part qu'il représente un processus de Markov : les expressions de gènes à des instants antérieurs à  $t - 1$  n'ont pas d'influence sur les expressions des gènes au temps  $t$ . Nous pouvons alors obtenir l'expression d'un gène à n'importe quel temps. Celle-ci s'écrit :

$$\forall n \in [1, N], \forall j \in [1, p], \forall t \in [1, T], \quad X_{j,t}^n = f_j(pa(j, t-1), \epsilon_t^j) \quad (1)$$

où  $\epsilon_t^j$  est un bruit. En itérant plusieurs fois cette équation, il est possible d'obtenir une expression en fonction des valeurs initiales des expressions de gènes à  $t = 0$ .

La structure du graphe joue un rôle prépondérant lors de la simulation du *silencing*. En effet, silencer un gène impactera directement le niveau d'expression de ses descendants tandis que les autres gènes seront préservés de toute modification. Deux remarques préalables peuvent alors être faites : le gène cible doit nécessairement se trouver dans la descendance du gène silencé et une inférence correcte des arêtes du graphe est primordiale pour une bonne estimation des effets du *silencing*.

## 2.2 Une première simulation du *silencing*

Le *silencing* a pour but de réduire à zéro le niveau d'expression d'un gène  $y$  précis dans les cellules. Néanmoins, il est expérimentalement impossible d'obtenir une disparition totale du gène considéré. C'est généralement pour cette raison qu'on considère qu'une fraction  $\alpha$  du nombre de copies d'ARNm codant pour ce gène  $y$  subsiste malgré tout dans l'échantillon après *silencing*. Mathématiquement, cela se traduit par l'ajout d'un coefficient multiplicateur  $\alpha$  dans l'équation (1) lorsque  $j = y$  :

$$\forall n \in [1, N], \forall j \in [1, p], \forall t \in [1, T], \quad X_{j,t}^n = \alpha_j f_j(pa(j, t-1), \epsilon_t^j) \quad (2)$$

avec  $\alpha_j = \alpha$  si  $j = y$ ,  $\alpha_j = 1$  sinon.

En utilisant les valeurs d'expression de gènes à  $t = 0$  du jeu de données et en appliquant la formule 2, on peut alors obtenir le niveau d'expression théorique de n'importe quel gène lorsqu' $y$  est silencé.

Plus concrètement, en notant  $z$  le gène cible qui nous intéresse et sur lequel on souhaite obtenir un effet,  $\hat{\beta}$  les paramètres du modèle inféré et  $x_0^n$  les données d'expressions de gènes à  $t = 0$  pour le patient  $n$ , nous calculons théoriquement ou estimons la quantité :

$$\mathbb{E}(z_T^{sil}(\hat{\beta}) | \mathbf{x}_0^n)$$

Lorsque les données initiales des différents patients suivent la même loi, nous pouvons également estimer les espérances moyennées :

$$\frac{1}{N} \sum_{n=1}^N \mathbb{E}(z_T^{sil}(\hat{\beta}) | \mathbf{x}_0^n)$$

### 2.3 S'adapter au faible nombre de données médicales, une deuxième simulation du *silencing*

Utiliser la méthode décrite ci-dessus nécessite d'avoir un modèle et un graphe bien estimés. En effet, les paramètres du modèle ainsi que les liens trouvés entre les gènes sont déterminants lors du calcul de l'espérance de  $z_T^{sil}$ . Or, nous sommes justement confrontés à une situation dans laquelle le faible nombre de données médicales risque d'influer énormément sur la qualité du modèle, et donc sur notre simulation du *silencing*.

Afin d'éviter de propager des erreurs en utilisant de façon systématique un modèle mal spécifié ou dont les paramètres sont potentiellement mal estimés, nous proposons une seconde méthode de simulation du *silencing* qui privilégie la réutilisation des données observées.

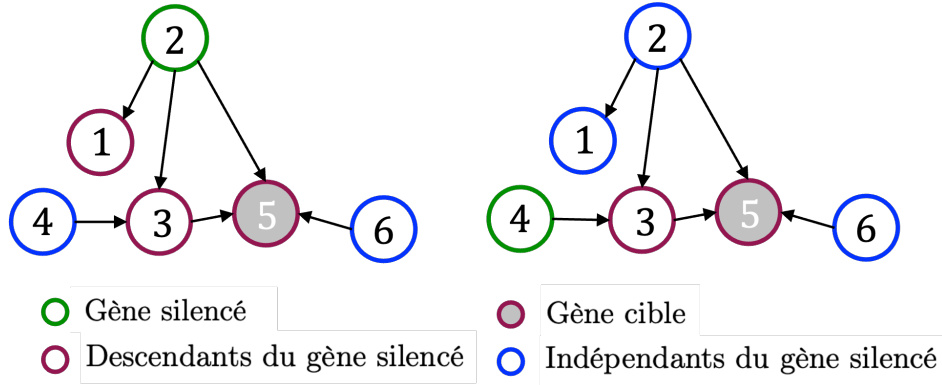


FIGURE 2 – Deux représentations d'un même réseau d'interactions de gènes. Selon le gène  $y$  considéré, on observe que les gènes impactés par son *silencing* sont différents : la structure du réseau permet de déterminer les liens de "parentés" entre les gènes. Pour un descendant de  $y$ , on utilisera le modèle inféré. Pour les autres, on exploitera directement le jeu de données.

A tout instant  $t$ , nous évaluons pour chaque gène si celui-ci appartient ou non à  $Enf^t(y)$ . Si oui, nous utilisons le modèle inféré pour déterminer le niveau d'expression. Si non, nous utilisons directement la donnée correspondante puisque le gène n'est pas impacté par le *silencing* de  $y$ . L'expression (2) devient alors :

$$\forall n \in [1, N], \forall j \in [1, p], \forall t \in [1, T], \quad X_{j,t,sil}^n = \begin{cases} \alpha_j f_j(pa(j, t-1), \epsilon_t^j) & \text{si } j \in \bigcup_{t'=0}^t Enf^{t'}(y) \\ x_{j,t}^n & \text{sinon} \end{cases} \quad (3)$$

Comme pour la première méthode de simulation du *silencing*, nous estimons la quantité :

$$\mathbb{E}(z_T^{sil}(\hat{\beta}) | \mathbf{x}_0^n, \mathbf{x}_{NE}^n)$$

où  $x_{NE}^n$  représente les données d'expressions des gènes qui ne sont pas des descendants de  $y$ . Nous pouvons également estimer :

$$\frac{1}{N} \sum_{n=1}^N \mathbb{E}(z_T^{sil}(\hat{\beta}) | \mathbf{x}_0^n, \mathbf{x}_{NE}^n)$$

## 2.4 Effet du *silencing*

Nous souhaitons quantifier l'effet du *silencing* du gène  $y$  sur le gène  $z$ . Cet effet peut être mesuré par la différence entre l'espérance du gène  $z$  calculée avec le modèle original et l'espérance de  $z$  calculée avec le modèle silencé, pour chacune des 2 méthodes décrites ci-dessus. Par exemple, pour la première méthode, l'influence du silencing est évaluée par :

$$\mathbb{E}(z_T(\hat{\beta}) | \mathbf{x}_0^n) - \mathbb{E}(z_T^{sil}(\hat{\beta}) | \mathbf{x}_0^n)$$

Calculer cette différence nous permet d'une part de vérifier si le *silencing* du gène  $y$  agit de façon positive (augmentation du niveau d'expression) ou négative (diminution du niveau d'expression) sur le gène cible  $z$ . D'autre part, elle permet également de mesurer *silencing* de  $y$  sur  $z$  en regardant sa valeur absolue.

## 3 Etude numérique d'un modèle linéaire

Nous effectuons des simulations numériques afin de vérifier l'intérêt des 2 méthodes de simulation du *silencing* présentées dans la section précédente.

### 3.1 Description du modèle d'inférence utilisé

Nous travaillons avec un modèle linéaire standard pour nos simulations numériques. Concrètement, le jeu de données est construit de la manière suivante :

1. Nous choisissons un réseau. Nous connaissons ainsi, pour chaque gène  $j \in [1, p]$ , son ascendance et sa descendance. Un poids  $\beta_k^j$  est également associé à chaque arête du graphe (donc à chaque lien de parenté).
2.  $\forall n \in [1, N], \forall t \in [1, T - 1], \forall j \in [1, p]$ , nous construisons de façon récursive les jeux de données grâce à la relation :

$$X_{j,t}^n = \sum_{k \in Pa^1(j)} \beta_k^j X_{k,t-1}^n + \epsilon_t^j \quad (4)$$

où  $\epsilon_t^j$  suit une loi normale centrée de variance  $\sigma_j$ . Cela revient à avoir  $p$  modèles de régressions linéaires, un pour chaque gène.

3. On obtient alors un jeu de données de  $N$  observations indépendantes, sous la forme de  $N$  matrices de taille  $p \times T$ .

Une fois le jeu de données construit, nous l'utilisons pour estimer la matrice des paramètres du modèle  $\beta$  par l'estimateur des moindres carrés et sélection de variables. Nous obtenons ainsi une matrice  $\hat{\beta}$ . Une arête de  $i$  vers  $j$  se trouve dans le graphe si  $\hat{\beta}_{ij} \neq 0$ .

### 3.2 Estimateurs de l'effet du silencing

Pour chacune des deux méthodes développées dans la section précédente, nous calculons les espérances du niveau d'expression du gène  $z$  au temps temps, avec ou sans silencing du gène  $y$ .

**Méthode 1 :** Par récurrence nous montrons que

$$\hat{\mathbf{X}}_t = \hat{\beta}^t \mathbf{X}_0 + \sum_{k=1}^t \hat{\beta}^{t-k} \epsilon_k$$

Son espérance est donc

$$\mathbb{E}(\hat{\mathbf{X}}_t | \mathbf{x}_0^i) = \hat{\beta}^t \mathbf{x}_0^i$$

De même, pour la **Méthode 2 :**

$$\hat{\mathbf{X}}_{t+1} = \hat{\beta}_{t+1} \hat{\mathbf{X}}_t + \mathbf{m}_{t+1} + \hat{\epsilon}_{t+1}$$

où

$$\begin{aligned} \text{--- } \hat{\beta}_{t+1} \text{ est la matrice telle que } (\hat{\beta}_{t+1})_{ij} &= \begin{cases} 0 & \text{si } i \notin \bigcup_{t'=0}^{t+1} \text{Enf}^{t'}(y) \\ \hat{\beta}_{ij} & \text{sinon} \end{cases} \\ \text{--- } \mathbf{m}_{t+1} \text{ est le vecteur tel que } (m_{t+1})_i &= \begin{cases} x_{i,t+1} & \text{si } i \notin \bigcup_{t'=0}^{t+1} \text{Enf}^{t'}(y) \\ 0 & \text{sinon} \end{cases} \\ \text{--- } \hat{\epsilon}_{t+1} \text{ est le vecteur tel que } (\hat{\epsilon}_{t+1})_i &= \begin{cases} 0 & \text{si } i \notin \bigcup_{t'=0}^{t+1} \text{Enf}^{t'}(y) \\ \hat{\sigma}_i & \text{sinon} \end{cases} \end{aligned}$$

L'ensemble des descendants du gène  $y$ ,  $\bigcup_{t'=0}^{t+1} \text{Enf}^{t'}(y)$ , peut être connu au préalable, ou inféré en même temps que le graphe et  $\hat{\beta}$ . Nous obtenons à nouveau par récurrence une formule



similaire à celle du Méthode 1 :

$$\hat{\mathbf{X}}_t = \prod_{k=1}^t \hat{\beta}_k \mathbf{X}_0 + \sum_{k=1}^t \left( \prod_{i=0}^{t-(k+1)} \hat{\beta}_{t-i} \right) (\epsilon_k + \mathbf{m}_k)$$

En prenant l'espérance :

$$\mathbb{E}(\hat{\mathbf{X}}_t | \mathbf{x}_0^n, \mathbf{x}_{NE}^n) = \prod_{k=1}^t \hat{\beta}_k \mathbf{x}_0^n + \sum_{k=1}^t \left( \prod_{i=1}^{t-(k+1)} \hat{\beta}_{t-i} \right) \mathbf{m}_k$$

### 3.3 Résultats

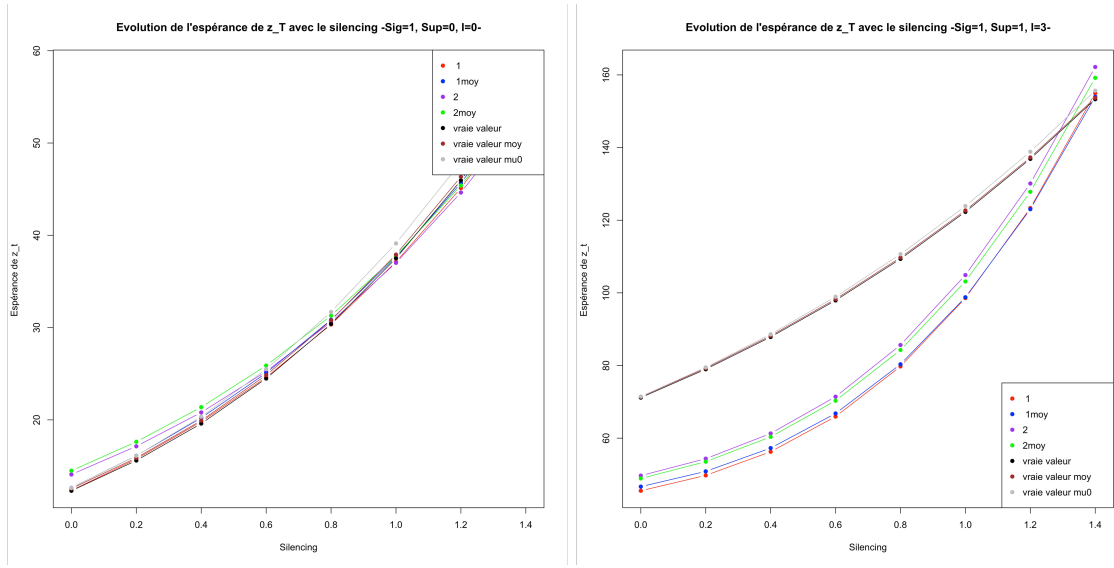


FIGURE 3 – Deux exemples d'évolution du niveau d'expression d'un gène en fonction du paramètre de *silencing*  $\alpha$ , avec  $N = 4$ ,  $p = 10$  et  $T = 5$ . Les deux méthodes présentées dans la section précédente ainsi que leur homologue moyenné sont représentés sur le graphique. Les valeurs théoriques attendues sont désignées sous l'appellation "vraie valeur". Les paramètres choisis pour le modèle linéaire sont indiqués dans le titre du graphique. A gauche, le modèle est bien spécifié, à droite, nous avons rajouté une constante pour qu'il ne le soit pas.

Les espérances ci-dessus sont comparées aux espérances théoriques calculées avec le vrai  $\beta$  dans différentes configurations du modèle linéaire pour analyser le comportement des deux méthodes. Nous pouvons par exemple inférer un modèle mal spécifié, notamment par l'ajout d'une constante, ou bien réaliser l'inférence des paramètres du modèle en connaissant ou non au préalable la structure du graphe. Les espérances sont également calculées et comparées lorsqu'on silencie le gène  $y$ . Pour une même configuration de modèle et de paramètres, 300 estimations de  $\hat{\beta}$  sont réalisées, sur 300 jeux de données de taille  $N = 4$  différents, de façon à éliminer l'influence de l'aléa des jeux de données sur l'estimation de  $\hat{\beta}$ . Nous

obtenons à la fois les valeurs théoriques des espérances que nous pouvons comparer, ainsi qu’une représentation graphique des variations du niveau d’expression du gène  $z$  en fonction de différents paramètres. Un exemple de résultat est donné Figure 3. On peut notamment observer l’importance des différents paramètres sur la performance de chacune des méthodes. A gauche, le modèle linéaire à estimer est bien spécifié, et aucune information à propos de la structure du graphe n’a été donnée. On remarque alors qu’il ne semble pas y avoir de différence flagrante entre les valeurs proposées par les 2 méthodes, ce qui s’explique certainement par la très bonne estimation des paramètres du modèle dans ce cas là. Cependant, à droite, nous avons considéré un modèle mal spécifié en modifiant (4) pour construire les données par  $X_{j,t}^n = \sum_{k \in Pa^1(j)} \beta_k^j X_{k,t-1}^n + \beta_0^j + \epsilon_t^j$ , supposé la structure du graphe connu et forcé le paramètre  $\beta_0^j$  à être nul lors de l’estimation. Ainsi, le modèle estimé sera nécessairement mal spécifié. Nous observons alors une meilleure estimation de l’espérance de  $z_T$  avec la méthode 2 utilisant les données d’expression de gènes : réutiliser ces données a permis de réduire l’emploi des paramètres estimés  $\hat{\beta}$  qui sont mal estimés.

## 4 Perspectives

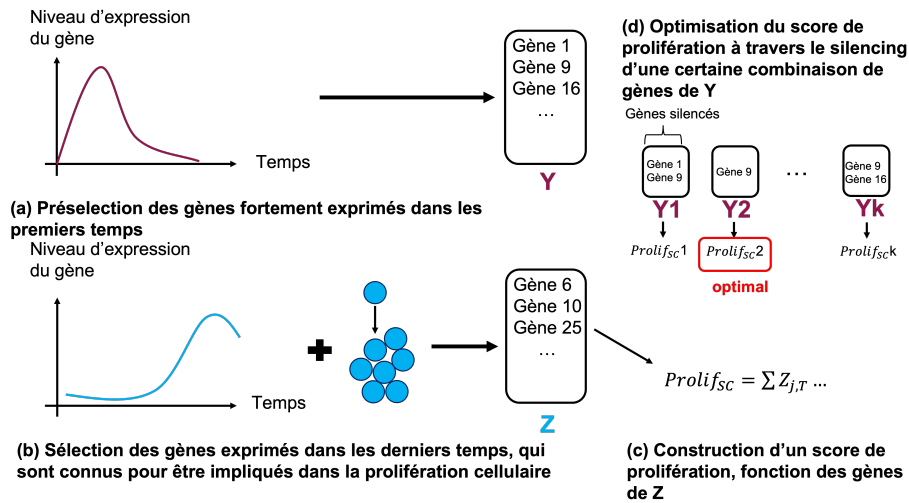


FIGURE 4 – Résumé des différentes étapes pour mesurer l’impact du *silencing* de gènes sur la prolifération cellulaire.

Afin d’aller plus loin dans la mise en concurrence des 2 méthodes proposées, nous souhaitons utiliser de nouveaux jeux de données simulés avec un modèle plus élaboré. De tels jeux de données existent déjà, l’un d’entre eux a notamment été créé par *Marbach et al* (2009) pour le challenge DREAM4, qui proposait aux participants de reconstruire des réseaux d’interactions de gènes. Ce jeu de données a également été utilisé afin de valider et de comparer différents méthodes d’inférence de réseaux, comme DynGENIE3 *Huynh-Thu and Geurts* (2018). L’utilisation de ces données nous serait utile pour déterminer le comportement de nos 2 méthodes lorsque le modèle sous-jacent du réseau et le modèle estimé sont complètement

différents. Autrement dit, nous souhaitons observer l'influence d'un modèle mal spécifié sur nos 2 méthodes. Nous espérons que la réutilisation des données dans la deuxième méthode lui fournira un avantage pour contrer la mauvaise estimation du modèle.

La suite du travail de modélisation du *silencing* consistera à travailler sur des jeux de données réels provenant du CHU de Strasbourg, *Schleiss and Vallat (2021)*, et à déterminer les gènes à silencer ainsi que les gènes cibles à partir d'un score de prolifération préalablement construit.

## Bibliographie

Sanguinetti, G. , Huynh-Thu, V. *et al* (2019), Gene Regulatory Networks : Methods and Protocols, *Springer New York*.

Marbach, D., Schaffter, T., Mattiussi, C., Floreano, D. (2009). Generating realistic in silico gene networks for performance assessment of reverse engineering methods, *Journal of computational biology*.

Schleiss, C. , Vallat, L. *et al* (2021), Temporal multiomic modeling reveals a B-cell receptor proliferative program in chronic lymphocytic leukemia, *Leukemia*.

Huynh-Thu, V., Geurts, P. (2018) dynGENIE3 : dynamical GENIE3 for the inference of gene networks from time series expression data, *Scientific Reports*

Vallat, L., Kemper, CA., Jung, N., Maumy-Bertrand, M., Bertrand, F., Meyer, N., Pocheville, A., Fisher, JW. 3rd, Gribben JG, Bahram S. (2013), Reverse-engineering the genetic circuitry of a cancer cell with predicted intervention in chronic lymphocytic leukemia.