

Model inference from protein time-course in Hematopoietic Stem Cells (HSC)

Pandu Raharja Rene Schoeffel

August 22, 2016

Abstract

Stochasticity of gene expression becomes apparent when studying the dynamics of single cells rather than a population of cells. These fluctuations in gene expression in fact reveal more information about the underlying mechanisms of transcription and translation than one could obtain from population averages.

In this paper we developed a particle filtering algorithm to infer the parameters of stochastic models from single cell time series data. In particular, we apply this method to time-lapse microscopy data of two transcription factors (*Pu.1* and *Gata.1*) in murine blood stem cells. These transcription factors are thought to play a major role in stem cell differentiation.

Our results provide several insights into the dynamics of blood stem cells maturation and specifically in the single cell environment, we managed to gain several valuable insights on the interaction dynamics between two transcription factors and the outcome of cell maturation.

1 Introduction

2 Experimental Setting

The stem cells were put in a culture for one week. Two videos of the cells, taken under bright-field and ultraviolet microscopy, would then be taken for the whole duration of experiment. An automatic detection program tracks down each cell uniquely within the video frame. The program also tracks down cell division and would track the newly created cells likewise. It also detects cell death.

Once every 20-30 minutes a protein expression levels of each cell would then be measured. The expressed protein of interests (in our case *Pu.1* and *Gata.1*) were

tagged with distinct fluorescence proteins, which would allow for the expression level measurement at each time point.

3 Models

3.1 Particle Filtering

3.1.1 Particle

We use particle filtering in our simulation. Particle filtering is a family of methods that use utilize the concept of *particle*. A particle \mathcal{K} is defined as a triple of trajectory X , parameter set θ and assumed model \mathcal{M} ,

$$\mathcal{K} := (X, \theta, \mathcal{M})$$

Specifically in our case, trajectory refers to the concentration of measured proteins Pu.1 and GATA1 while parameter set encompasses all variables that were involved in the reaction such as propensity and its corresponding auxiliary variables (see subsection **Reaction Models**). Note that a trajectory is different from the experimental data \mathcal{D} . A trajectory is the result of the simulation done by applying the parameters onto our model. In functional notation we would assume the i -th value of the trajectory X_i to be a function of the i -th value of the trajectory, the model and the parameters,

$$X_i = f(X_{i-1}, \mathcal{M}, \theta_{i-1})$$

3.1.2 Posterior

Our posterior describes the probability of having the trajectory X and parameter θ given the observation \mathcal{D} ,

$$P(X, \theta | \mathcal{D})$$

We could understand this as the probability of having our simulation return a given set of values *and* having the parameter set θ given that we previously observed the experimental data \mathcal{D} . Using Bayes' Theorem we could further expand our posterior into an update rule,

$$P(X, \theta | \mathcal{D}) = \frac{P(\mathcal{D} | X, \theta) P(X, \theta)}{P(\mathcal{D})}$$

In our simulation we know that, to compute prior $P(\mathcal{D}|X, \theta)$, only knowledge about the trajectory of the simulation is needed. Hence, we could simplify our prior,

$$P(\mathcal{D}|X, \theta) = P(\mathcal{D}|X)$$

Not that the above equation inherently assumes that \mathcal{D} is only directly dependent on X and X is in turn only directly dependent on θ . Incorporating this onto our update rule, and expanding the definition of $P(X, \theta)$ using chain rule, we get

$$P(X, \theta|\mathcal{D}) = \frac{P(\mathcal{D}|X)P(X|\theta)\pi(\theta)}{P(\mathcal{D})}$$

One of the interesting aspects of our is the fact that the i -th simulation result is only dependent on previous simulation, a property known as *Markov property*. We could thus rewrite the update rule as follow,

$$P(X, \theta|\mathcal{D}) = \frac{\prod_{i=0}^N P(\mathcal{D}_i|X_i) \cdot P(X_0) \cdot \prod_{l=1}^N P(X_{[t-1, ti]}|X_i, \theta) \cdot \pi(\theta)}{P(\mathcal{D})}$$

3.1.3 Parameters

During the simulation we assumes certain parameters that would influence the trajectory of the simulation. The parameter set θ could then mathematically be understood as P-tuple containing all the parameters that are assumed in the simulation,

$$P := (P_0, P_1, \dots, P_P)$$

3.1.4 Prior

Our prior $\pi(\theta)$ could be expanded by assuming the independence of each parameter θ_i within the parameter set θ ,

$$\pi(\theta) = \prod_{i=1}^P \pi(\theta_i)$$

Specifically for this simulation, we assume our prior to be Gamma distributed with the parameters α and β ,

$$\pi(\theta) = \prod_{i=1}^P Ga(\theta_i, \alpha_i, \beta_i)$$

We used Gamma distribution as our prior due to the fact that a conjugate prior of a Gamma distributed random variable is also Gamma distributed.

3.2 Simulation Process

The simulation is run in roughly following high level steps:

1. Initialization of parameters θ .
2. Input of data \mathcal{D} .
3. Particle filtering routine:
 - (a) Generation of initial particles for step i

$$Ki := (K_{i1}, K_{i2}, \dots, K_{im})$$

- (b) Simulation run of each particle K_{ij}
- (c) Weighting of each particle. The weight is a function of probability of observing the data given the simulation result.

$$w_i^k = P(D_i|X_i^k) = \mathcal{N}(D_i|X_i^k)$$

- (d) Parameter update for every K,

$$\theta^k \propto P(\theta|X_{[to,ti]}^k)$$

4. Model comparison. This could be done by measure such as AIC, BIC, Bayes Factor etc.

3.2.1 Particle Generation

As mentioned before, A particle \mathcal{K} is defined as a 3-tuple of trajectory X , parameter set θ and assumed model \mathcal{M} . For a specific model, j-th particle at i-th time particle is then just a tuple of trajectory and parameter set $k_{ij} = (x_i, \theta_{ij})$.

For initial time, a particle could be constructed by combining initial parameters, which were arbitrarily defined by a pre-defined prior, and initial trajectory value X_0 . There are three ways to define the initial trajectory value. First is to take 0 as initial value. Second is to take the first measurement data \mathcal{D}_0 . Third is to model initial data as normal distributed instances around \mathcal{D}_0 , i.e. $X_0 \propto \mathcal{N}(\mathcal{D}_0)$. It is worth noting that the first strategy may not always be applicable.

M particles would then be created and for each iteration of the simulation with each particle having weight which calculated from previous iteration w_{i-1}^k . In the

next iteration, M particles are to be chosen. Note that this implies that a particle from previous iteration could be chosen more than once. Corollarily, a particle from previous iteration may also not be chosen at all. The draw is done by first drawing a random $z \propto \mathcal{U}(0, 1)$. A particle k_{ij} is chosen with j denoting the smallest index so that the sum of all weight of particles with index smaller than j is larger than z , i.e.

$$\sum_{l=1}^j w_{i-1}^l \geq z$$

3.2.2 Simulation Run

Simulation is done using Gillespie’s SSA algorithm. There are possible improvements that we thought might be useful in our project such as τ -leaping and Bayesian Approximation Method. Such improvements are generally essential since it would reduce computational cost of the simulation and in turn enable the simulation to scale further to accommodate more leveraged simulation – in our case we could simulate more cells at the same time in possible more fine-grained time-resolved manner. In our case especially, some kind of improvement would be needed to accommodate the explosion of the number of cells within our culture. Without some kind of approximative method there would be a bottleneck in simulation caused by the exponential increase of the number of living cells that have to be simulated.

Tau Leaping

Generally, tau-leaping works by performing all reactions happening for an interval of length tau before updating the parameters and propensity function. By allowing this kind of approximation we potentially allow more efficient simulation and thus increase the capability of the simulation to cope with larger systems. In our case, we use following tau,

$$\tau = \frac{1}{a_0} \ln \Gamma$$

with a_0 referring to total sum of reaction propensities and Γ being uniformly distributed between 0 and 1.

3.2.3 Particle Weighting

Upon the completion of all simulations within one iteration, the quality of each simulation (and in turn, the particle) will be quantified. This quantification is implied in the weighting of the particle at the next iteration w_i^k . We assume the Gaussian distri-

bution of particle around the measurement time at time $t = i$. Using this assumption, we can quantify our weight in following manner,

$$w_i^k = P(\mathcal{D}_i | X_i^k) = \mathcal{N}(\mathcal{D}_i | X_i^k)$$

3.2.4 Parameters Update and Gamma Distribution

After each run, the parameters would then be updated using Gamma Distribution. For k -th particle, the updated parameters at time i are dependent on previous parameters given the trajectory of data in $[t_0, t_i]$,

$$\Theta^k \propto P(\Theta | X_{[t_0, t_i]})$$

Applying Bayes' Rule on the probability and independent assumption, we get following equation,

$$P(\Theta | X) = \frac{P(X | \Theta)}{P(X)} \prod_i P(\Theta_i)$$

Gamma Distribution is particularly favored for the update since a conjugate prior of Gamma Distribution is in turn Gamma distributed,

$$P(\Theta | X) = \frac{P(X | \Theta)}{P(X)} \prod_{d=1}^P Ga(\Theta_d, \alpha_d, \beta_d)$$

With α and β referring to both the slope and skewness of the distribution respectively. Owing to the fact that conjugate prior of Gamma Distribution is also Gamma distributed, this could be then further simplified as an update rule of the Gamma distribution,

$$P(\Theta | X) = \frac{P(X | \Theta)}{P(X)} \prod_{d=1}^P Ga(\Theta_d, \alpha_d + r_d, \beta_d + G_d)$$

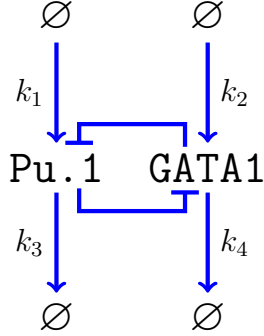
r_d here refers to the number of times the d -th reaction was fired during the run and G_d is defined as follows,

$$G_d := \frac{1}{k_d} \sum_{l=0}^i a_l(X(S))$$

Here, k_d and $a_l(X(S))$ refer to reaction rate of the d -th reaction and trajectory-dependent propensity function, respectively.

3.3 Reaction Models

We assume following putative cross inhibition between genes Pu.1 and GATA1 playing roles in cell differentiation dynamics:



There are four possible reactions happening within our model: k_1 , k_2 , k_3 and k_4 . In our simulation, each reaction was assigned propensity value which roughly is proportional to the likelihood of the reaction happening in a given time:

$$P(R_i = k_l) \propto a_l$$

The propensity a of decay reactions k_3 and k_4 roughly follows the mass action law of chemical reaction:

$$a_l = k_l \cdot A; l \in [3, 4]$$

For production reactions k_1 and k_2 we assume Michaelis-Menten inhibition kinetics that influences the production of the species. In this assumption, the inhibiting characteristics of a species on another species negatively influences the creation rate of the species it inhibits. The propensity of the reactions are thus,

$$a_1 = k_1 \cdot \frac{\text{GATA1}^{-n}}{\text{GATA1}^{-n} + K_{P1}^{-n}}$$

$$a_2 = k_2 \cdot \frac{\text{Pu.1}^{-n}}{\text{Pu.1}^{-n} + K_{P2}^{-n}}$$

4 BlaBla