

# Identification of cell signaling pathways based on biochemical reaction kinetics repositories

---

Student: Gustavo Estrela

Advisor: Marcelo da Silva Reis (Butantan Institute)

May 2019

Instituto de Matemática e Estatística  
Centro de Toxinas, Resposta-imune e Sinalização Celular (CeTICS)  
Laboratório Especial de Ciclo Celular, Instituto Butantan  
This project receives funding from FAPESP

# Introduction

---

# Cell Signaling

Cell signaling allows cells to respond to signals that come from its environment changing its behaviour accordingly.

# Cell Signaling

Cell signaling allows cells to respond to signals that come from its environment changing its behaviour accordingly.

This mechanism is essential for many cell functions, including reproduction, growth and death.

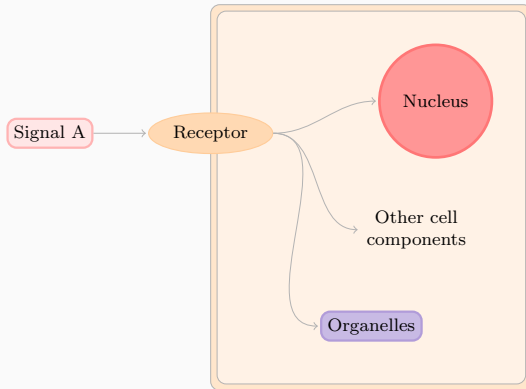
# Cell Signaling

Cell signaling allows cells to respond to signals that come from its environment changing its behaviour accordingly.

This mechanism is essential for many cell functions, including reproduction, growth and death.

Understanding the functioning of cell signaling is important in many biological areas.

# Cell Signaling



A signal propagates in an organism through chemical reactions that are caused by the change of concentration of chemical species.

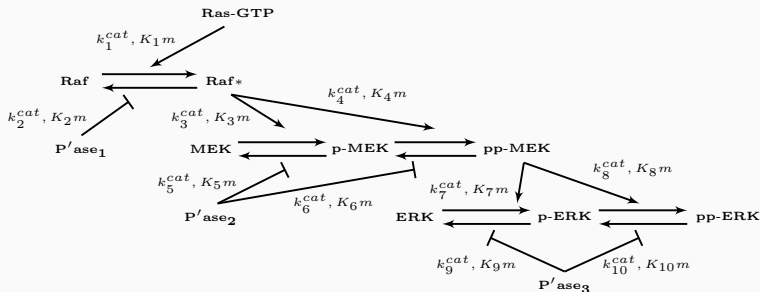
A signal propagates in an organism through chemical reactions that are caused by the change of concentration of chemical species.

We call the path of a signal a **cell signaling pathway**.



# Cell Signaling Pathways

A cell signaling network can be characterized by a sequence of chemical reactions



# Mathematical Models of Signaling Networks

We can summarize the state of the cell with measurements based on the concentration of some chemical species.

# Mathematical Models of Signaling Networks

We can summarize the state of the cell with measurements based on the concentration of some chemical species.

Using biochemical and enzymatic kinetics, we can model the concentration change of chemical species over time of a pathway.

# Identification of Cell Signaling Pathways

The problem of identification of cell signaling pathways is the problem of finding the components of a signaling pathway and how they interact given a set of experimental measurement.

# Identification of Cell Signaling Pathways

The problem of identification of cell signaling pathways is the problem of finding the components of a signaling pathway and how they interact given a set of experimental measurement.

As the input, a description of a biological experiment and a set of experimental measurements are given.

# Identification of Cell Signaling Pathways

The problem of identification of cell signaling pathways is the problem of finding the components of a signaling pathway and how they interact given a set of experimental measurement.

As the input, a description of a biological experiment and a set of experimental measurements are given. A possible output to the problem is composed by:

# Identification of Cell Signaling Pathways

The problem of identification of cell signaling pathways is the problem of finding the components of a signaling pathway and how they interact given a set of experimental measurement.

As the input, a description of a biological experiment and a set of experimental measurements are given. A possible output to the problem is composed by:

- a model composed by a set of chemical reactions that are relevant for the biological experiment;

# Identification of Cell Signaling Pathways

The problem of identification of cell signaling pathways is the problem of finding the components of a signaling pathway and how they interact given a set of experimental measurement.

As the input, a description of a biological experiment and a set of experimental measurements are given. A possible output to the problem is composed by:

- a model composed by a set of chemical reactions that are relevant for the biological experiment;
- information about the reaction rate constants of the model.



## Identification of Cell Signaling Pathways

One can search for the set of chemical reactions relevant for a biological experiment in repositories like the Kyoto Encyclopedia of Genes and Genomes (KEGG).

## Identification of Cell Signaling Pathways

One can search for the set of chemical reactions relevant for a biological experiment in repositories like the Kyoto Encyclopedia of Genes and Genomes (KEGG). However, the pathway maps from KEGG may be incomplete or have impertinent reactions for the biological experiment of interest.

# Identification of Cell Signaling Pathways

One can search for the set of chemical reactions relevant for a biological experiment in repositories like the Kyoto Encyclopedia of Genes and Genomes (KEGG). However, the pathway maps from KEGG may be incomplete or have impertinent reactions for the biological experiment of interest.

Hence, it is desirable to construct a method that can systematically modify these models and choose the one that better represents the experiment.

# Identification of Cell Signaling Pathways

Lulu Wu (2015) presented in her master dissertation a methodology that proposes to systematically modify models of signaling network in order to better represent experiments.

# Identification of Cell Signaling Pathways

Lulu Wu (2015) presented in her master dissertation a methodology that proposes to systematically modify models of signaling network in order to better represent experiments.

On her work, the problem of identification of cell signaling pathways is treated as a feature selection problem.

# Feature Selection for Identification of Signaling Pathways

The methodology proposed by Wu defines the set of features as a set of chemical reactions that can be added to a starting model.

# Feature Selection for Identification of Signaling Pathways

The methodology proposed by Wu defines the set of features as a set of chemical reactions that can be added to a starting model. This set of chemical reactions is fetched from KEGG and stored in a database of interactions.

## Wu's Search Algorithm for Feature Selection

The search algorithm used by Wu is the Sequential Forward Selection (SFS).



## Wu's Cost Function for Feature Selection

Wu defines the cost function as the minimum distance between experimental and simulated data.

## Wu's Cost Function for Feature Selection

Wu defines the cost function as the minimum distance between experimental and simulated data. The minimum distance is found using a Simmulated Annealing that traverses the parameter space.

## Results of Wu's Methodology

Lulu Wu tested her methodology by trying to recreate models given a cut of the original model.

## Results of Wu's Methodology

Lulu Wu tested her methodology by trying to recreate models given a cut of the original model. However, the methodology worked satisfactorily only when the cut was similar to the original model.

## Difficulties of Wu's Methodology

We can point three aspects of Wu's work that could explain its limitations.

## Difficulties of Wu's Methodology

We can point three aspects of Wu's work that could explain its limitations.

- the database of interactions used could be more nearly complete;

# Difficulties of Wu's Methodology

We can point three aspects of Wu's work that could explain its limitations.

- the database of interactions used could be more nearly complete;
- the search algorithm could also consider removing interactions;

# Difficulties of Wu's Methodology

We can point three aspects of Wu's work that could explain its limitations.

- the database of interactions used could be more nearly complete;
- the search algorithm could also consider removing interactions;
- the cost function could implement a proper penalization of models;



## What we Propose on this Project

We propose to create a methodology that uses a feature selection approach for identification of signaling pathways, tackling the difficulties encountered by Wu.

## What we Propose on this Project

To get a more nearly complete database of interactions, we should fetch information from KEGG and other databases,

## What we Propose on this Project

To get a more nearly complete database of interactions, we should fetch information from KEGG and other databases, such as STRING and SABIO-RK.

## What we Propose on this Project

To use new search algorithms,

## What we Propose on this Project

To use new search algorithms, we intend to use more general algorithms that can also remove interactions.

## What we Propose on this Project

To define new cost functions,

## What we Propose on this Project

To define new cost functions, we intend to use Bayesian approaches of model selection that allow us to create estimates of probabilities such as  $p(M|D)$  or  $p(D|M)$ .

## Objectives of this Project



## Objectives of this Project

- Build a database of interactions.

## Objectives of this Project

- Build a database of interactions.
- Create a cost function for models of signaling pathways.

## Objectives of this Project

- Build a database of interactions.
- Create a cost function for models of signaling pathways.
- Formulate systematic modifications to a model as the search space of a feature selection model.

# Objectives of this Project

- Build a database of interactions.
- Create a cost function for models of signaling pathways.
- Formulate systematic modifications to a model as the search space of a feature selection model.
- Test the methodology on known signaling pathways and also on pathways of interest in our lab.

# Fundamental Concepts

---

## Kinetics Modeling of Chemical Reactions

# Mathematical Modeling of Reactions

In this project we use three possible models of kinetics of an interaction:

In this project we use three possible models of kinetics of an interaction:

- first order interaction kinetics;



In this project we use three possible models of kinetics of an interaction:

- first order interaction kinetics;
- second order interaction kinetics;

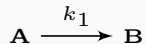
# Mathematical Modeling of Reactions

In this project we use three possible models of kinetics of an interaction:

- first order interaction kinetics;
- second order interaction kinetics;
- Michaelis-Menten enzymatic kinetics.

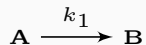
# Kinetic Modeling of First Order Iteration

A first order reaction:



# Kinetic Modeling of First Order Iteration

A first order reaction:

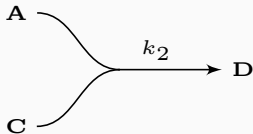


has rate of:

$$k_1[\mathbf{A}].$$

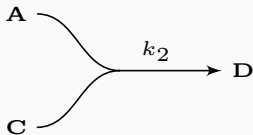
# Kinetic Modeling of Second Order Iteration

A second order reaction:



# Kinetic Modeling of Second Order Iteration

A second order reaction:

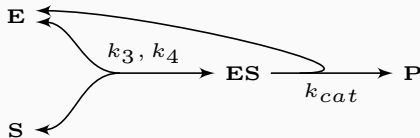


has rate of:

$$k_2[A][C].$$

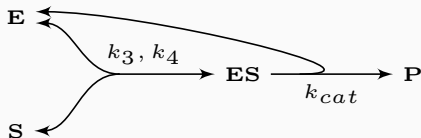
# Kinetic Modeling of Enzymatic Reactions

An enzymatic reaction:



# Kinetic Modeling of Enzymatic Reactions

An enzymatic reaction:

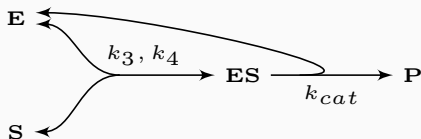


Can be divided in two first order reactions plus a second order reaction.



# Kinetic Modeling of Enzymatic Reactions

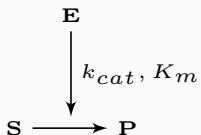
An enzymatic reaction:



Can be divided in two first order reactions plus a second order reaction. However, with the appropriate assumptions, it is possible to use a Michaelis-Menten simplification of this reaction.

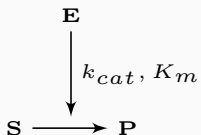
# Michaelis-Menten Kinetics

We denote Michaelis-Menten simplification of the last enzymatic reaction as



# Michaelis-Menten Kinetics

We denote Michaelis-Menten simplification of the last enzymatic reaction as

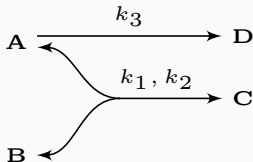


and it has rate of:

$$k_{cat} \frac{[\mathbf{E}][\mathbf{S}]}{K_M + [\mathbf{S}]}.$$

# Kinetics of a System of Reactions

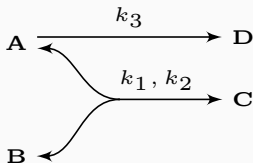
Suppose we want to model the kinetics of A on these reactions:



This system can be divided in three reactions:

# Kinetics of a System of Reactions

Suppose we want to model the kinetics of A on these reactions:

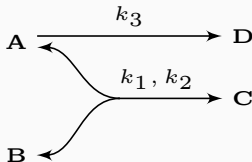


This system can be divided in three reactions:

- $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ ,

# Kinetics of a System of Reactions

Suppose we want to model the kinetics of A on these reactions:

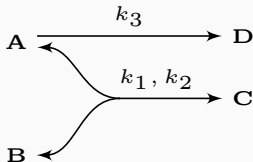


This system can be divided in three reactions:

- $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ ,
- $C \longrightarrow A + B$ , with rate  $k_2[C]$ ,

# Kinetics of a System of Reactions

Suppose we want to model the kinetics of A on these reactions:



This system can be divided in three reactions:

- $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ ,
- $C \longrightarrow A + B$ , with rate  $k_2[C]$ ,
- $A \longrightarrow D$ , with rate  $k_3[A]$ .

# Kinetics of a System of Reactions

In  $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ , A is a reactant.



# Kinetics of a System of Reactions

In  $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ , A is a reactant.

In  $C \longrightarrow A + B$ , with rate  $k_2[C]$ , A is a product.

# Kinetics of a System of Reactions

In  $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ , A is a reactant.

In  $C \longrightarrow A + B$ , with rate  $k_2[C]$ , A is a product.

In  $A \longrightarrow D$ , with rate  $k_3[A]$ , A is a reactant.

# Kinetics of a System of Reactions

In  $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ , A is a reactant.

In  $C \longrightarrow A + B$ , with rate  $k_2[C]$ , A is a product.

In  $A \longrightarrow D$ , with rate  $k_3[A]$ , A is a reactant.

Then, the differential equation that models the concentration change of A is:

# Kinetics of a System of Reactions

In  $A + B \longrightarrow C$ , with rate  $k_1[A][B]$ , A is a reactant.

In  $C \longrightarrow A + B$ , with rate  $k_2[C]$ , A is a product.

In  $A \longrightarrow D$ , with rate  $k_3[A]$ , A is a reactant.

Then, the differential equation that models the concentration change of A is:

$$\frac{d[A]}{dt} = -k_1[A][B] + k_2[C] - k_3[A].$$

## Bayesian Methods for Biochemical Model Selection

# State of the Art Methods for Model Selection

There are two main Bayesian methods available for biochemical model selection:

# State of the Art Methods for Model Selection

There are two main Bayesian methods available for biochemical model selection:

- Approximate Bayesian Computation;

# State of the Art Methods for Model Selection

There are two main Bayesian methods available for biochemical model selection:

- Approximate Bayesian Computation;
- Marginal likelihood estimation through Thermodynamic Integration.



# State of the Art Methods for Model Selection

There are two main Bayesian methods available for biochemical model selection:

- Approximate Bayesian Computation;
- Marginal likelihood estimation through Thermodynamic Integration.

For both methods, we resort to Metropolis-Hastings algorithm to generate samples of distributions.

# Metropolis-Hastings algorithm

With Metropolis-Hastings, we can generate a sample of a distribution  $p(\lambda)$  doing the following:

# Metropolis-Hastings algorithm

With Metropolis-Hastings, we can generate a sample of a distribution  $p(\lambda)$  doing the following:

# Metropolis-Hastings algorithm

With Metropolis-Hastings, we can generate a sample of a distribution  $p(\lambda)$  doing the following:

1. Choose some  $\lambda_0$  for which  $p(\lambda_0) > 0$ , and set  $t = 1$ ;

# Metropolis-Hastings algorithm

With Metropolis-Hastings, we can generate a sample of a distribution  $p(\lambda)$  doing the following:

1. Choose some  $\lambda_0$  for which  $p(\lambda_0) > 0$ , and set  $t = 1$ ;
2. Sample a candidate point  $\lambda^*$  from a jump distribution,  $J(\lambda|\lambda_{t-1})$ ;

# Metropolis-Hastings algorithm

With Metropolis-Hastings, we can generate a sample of a distribution  $p(\lambda)$  doing the following:

1. Choose some  $\lambda_0$  for which  $p(\lambda_0) > 0$ , and set  $t = 1$ ;
2. Sample a candidate point  $\lambda^*$  from a jump distribution,  $J(\lambda|\lambda_{t-1})$ ;
3. Calculate the ratio  $r = \frac{p(\lambda^*)J_t(\lambda^{t-1}|\lambda^*)}{p(\lambda^{t-1})J_t(\lambda^*|\lambda^{t-1})}$ ;

# Metropolis-Hastings algorithm

With Metropolis-Hastings, we can generate a sample of a distribution  $p(\lambda)$  doing the following:

1. Choose some  $\lambda_0$  for which  $p(\lambda_0) > 0$ , and set  $t = 1$ ;
2. Sample a candidate point  $\lambda^*$  from a jump distribution,  $J(\lambda|\lambda_{t-1})$ ;
3. Calculate the ratio  $r = \frac{p(\lambda^*)J_t(\lambda^{t-1}|\lambda^*)}{p(\lambda^{t-1})J_t(\lambda^*|\lambda^{t-1})}$ ;
4. Set  $\lambda_t = \lambda^*$  with probability  $\min(1, r)$  and  $\lambda_t = \lambda_{t-1}$  otherwise;

# Metropolis-Hastings algorithm

With Metropolis-Hastings, we can generate a sample of a distribution  $p(\lambda)$  doing the following:

1. Choose some  $\lambda_0$  for which  $p(\lambda_0) > 0$ , and set  $t = 1$ ;
2. Sample a candidate point  $\lambda^*$  from a jump distribution,  $J(\lambda|\lambda_{t-1})$ ;
3. Calculate the ratio  $r = \frac{p(\lambda^*)J_t(\lambda^{t-1}|\lambda^*)}{p(\lambda^{t-1})J_t(\lambda^*|\lambda^{t-1})}$ ;
4. Set  $\lambda_t = \lambda^*$  with probability  $\min(1, r)$  and  $\lambda_t = \lambda_{t-1}$  otherwise;
5. Increase  $t$  by one and repeat from Step 2 if not reached iteration limit.



# Model Selection

---

Ranking with Marginal Likelihood Estimation

## Likelihood of Data Given Model and Parameters

If we consider that a model  $M$  with parameters  $\theta$  correctly represent the signaling pathway

## Likelihood of Data Given Model and Parameters

If we consider that a model  $M$  with parameters  $\theta$  correctly represent the signaling pathway and that there is a Gaussian observation error on  $D$ .

## Likelihood of Data Given Model and Parameters

If we consider that a model  $M$  with parameters  $\theta$  correctly represent the signaling pathway and that there is a Gaussian observation error on  $D$ . Then, the likelihood of observing experimental data  $D$  is:

$$p(D|M, \theta) =$$

## Likelihood of Data Given Model and Parameters

If we consider that a model  $M$  with parameters  $\theta$  correctly represent the signaling pathway and that there is a Gaussian observation error on  $D$ . Then, the likelihood of observing experimental data  $D$  is:

$$p(D|M, \theta) = p_{\mathcal{N}(\bar{0}, \Sigma)}(\phi(M, \theta) - D).$$

Where  $\phi(M, \theta)$  is the simulated observation.

# Marginal Likelihood of Data

We can marginalize the likelihood to obtain:

$$p(D|M) = \int_{\Theta} p(D|M, \theta) p(\theta|M) d\theta$$

# Marginal Likelihood of Data

We can marginalize the likelihood to obtain:

$$p(D|M) = \int_{\Theta} p(D|M, \theta) p(\theta|M) d\theta$$

Calculating this integral is hard, therefore we resort to estimating another integral.



## Power-posterior distributions

We define a power-posterior distribution as:

$$p_{\beta}(\theta) = \frac{p(D|\theta, M)^{\beta} p(\theta|M)}{\int_{\Theta} p(D|\theta, M)^{\beta} p(\theta|M) d\theta},$$

# Power-posterior distributions

We define a power-posterior distribution as:

$$p_{\beta}(\theta) = \frac{p(D|\theta, M)^{\beta} p(\theta|M)}{\int_{\Theta} p(D|\theta, M)^{\beta} p(\theta|M) d\theta},$$

Note that:

$$p_{\beta=0}(\theta) = p(\theta|M),$$

## Power-posterior distributions

We define a power-posterior distribution as:

$$p_{\beta}(\theta) = \frac{p(D|\theta, M)^{\beta} p(\theta|M)}{\int_{\Theta} p(D|\theta, M)^{\beta} p(\theta|M) d\theta},$$

Note that:

$$p_{\beta=0}(\theta) = p(\theta|M),$$

and that

$$p_{\beta=1}(\theta) = \frac{p(D, \theta|M)}{\int_{\Theta} p(D, \theta|M) d\theta} = \frac{p(\theta|D, M)p(D|M)}{p(D|M)} = p(\theta|D, M).$$

# The Thermodynamic Integral

Using power-posteriors distributions, it is possible to show that

# The Thermodynamic Integral

Using power-posteriors distributions, it is possible to show that

$$\ln p(D|M) = \int_0^1 \mathbb{E}_{p_\beta(\theta)} [\ln p(D|\theta, M)] d\beta.$$

## Estimating the Thermodynamic Integral

It is possible to estimate the Thermodynamic Integral using the trapezoidal rule.

# Estimating the Thermodynamic Integral

It is possible to estimate the Thermodynamic Integral using the trapezoidal rule. Setting  $0 = \beta_0 < \beta_1 < \dots < \beta_T = 1$ , the marginal likelihood is approximately equal to:

# Estimating the Thermodynamic Integral

It is possible to estimate the Thermodynamic Integral using the trapezoidal rule. Setting  $0 = \beta_0 < \beta_1 < \dots < \beta_T = 1$ , the marginal likelihood is approximately equal to:

$$\sum_{t=0}^{T-1} (\beta_{t+1} - \beta_t) \frac{\mathbb{E}_{p_{\beta_{t+1}}(\theta)}[\log p(D|M, \theta)] + \mathbb{E}_{p_{\beta_t}(\theta)}[\log p(D|M, \theta)]}{2}$$



# Estimating the Thermodynamic Integral

To produce the estimates of

$$\mathbb{E}_{p_{\beta_t}(\theta)}[\log p(D|M, \theta)] \text{ for } t \in \{0, \dots, T\}$$

we need to produce samples of the power-posteriors  $p_{\beta_t}(\theta)$ .

## Sampling from Power-posteriors

The sampling of the power-posteriors are generated using Metropolis-Hastings algorithms in three steps.

## Sampling from Power-posteriors

The sampling of the power-posteriors are generated using Metropolis-Hastings algorithms in three steps. In all of the steps, the proposal distribution used is a truncated multivariate normal.

## Sampling from Power-posteriors

On the first step, the jump distribution has a diagonal covariance matrix.

## Sampling from Power-posteriors

On the first step, the jump distribution has a diagonal covariance matrix. This matrix is updated according to the rate of acceptance of parameters.

## Sampling from Power-posteriors

On the first step, the jump distribution has a diagonal covariance matrix. This matrix is updated according to the rate of acceptance of parameters.

- if the acceptance rate is high, then increase the variance of the jump;

## Sampling from Power-posteriors

On the first step, the jump distribution has a diagonal covariance matrix. This matrix is updated according to the rate of acceptance of parameters.

- if the acceptance rate is high, then increase the variance of the jump;
- if the acceptance rate is low, then decrease the variance of the jump.

## Sampling from Power-posteriors

On the first step, the jump distribution has a diagonal covariance matrix. This matrix is updated according to the rate of acceptance of parameters.

- if the acceptance rate is high, then increase the variance of the jump;
- if the acceptance rate is low, then decrease the variance of the jump.

On the second and third step, the covariance matrix of the jump distribution is estimated with the current sample of the posterior.



## Sampling from the Power-posteriors

On the third step a Populational Monte Carlo Markov Chain is performed. This algorithm allows us to mix samples from different temperatures.

## Ranking with Approximate Bayesian Computation

# Approximate Bayesian Computation

Approximate Bayesian Computation (ABC) is a method that allows one to obtain samples of a distribution close to  $p(\theta, M|D)$ .

# Approximate Bayesian Computation

Approximate Bayesian Computation (ABC) is a method that allows one to obtain samples of a distribution close to  $p(\theta, M|D)$ . A general ABC implementation works as follow:

# Approximate Bayesian Computation

Approximate Bayesian Computation (ABC) is a method that allows one to obtain samples of a distribution close to  $p(\theta, M|D)$ . A general ABC implementation works as follow:

1. Sample a parameter candidate  $(\theta^*, M^*)$  from some proposal distribution.

# Approximate Bayesian Computation

Approximate Bayesian Computation (ABC) is a method that allows one to obtain samples of a distribution close to  $p(\theta, M|D)$ . A general ABC implementation works as follow:

1. Sample a parameter candidate  $(\theta^*, M^*)$  from some proposal distribution.
2. Generate simulations  $\phi(M^*, \theta^*) = D^*$ .

# Approximate Bayesian Computation

Approximate Bayesian Computation (ABC) is a method that allows one to obtain samples of a distribution close to  $p(\theta, M|D)$ . A general ABC implementation works as follow:

1. Sample a parameter candidate  $(\theta^*, M^*)$  from some proposal distribution.
2. Generate simulations  $\phi(M^*, \theta^*) = D^*$ .
3. Calculate  $d(D^*, D)$ . If  $d(D^*, D) < \epsilon$  for some previously specified  $\epsilon$ , then add  $(\theta^*, M^*)$  to the sample.

# Approximate Bayesian Computation

Approximate Bayesian Computation (ABC) is a method that allows one to obtain samples of a distribution close to  $p(\theta, M|D)$ . A general ABC implementation works as follow:

1. Sample a parameter candidate  $(\theta^*, M^*)$  from some proposal distribution.
2. Generate simulations  $\phi(M^*, \theta^*) = D^*$ .
3. Calculate  $d(D^*, D)$ . If  $d(D^*, D) < \epsilon$  for some previously specified  $\epsilon$ , then add  $(\theta^*, M^*)$  to the sample.
4. Repeat until some iteration limit.



The result of the algorithm is a sample of the distribution

$$p(\theta, M | d(\phi(M, \theta), D) < \epsilon).$$

ABC sequential Monte Carlo improves a simple ABC algorithm by using a sequence  $\epsilon_0 > \dots > \epsilon_T$  acceptance tolerances.

ABC sequential Monte Carlo improves a simple ABC algorithm by using a sequence  $\epsilon_0 > \dots > \epsilon_T$  acceptance tolerances. The sample for a tolerance  $\epsilon_i$  is used to generate candidates for sample of tolerance  $\epsilon_{i+1}$ .

ABC sequential Monte Carlo improves a simple ABC algorithm by using a sequence  $\epsilon_0 > \dots > \epsilon_T$  acceptance tolerances. The sample for a tolerance  $\epsilon_i$  is used to generate candidates for sample of tolerance  $\epsilon_{i+1}$ .

We can use the accepted parameters of tolerance  $\epsilon$  and model  $M$  to estimate

$$p(M|d(\phi(M, \theta)) < \epsilon).$$

# Experiments on Model Selection

---

## Software used for Model Ranking Experiments

We tested the two methods of model ranking using the software:

# Software used for Model Ranking Experiments

We tested the two methods of model ranking using the software:

- **SigNetMS:** an implementation of the Marginal Likelihood method created in this project.

# Software used for Model Ranking Experiments

We tested the two methods of model ranking using the software:

- **SigNetMS**: an implementation of the Marginal Likelihood method created in this project.
- **ABC-SysBio**: an implementation of ABC-SMC.



## Experiment description

We ran two experiments based on the same procedure:

## Experiment description

We ran two experiments based on the same procedure:

- Create 4 candidate models.

## Experiment description

We ran two experiments based on the same procedure:

- Create 4 candidate models.
- For one of the models, choose a set of parameter values and time steps and simulate data.

## Experiment description

We ran two experiments based on the same procedure:

- Create 4 candidate models.
- For one of the models, choose a set of parameter values and time steps and simulate data.
- Add Gaussian noise to the simulations. Repeat two more times to generate three observations of the system.

## Experiment description

We ran two experiments based on the same procedure:

- Create 4 candidate models.
- For one of the models, choose a set of parameter values and time steps and simulate data.
- Add Gaussian noise to the simulations. Repeat two more times to generate three observations of the system.
- Neglect chosen parameter values and define prior distributions for every parameter.

## Experiment description

We ran two experiments based on the same procedure:

- Create 4 candidate models.
- For one of the models, choose a set of parameter values and time steps and simulate data.
- Add Gaussian noise to the simulations. Repeat two more times to generate three observations of the system.
- Neglect chosen parameter values and define prior distributions for every parameter.
- Rank the four models.

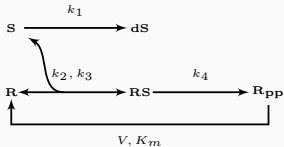
## Experiment #1

# Experiment #1

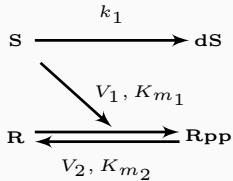
This experiment is originally from Vyshemirsky and Girolami (2007), in which they present results of Annealing Melting Integration.



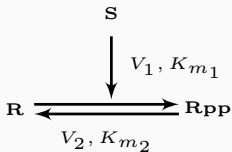
# Experiment #1



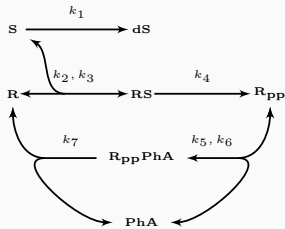
Model 1



Model 2



Model 3



Model 4

## Experiment #1

The model used to create the observations was Model 1.

## Experiment #1

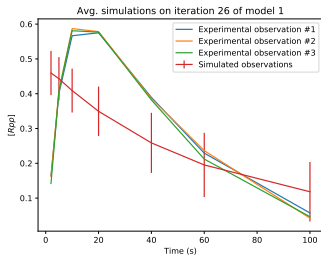
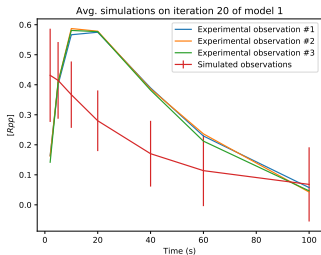
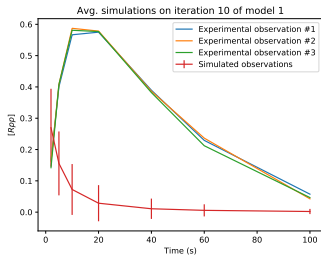
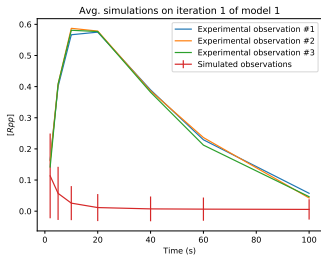
The model used to create the observations was Model 1.

The priors distribution used for all parameters is  $\text{Gamma}(1, 3)$ .

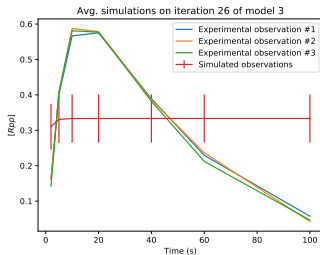
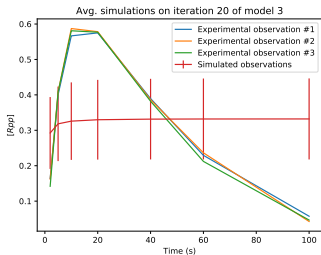
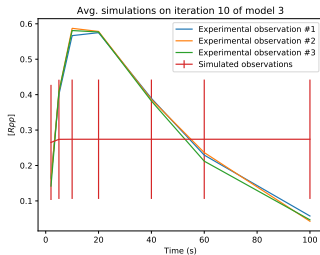
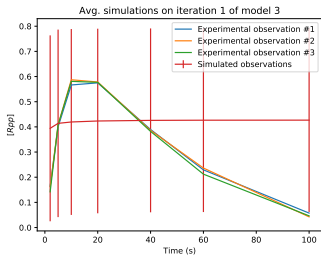
The ABC-SysBio software returned the following ranking of models:

$$3 > 2 > 4 > 1$$

# Results on ABC-SysBio



# Results on ABC-SysBio



The ranking returned by SigNetMS on the first experiment is:

$$1 > 2 > 4 > 3;$$

The ranking returned by SigNetMS on the first experiment is:

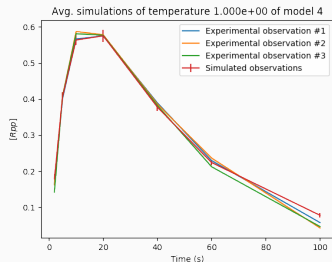
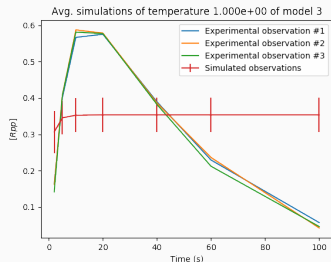
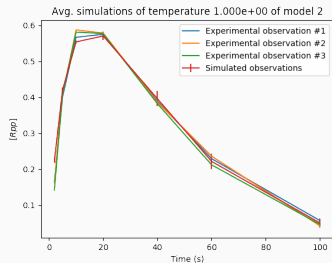
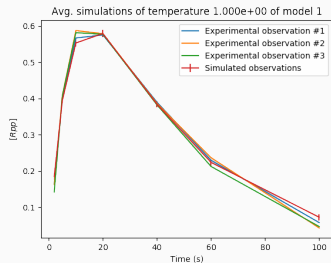
$$1 > 2 > 4 > 3;$$

which is very similar to the ranking presented originally by Vyshemirsky and Girolami (2007):

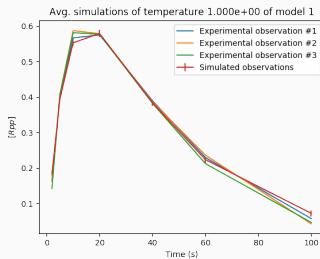
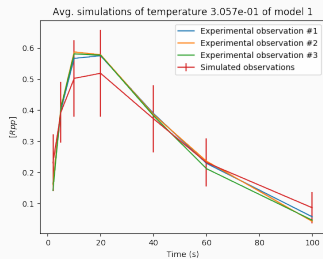
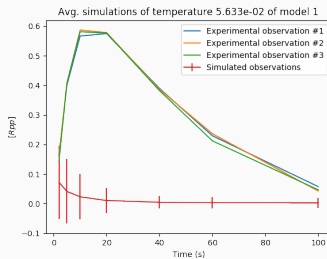
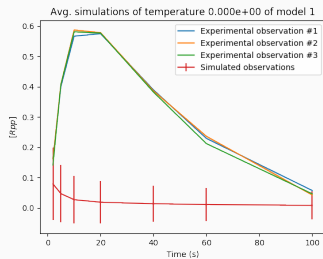
$$1 > 4 > 2 > 3.$$



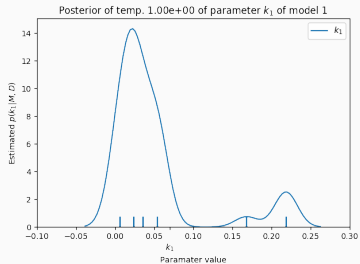
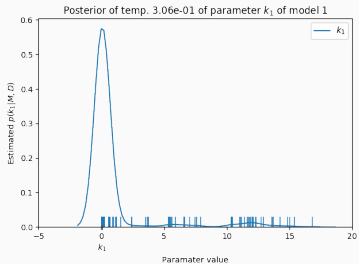
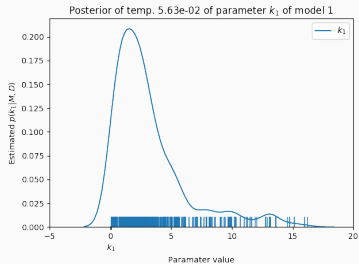
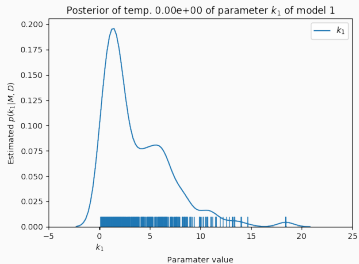
# Results on SigNetMS



# Results on SigNetMS



# Results on SigNetMS

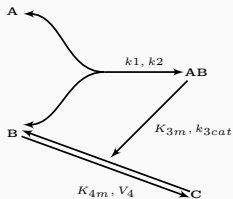


## Experiment #2

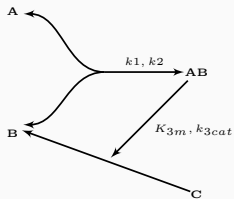
## Experiment #2

This experiment is very similar to the later and it was designed by us.

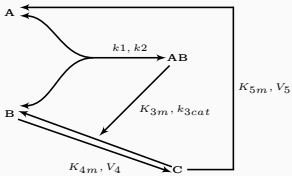
## Experiment #2



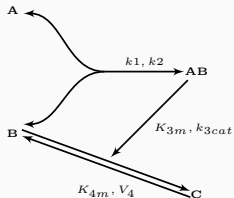
Model 1



Model 2



Model 3



Model 4

## Experiment #2

We used the following prior distributions for model parameters:

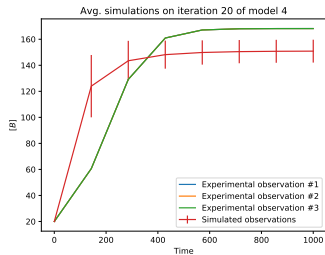
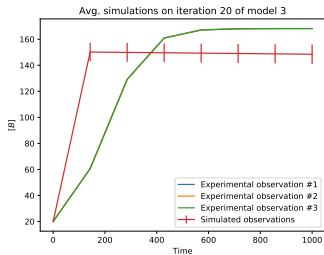
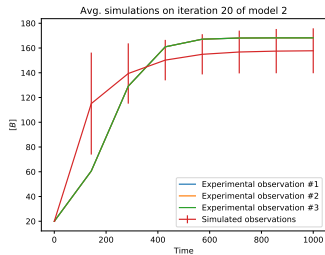
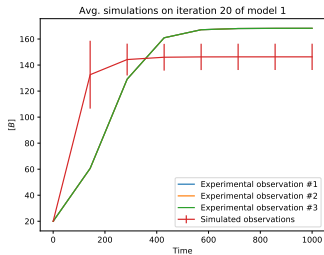
Parameter	Models	Prior
$k_1$	1, 2, 3, 4	$\text{Gamma}(1, 0.01)$
$k_2$	1, 2, 3, 4	$\text{Gamma}(2, 0.5)$
$k_{3cat}$	1, 2, 3, 4	$\text{Gamma}(4, 1)$
$K_{3m}$	1, 2, 3, 4	$\text{Gamma}(2, 1500)$
$V_4$	1, 3, 4	$\text{Gamma}(2, 1)$
$K_{4m}$	1, 3, 4	$\text{Gamma}(2, 100)$
$V_5$	3	$\text{Gamma}(2, 0.4)$
$K_{5m}$	3	$\text{Gamma}(2, 100)$

The ABC-SysBio software returned the following ranking of models:

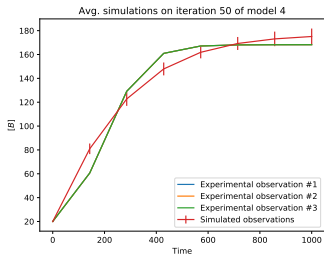
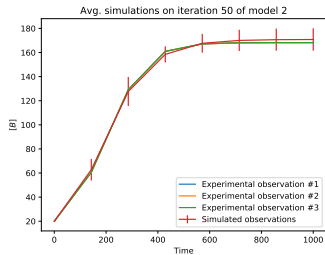
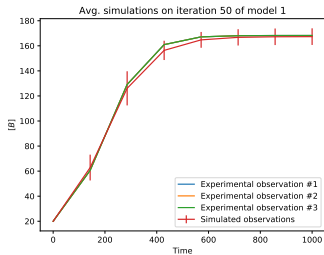
$$2 > 1 > 4 > 3$$



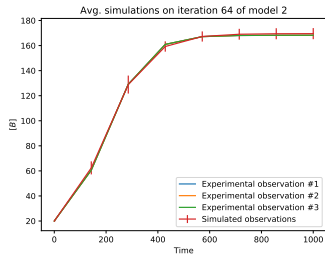
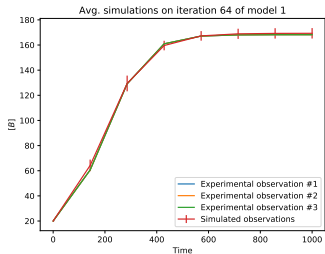
# Results on ABC-SysBio



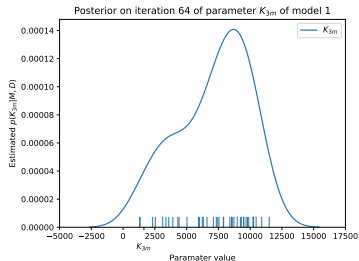
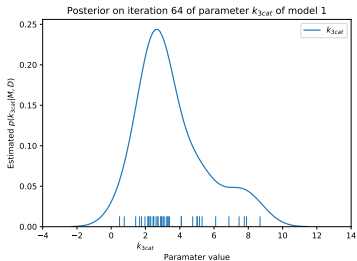
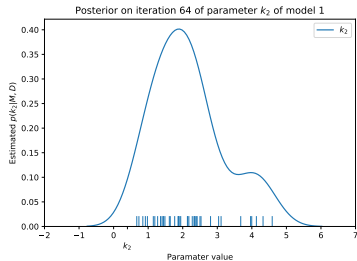
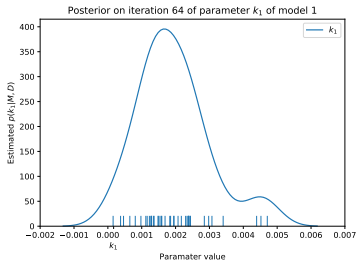
# Results on ABC-SysBio



# Results on ABC-SysBio



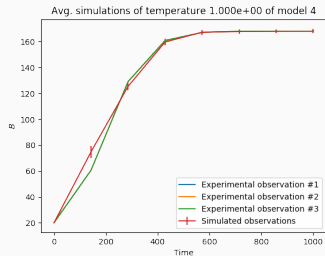
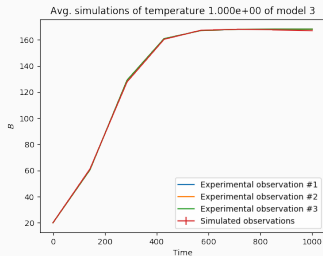
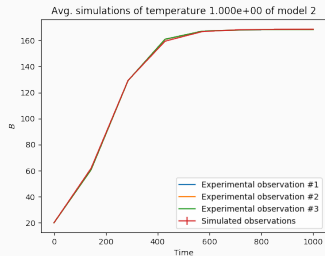
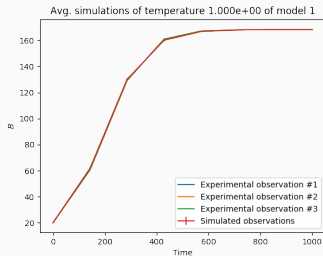
# Results on ABC-SysBio



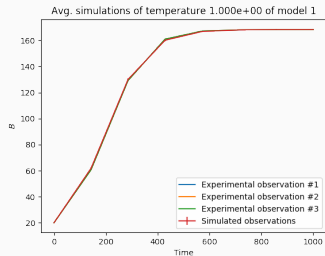
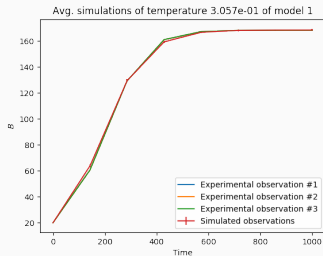
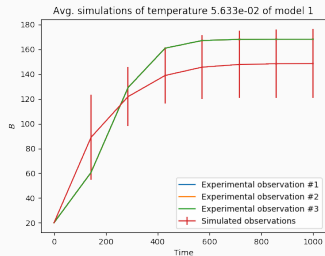
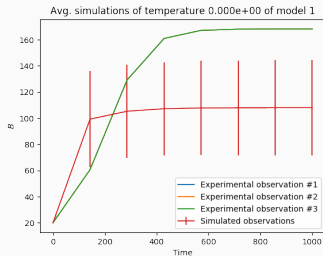
SigNetMS returned the following ranking of models:

$$2 > 1 > 3 > 4$$

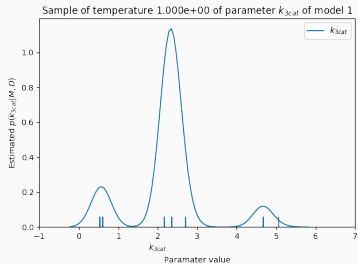
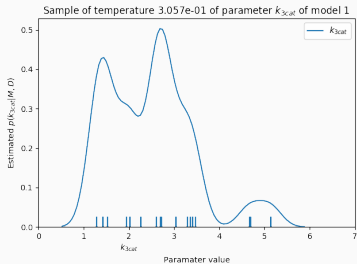
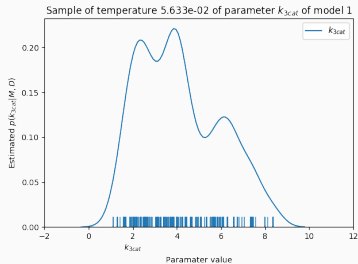
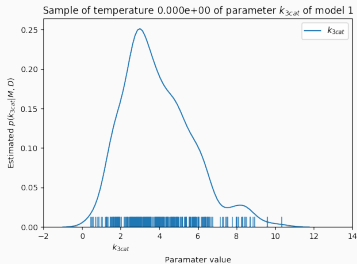
# Results on SigNetMS



# Results on SigNetMS



# Results on SigNetMS





## Future Work

---

More activities are expected to be completed in this project, mainly the follow:

## Future Work

More activities are expected to be completed in this project, mainly the follow:

- Complete program credits requirement.

More activities are expected to be completed in this project, mainly the follow:

- Complete program credits requirement.
- Improve computationally SigNetMS, since we chose it as our cost function for models.

More activities are expected to be completed in this project, mainly the follow:

- Complete program credits requirement.
- Improve computationally SigNetMS, since we chose it as our cost function for models.
- Studies of databases of chemical kinetics such as SABIO-RK and BRENDA.

More activities are expected to be completed in this project, mainly the follow:

- Complete program credits requirement.
- Improve computationally SigNetMS, since we chose it as our cost function for models.
- Studies of databases of chemical kinetics such as SABIO-RK and BRENDA.
- Creation of a relational database of chemical interactions focused on our further applications.

- Define a wrapper on featsel that allows the use of SigNetMS as a cost function.

- Define a wrapper on featsel that allows the use of SigNetMS as a cost function.
- Choice of a feature selection search algorithm.



- Define a wrapper on featsel that allows the use of SigNetMS as a cost function.
- Choice of a feature selection search algorithm.
- Apply the method in ERK signaling pathways of tumor cell lines Y1 and HEK293.

Thank you!