

infer_sde user manual

Cyril Galitzine

May 30, 2018

1 Introduction

infer_sde is a Python code that allows the inference of the distribution of the rates of a stochastic differential equation (SDE) from a single or multiple noisy observation(s) of its trajectory. It is based on a Markov chain Monte Carlo (MCMC) method which samples from the posterior rate distribution. The likelihood of the data is estimated via a particle filter method. Further details about the inference method can be found in [1]

It was initially created to infer the rates governing peroxisome dynamics [1]. In such situation, we are interested in inferring three rates k_d, k_f, γ and the standard deviation of the measurement error σ . The peroxisome count X_t at time t is governed by the following SDE:

$$dX_t = [k_d + (k_f - \gamma) X_t] dt + [k_d + (k_f + \gamma) X_t]^{1/2} dW_t \quad (1)$$

which corresponds to a birth-death-immigration (BDI) stochastic process. We simultaneously observe multiple *realizations* of this SDE by measuring the X_t time course in multiple cells. Each realization of the stochastic process is called a *replicate* in the following the biology terminology.

The code can readily be extended to other types of SDEs because of its modular object oriented design.

2 Running the code

2.1 Input files

2.1.1 sim.dat

This file is used to simulate data which is then stored the data.csv file. A typical sim.dat file is reproduced below:

```
Equation_sim BDI
Error_model_sim Normal
Nrep_sim 4
param_sim [7.75e-4,4.0e-5,4.0e-5]
error_sim [7.5]
Rate_heterogeneity_model_sim Homogeneous
Error_heterogeneity_model_sim Homogeneous
X0_sim [413]
Ntime_sim 120
T_sim 40000
```

The different lines of the file are described below:

- **Equation_sim**: Name of the SDE to simulate. Currently the following SDE are available:
 - BDI: Eq. (1) with parameters k_d, k_f and γ (in that order). In this case the equivalent master equation is simulated via the Gillespie algorithm. This potentially avoid any error in the simulation of the SDE.
- **Error_model_sim**: Error model to use to simulate data. Currently the following models are available:
 - Normal**: Measured values x_{measured} are obtained from the true hidden values, x_{hidden} following $x_{\text{measured}} \sim \text{Normal}(\mu = x_{\text{hidden}}, \sigma = \text{error_sim})$. The standard deviation is assumed to be constant with time and the same for all replicates.

- **Nrep_sim**: Number of replicates to simulate
- **param_sim**: Value of the parameters used to simulate data. The order of the parameters has to correspond to that required by the SDE type.
- **Rate_heterogeneity_model_sim**: Heterogeneity model for the rates between replicates

Homogeneous: The data of all replicates is generated with the exact same rates

Heterogeneous: Rates vary between replicates following a gamma distribution for each rate. When this option is used **param_sim** does not specify directly the rates of the SDE but instead the mean and standard deviation of the gamma distribution of each rate. Assuming that the standard deviation is 10% of the rate mean, **param_sim** would become instead `[7.75e-4,7.75e-5,4.0e-5,4.0e-6,4.0e-5,4.0e-6]` of `[7.75e-4,4.0e-5,4.0e-5]` for Homogeneous rate model.

- **Error_heterogeneity_model_sim**:

Homogeneous: The data of all replicates is generated with the same error standard deviation.

- **X0_sim**: Initial value at time $t = 0$. This value should be specified for each **Nrep_sim** replicate. If it is not the initial value is determined via Poisson distribution with $\lambda = \text{X0_sim}$.
- **Ntime_sim**: Number of time points (A constant Δt is used)
- **T_sim**: Overall simulation time (has to be in the same time unit as the rate specified with **param_sim**).

2.1.2 data.csv

data.csv contains the time course data that is used for the inference. It can be either generated by simulation or obtained through experimental measurement. The first few lines of typical data.csv file are reproduced below:

```
,replicate,t,x
0,0,0.0,420.0
1,0,336.1344537815126,427.0
2,0,672.2689075630252,421.0
3,0,1008.4033613445379,438.0
4,0,1344.5378151260504,456.0
5,0,1680.672268907563,427.0
6,0,2016.8067226890757,409.0
7,0,2352.9411764705883,417.0
8,0,2689.075630252101,416.0
9,0,3025.2100840336134,414.0
10,0,3361.344537815126,416.0
11,0,3697.4789915966385,400.0
12,0,4033.6134453781515,405.0
13,0,4369.747899159664,412.0
```

The first column corresponds to the index of the DataFrame is not important. It does not need to be specified. The following comma separated columns: **replicate**, **t**, **x** should always be present. The number of replicate starts at 0 and should always be specified even if there is only a single replicate.

2.1.3 inference.dat

The inference.dat file controls the inference procedure. A typical inference.dat file (the one used in example BDI_4rep_homogeneous_rates_parallel) is detailed below:

```
Simulate_data 0
Equation BDI
Ndisc 1
Error_model Normal
Nsamp 12000
Npart 1000
param [kd=3.0e-4,kf=3.0e-5,gamma=5.0e-5]
param_infer [1,1,1]
param_error [sigma_error1=4,sigma_error2=4,sigma_error3=4,sigma_error4=4]
param_error_infer [1,1,1,1]
Rate_heterogeneity_model Homogeneous
Error_heterogeneity_model Homogeneous
sd_MH 0.1
MH_step_scaling exponential
output_freq 20
```

- **Simulate_data:** 0: Read existing data.csv file (if it exists), 1: Simulate new data
- **Equation:** Name of the SDE to simulate. Currently the following SDE(s) are available:
BDI: Eq. (1) with parameters k_d, k_f and γ (in that order).
- **Error_model:** Error model for the data. Currently the following models are available:
Normal: Measured values x_{measured} are obtained from the true hidden values, x_{hidden} following $x_{\text{measured}} \sim \text{Normal}(\mu = x_{\text{hidden}}, \sigma = \text{error_sim})$. The standard deviation is assumed to be constant with time for each replicate.
- **Nsamp:** Number of accepted MCMC samples to calculate. MCMC sample proposals that are not accepted are not counted.
- **Npart:** Number of particles to use for the particle filter
- **param:** Initial starting value for the parameters at step 0. The parameters have to be in the correct order dictated by Equation.
- **param_infer:** Infer parameter of **param** 0/1 = Yes/No. In case param_infer=0 for a particular parameter, its value is assumed constant and equal to that specified in **param**.
- **param_error:** Initial starting value for the error parameter at step 0. The names specified for the error parameters do not matter. Typically one error parameter is inferred for each replicate such as here.
- **param_error_infer:** Infer parameter of **param_error** 0/1 = Yes/No. In case param_error_infer=0 for a particular parameter, its value is assumed constant and equal to that specified in **param_error**.
- **Rate_heterogeneity_model:** Heterogeneity model for the rates between replicates
Homogeneous: Assume that all replicates are governed by the exact same rates.
Heterogeneous: Assume that rates vary between replicates following a gamma distribution for each rate. When this option is used **param** does not specify directly the initial rate values of the SDE but instead the initial mean and standard deviation of the gamma distribution of each rate. Assuming that initially the standard deviation is 10% of the rate mean, **param_sim** would become instead [7.75e-4,7.75e-5,4.0e-5,4.0e-6,4.0e-5,4.0e-6] of [7.75e-4,4.0e-5,4.0e-5] for Homogeneous rate model.
- **Error_heterogeneity_model:** Heterogeneity model for the measurement error between replicates

Homogeneous: Assumes that each replicate has a different error standard deviation. This terminology might be misleading but it is used for consistency with the rate heterogeneity model. When this model is used, **param_error** needs to contain one error standard deviation for each replicate.

- **sd_MH:** Parameter controlling the step size of the Metropolis Hastings algorithm. It should be adjusted so that the average acceptance rate is around 0.3.
- **MH_step_scaling:** Metropolis Hasting stepping method. The following methods are available:

exponential: New samples are generated following a log-normal distribution centered on the current value, i.e. $x^{n+1} = x^n \exp(Z)$ with $Z \sim \mathcal{N}(0, \sigma_{\text{MH_step_scaling}})$. This stepping method ensures that each rate step is always properly scaled that the rates remain positive (rates are always defined to be positive).

2.2 Output files

2.2.1 out.dat

This file contains the accepted MCMC samples (rejected samples are not recorded) for each of the parameters listed in the **inference.dat** file. It also contains the overall likelihood of the data. The **out.dat** file obtained with the **inference.dat** file shown above looks like this: (in small fonts to show everything...)

```
iter kd kf gamma sigma_error1 sigma_error2 sigma_error3 sigma_error4 L
0.0000000000000000e+00 3.205407301167742941e-04 2.834254282171117851e-05 5.280902176267724725e-05 4.179274617981517004e+00 3.595000823783015775e+00 3.607774365606607248e+00 4.096877671823151879e+00 -2.20412029
1.0000000000000000e+00 3.189495911384856671e-04 2.933580409225989132e-05 5.478572029166300960e-05 4.904197855752745028e+00 3.607682155543716185e+00 3.782985296115713325e+00 4.250275344357136831e+00 -2.19170308
2.0000000000000000e+00 3.410004245321491753e-04 2.876720992375673006e-05 4.670560555282901182e-05 4.990894925318219322e+00 3.174330773270184647e+00 3.292475942392705690e+00 4.006153697471441788e+00 -2.11997092
3.0000000000000000e+00 3.869474376134196273e-04 2.905421001429177460e-05 4.405738749690098953e-05 5.801727455990557836e+00 2.931393794255999552e+00 3.327187839022418991e+00 3.392726216337757617e+00 -2.07616929
4.0000000000000000e+00 3.485198950970790186e-04 3.392397441063243639e-05 4.874383552471370364e-05 6.023682984290953613e+00 2.223490375545675413e+00 3.268900961215036816e+00 3.415251632976925134e+00 -2.03134317
5.0000000000000000e+00 3.031763779026777367e-04 3.241332428943650030e-05 4.446249326350097653e-05 6.248262705057657129e+00 3.106558659144453394e+00 2.859808957460676204e+00 3.735816211063907932e+00 -2.03167913
6.0000000000000000e+00 3.108710951466856680e-04 3.506106830792371141e-05 4.132333961480131213e-05 6.564705128646511234e+00 3.212402607545694622e+00 2.680282360484247839e+00 3.864462102702313384e+00 -1.98904723
7.0000000000000000e+00 3.399313843000589743e-04 3.696823900838415867e-05 4.084073293389663443e-05 7.597869826320112985e+00 4.072244487243427358e+00 2.878498480183227848e+00 3.917349070467505356e+00 -1.88805845
8.0000000000000000e+00 3.002088842308953684e-04 3.326782355569262752e-05 3.343909947130360629e-05 7.863116032613279849e+00 4.706588199617195656e+00 2.889220926724664462e+00 3.759765641963967564e+00 -1.88807780
9.0000000000000000e+00 2.954710936078864202e-04 3.547788661151836663e-05 3.563202246606440957e-05 8.227071329933828281e+00 4.722409352046546438e+00 2.923040167803822253e+00 3.945890410125978409e+00 -1.86992863
1.0000000000000000e+01 3.229093062367719108e-04 3.989732697050201886e-05 4.436819919805105205e-05 7.976615390598694511e+00 5.040607645098276635e+00 3.184090324627623758e+00 4.118143468436316290e+00 -1.84014224
```

The columns correspond to the samples for rates specified in **parama** and **param_error** (in the same order). The last column is the log likelihood calculated for that particular set of parameters.

3 Running the code

3.1 Overall workflow

3.2 Simulating data with **infer_sde_serial.py**

To simulate data, remove any existing **data.csv** file and run the serial inference code: **python infer_sde_serial.py**. You should have all the *.py files in the directory where you execute this command as well as the following required input files: **sim.dat**, **inference.dat**. If you do not wish to perform inference simply stop the executing of the code after a few steps once you are sure that **data.csv** has been generated.

You can plot the trajectories that you just generated in the **data.csv** file by running **python plot_traj.py** which will produce a plot: **data.png**

It is also likely that you will have to modify the code in **time_series.py** to generate simulated data that conforms to your desired specifications. This is because simulating data according to **sim.dat** can only produce a very narrow set of data types. For instance, it cannot generate replicates with different numbers of time points or overall times.

3.3 Running the serial inference code with **infer_sde_serial.py**

The inference code is executed by typing **python infer_sde_serial.py**. You should have all the *.py files in the directory where you execute this command as well as the following required input files: **sim.dat**, **inference.dat**. Of course the data file **data.csv**. When the code runs, it will output for each accepted sample the value of the parameters and measurement errors, the log likelihood (L), the max log likelihood (called **Lmin**) and the average acceptance ratio **AR**. The value of **sd_MH** in **inference.dat** should be adjusted so that the acceptance ratio is around 0.3 (i.e. the optimum acceptance ratio for the Metropolis Hastings algorithm).

3.4 Running the parallel inference code with `infer_sde_parallel.py`

The inference code was parallelized across replicates using a python version of the MPI library `mpi4py`. This allows us to split the calculation of the overall log likelihood across multiple CPUs to speed up the sampling procedure (but only when multiple replicates are considered). The N_{rep} replicates are distributed as uniformly as possible across CPUs. Each CPU then calculates the sum of the log likelihoods of all the replicates that are assigned to it. The log likelihoods sums are then summed again across cpus to calculate the overall log likelihood. When using N_{CPU} to run the inference, one CPU (the master CPU) will not be performing calculations. As such the calculation will be distributed among the remaining $N_{\text{CPU}} - 1$ CPUs. For instance, if want to infer the rates for $N_{\text{rep}} = 4$ replicates and use $N_{\text{CPU}} = 5$, each CPU will calculate the likelihood of one replicate which leads to the optimal speedup. If $N_{\text{CPU}} = 4$ are used for the same 4 replicates, the replicates will be split as follows across

CPU 2 CPU 3 CPU 4
 CPUs: $\overbrace{\text{Rep1}}^{\text{CPU 2}} \mid \overbrace{\text{Rep2}}^{\text{CPU 3}} \mid \overbrace{\text{Rep3, Rep4}}^{\text{CPU 4}}$ The fourth CPU will calculate the likelihood of two replicates) while CPU 1, the master CPU, is not involved in the calculation of likelihoods. If $N_{\text{CPU}} = 3$, we would obtain the following:
 CPU 2 CPU 3
 $\overbrace{\text{Rep1, Rep2}}^{\text{CPU 2}} \mid \overbrace{\text{Rep3, Rep4}}^{\text{CPU 3}}$. $N_{\text{CPU}} = 2$ should not be employed since the it will less efficient than the serial code the due to the master CPU.

To run the inference with, e.g. $N_{\text{CPUS}} = 4$ the following should be executed:

```
mpirun -np 4 python3 infer_sde_parallel.py
```

The inputs and outputs of the parallel code are identical to those of the serial code.

3.5 Analyzing results with `analyze_mcmc.py`

The inference code (serial or parallel) produces an `out.dat` file with the posterior rate samples. Plots of the distribution of the rates can be obtained by running `python analyze_mcmc.py`. This produces density (`*_density.pdf`) and correlation (`*_corr.pdf`) plots for the inferred parameters. The sampling frequency of the parameters (`frequency` variable) in `analyze_mcmc.py` should be adjusted so that the consecutive samples are not too correlated which otherwise results in an inaccurate estimation of the densities.

4 Test cases

4.1 BDI (birth death immigration) SDE: $dX_t = [k_d + (k_f - \gamma) X_t] dt + [k_d + (k_f + \gamma) X_t]^{1/2} dW_t$

4.1.1 1 replicate, serial (EXAMPLES/BDI_1rep)

All the input, output files and results of the inference can be found in folder EXAMPLES/BDI_1rep. Data was simulated according to the following sim.dat for a single replicate:

```
Equation_sim BDI
Error_model_sim Normal
Nrep_sim 1
Rate_variation_model_sim None
param_sim [7.75e-4,4.0e-5,4.0e-5]
error_sim [6.0]
Rate_heterogeneity_model_sim Homogeneous
Error_heterogeneity_model_sim Homogeneous
X0_sim [413]
Ntime_sim 120
T_sim 40000
```

The serial inference code, `infer_sde_serial.py`, was run for 24 hours and the results were analyzed with `analyze_mcmc.py`.

4.1.2 4 replicates, Homogeneous rates, serial (EXAMPLES/BDI_4rep_homogeneous_rates_serial)

All the input, output files and results of the inference can be found in folder EXAMPLES/BDI_4rep_homogeneous_rates_serial. Data was simulated according to the following sim.dat for 4 replicates (homogeneous rates):

```
Equation_sim BDI
Error_model_sim Normal
Nrep_sim 4
Rate_variation_model_sim None
param_sim [7.75e-4,4.0e-5,4.0e-5]
error_sim [7.5]
Rate_heterogeneity_model_sim Homogeneous
Error_heterogeneity_model_sim Homogeneous
X0_sim [413]
Ntime_sim 120
T_sim 40000
```

The serial inference code, `infer_sde_serial.py`, was run for 24 hours and the results were analyzed with `analyze_mcmc.py`.

4.1.3 4 replicates, Homogeneous rates, parallel (EXAMPLES/BDI_4rep_homogeneous_rates_parallel)

The same data as for case BDI_4rep_homogeneous_rates_serial was used. The parallel inference code was run for 24 hours using 5 CPUs with `mpirun -np 5 python infer_sde_parallel.py`. Results were then analyzed with `analyze_mcmc.py`.

4.1.4 4 replicates, Heterogeneous rates, parallel (EXAMPLES/BDI_4rep_heterogeneous_rates_parallel)

All the input, output files and results of the inference can be found in folder EXAMPLES/BDI_4rep_heterogeneous_rates_parallel. Data was simulated according to the following sim.dat for 4 replicates (heterogeneous rates between replicates):

```
Equation_sim BDI
Error_model_sim Normal
Nrep_sim 4
Rate_variation_model_sim None
param_sim [2.75e-4,2.75e-5,3.0e-5,3.0e-6,5.0e-5,5.0e-6]
error_sim [7.5]
Rate_heterogeneity_model_sim Heterogeneous
Error_heterogeneity_model_sim Homogeneous
```

X0_sim [413]
Ntime_sim 120
T_sim 40000

References

- [1] C. Galitzine, Pierre M. Jean Beltran, I. M. Cristea, and O. Vitek. Statistical inference of peroxisome dynamics. In Benjamin J. Raphael, editor, *RECOMB 2018*, pages 54–74. Springer, 2018.