

# Stat440 Final Project

## Dynamical Analysis of Molecular Interactions

Karme Koo (kmkoo) Wanxin Li (w328li) Jayden Luo (j57luo) Yi Xiang (y25xiang)

### Abstract

Fluorescence resonance energy transfer (FRET) has been used to study biological structures and monitor the activities of molecules (Lemke, Deniz, and Groarke 2016). However, some conformational changes are difficult to detect using ensemble FRET. In order to understand the interaction at single molecule level, single-molecule fluorescence resonance energy transfer (smFRET) is developed to probe the detailed kinetics of structure changes within a single molecule (Ha 2001). By using this revolutionized technique, scientists unlock the opportunities to research more about protein folding-unfolding, protein conformation dynamics, ion channel dynamics, receptor-ligand interactions, nucleic acid structure and conformation, vesicle fusion, and force induced conformational changes (Sasmal et al. 2016). All the new discoveries would reveal the mechanism of diseases at the cellular level (Ideker and Sharan 2008). Current smFRET research employs Brownian Motion (BM) in modelling free diffusion state (Wallace and Atzberger 2017) and Ornstein-Uhlenbeck (OU) process to model the bonding state where there is a consistent bond with stochastic perturbations between donor and acceptor (Yang, Witkoskie, and Cao 2002). One of the remaining challenges of smFRET modelling is whether it is possible to differentiate between free diffusion and bound donor-acceptor pairs. Another challenge is whether we can detect binding and unbinding events. To do so, a simplified experiment set up as follows. The number of photons at time follows a distribution:

$$Y_n \stackrel{ind}{\sim} \text{Poisson}(\exp(\beta_0 + \beta_1 X_t)) \quad (1)$$

where  $X_t$  is the true donor-acceptor distance modelled by BM and OU process.

The challenges of this problem are the following:

- Challenge 1: How to set up experiments so that OU parameters can be accurately estimated
- Challenge 2: Investigate on whether and if applicable, when it is possible to identify the simulation model by combining parameter inferences and model selection criteria.

By a simulation study, our contributions include:

- Implemented parameter inferences for BM and OU models and tested likelihoods and estimates by unit tests
- Found a set of  $\beta$ 's to accurately estimate OU parameters and explained how  $\beta$ 's affect OU parameter estimation
- Found a threshold to distinguish BM simulated data from OU model, vice versa, and explained why the thresholds are reasonable mathematically
- Enabled random start in optimization algorithms

# Introduction

## Brownian Motion

BM is widely used to model the donor-acceptor distance during free diffusion in the molecular environment. The random flow of molecular motion is subordinated by collision of moving molecules. In our dynamical analysis,  $X_t$  is denoted as the donor-acceptor distance at time  $t$ . It is a Gaussian Markov process with the following transition density:

$$X_{t+\Delta t}|X_t \sim N(X_t, \sigma^2 \Delta t) \quad (2)$$

Every subsequent variation in distance by an increment of  $\Delta t$  has a mean-reverting nature. That is, the donor-acceptor distance in the next instance revolves around the mean of the process in various directions. In the BM model, the mean is set as the current-time donor-acceptor distance. With the process involving Markov elements, the sequence of possible molecular motions is only current-state dependent, and it is often referred to as memoryless. The collection of the independent processes at each instance  $\Delta t$ , are recognised as jointly Gaussian (normal).

The variation is determined by the diffusion rate  $\sigma$ , multiplied by  $\Delta t$ . The diffusion rate is set constant in our study, which implies the thermal temperature and pressure are kept unchanged and no external force is applied into the system.

However, the many-body interactions that yield the random pattern cannot be solved by the BM model that accounts every involved molecule. Therefore, BM itself is not capable to model the full motion of molecules especially during the photon exchange between donor and acceptor. This is ultimately a downside of the BM model in our dynamical analysis.

## Ornstein-Uhlenbeck Process

Ornstein-Uhlenbeck process is further developed from BM by L. S. Ornstein and G. E. Uhlenbeck (G. E. and Ornstein 1930). The OU process is a stochastic Gaussian process with continuous paths. The OU process is defined as the following stochastic differential equation:

$$dX_t = \gamma(\mu - X_t)dt + \sigma dW_t \quad (3)$$

Where  $W_t$  is a standard BM on  $t \in (0, \infty)$ ,  $\gamma > 0$ ,  $\sigma > 0$ .  $X_t$  is a stationary Gaussian Markov process with transition density:

$$X_{t+\Delta t}|X_t \sim N(\mu + \omega_{\Delta t}(X_t - \mu), \tau^2(1 - \omega_{\Delta t}^2)) \quad (4)$$

where  $\omega_{\Delta t} = \exp(-\gamma \Delta t)$  and  $\tau^2 = \sigma^2/(2\gamma)$ . Compared to the BM model, the OU process includes more factors from the environment that can impact the change of the distance. The OU process has four components to describe the molecule interactions which are  $\gamma$ ,  $\mu$ ,  $\sigma$ , and  $dW_t$ . The interpretation of  $\gamma$  is the rate of mean reversion. The OU process shares mean reversion nature with BM. This property allows that the distance will eventually revert to the long-run average. The rate indicates how strong the distance will react toward the attractor.  $\mu$  is the asymptotic mean of "bond" length.  $\sigma$  is the parameter showing the volatility of noise (Hirotaka and Kiryu 2016). The OU process represents how molecules move towards the attractors  $\mu$  with BM.

The OU process considers the current state and the speed of molecule interaction, which is preferred when modeling the many-body interactions. Since with different types of donor-acceptor interactions, the molecule will react differently. One of the downsides of the OU process is that a very small amount of variation such as measurement error can profoundly affect the performance of the model. Also the distance is always positive. However, the donor-acceptor in OU allows negative distance.

## Methodology

### Laplace Implementation

Using MLE, we aim to estimate the  $\hat{\theta}$  that maximize  $l(\theta|X, Y)$ , which can be approximated by

$$l(\theta|X, Y) \approx l_{LAP}(\theta|Y) = l(\theta|\hat{X}_{\theta}, Y) - 1/2 * \log|H_{\theta}|$$

$l_{Lap}(\theta|Y)$ .  $\hat{\theta}$  is dependent upon  $\hat{X}_{\theta}$ . In addition, in  $X_{\theta} = \operatorname{argmax}_X l(\theta|X, Y)$ , we need the value of  $X_{\theta}$  for  $\hat{\theta}$ , thus  $X_{\theta}$  and  $\theta$  have to be optimized simultaneously; we use Newton's method to in Laplace approximation to achieve this.

Also, we compare the results of our self-implemented Laplace approximation with TMB's built-in Laplace approximation. Despite of similar accuracies as demonstrated in Appendix, the built-in Laplace implementation is much faster than ours. Thus, we decide to work with the built-in implementation.

### Multi-start Optimization

We initially implemented our likelihood function under the OU model in TMB using  $(\mu, \sigma, \gamma)$  as parameters. However, from our experiments we found that using the R method `optim()` does not produce satisfiable inference result regardless of the choice of method e.g. BFGS and Nelder-Mead. In addition, we noticed that the passing different initial parameter values to `optim()` can caused very different inference results. Therefore we implemented a multi-start on top of `optim()`. We believed that performing multi-start with a parameter that has a lower bound and upper bound such as  $\omega$  is more effective as the boundary allow us to find maximum gap between different starting points. Therefore we implemented another likelihood implementation using  $(\omega, \sigma, \tau)$  as parameters. Below is a RMSE comparison between single start and multi-start.

```
multistart_result <-read.csv("multi_start_result.csv")
multistart_result[, c("num_multistart","rmse.omega","rmse.mu","rmse.tau",
                     "rmse.gamma","rmse.t","rmse.sigma")]

##   num_multistart rmse.omega rmse.mu rmse.tau rmse.gamma rmse.t rmse.sigma
## 1              1       1.50  1.2000  11.000       0.88 330.00      0.190
## 2             10       0.15  0.0095   0.049       0.15  0.15      0.065
```

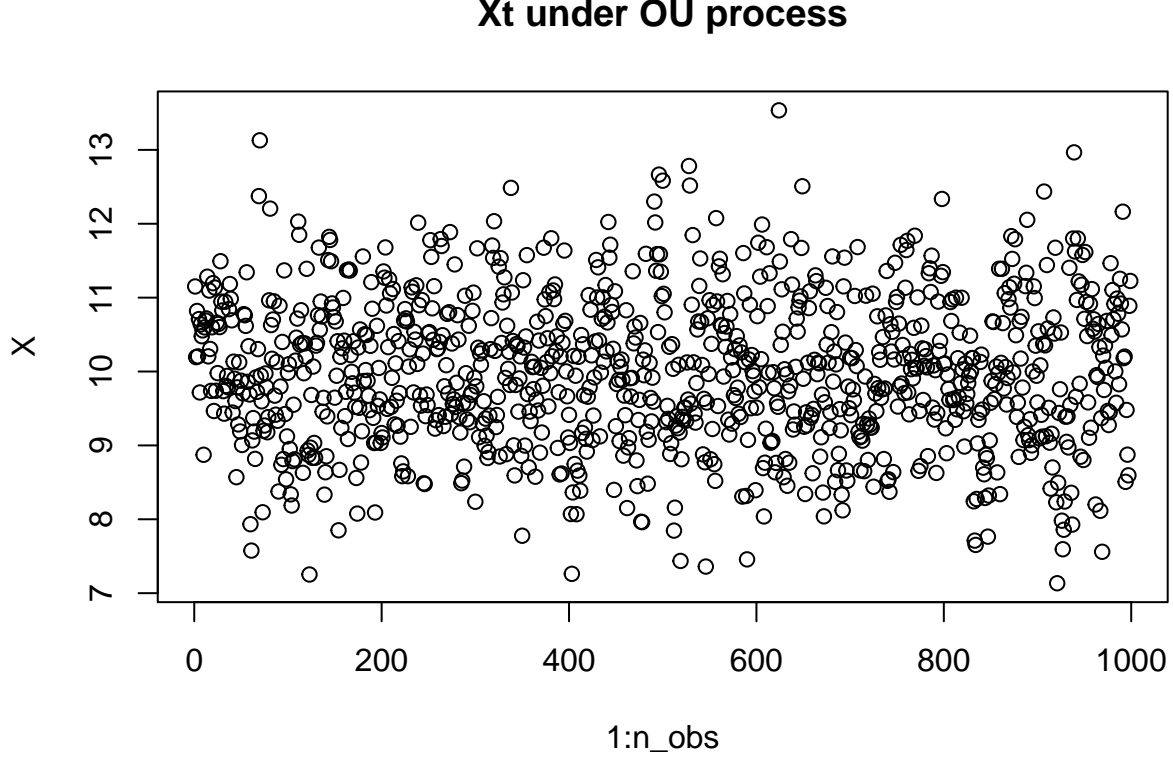
### NA Handling

For Challenge 2, whenever NA is simulated in  $Y_t$  due to large  $\lambda$  in *rpois*, we disgard it together with corresponding  $X_t$ .

### Simulation for Challenge 1

The goal here is to investigate on how the value of  $(\beta_0, \beta_1)$  could affect the accuracy of inference under the OU model. We setup a experiment with true parameters being  $\mu = 10$ ,  $\tau = 1$  and  $\gamma \in \{0.1, 1, 10\}$ . In order to evaluate the accuracy of inference, we simulate 100 dataset for each pair of  $(\beta_0, \beta_1)$  and calculated the root mean-squared error(RMSE) of  $t_{dec}$ ,  $\mu$  and  $\tau$  as measure of accuracy. We first find a initial pair of  $(\beta_0, \beta_1)$  such that the RMSE reasonably good. We looked into the simulated dataset when  $\gamma = 1$ , shown below. We observed that the values of  $X_n$  are mostly below 13 and hence set a initial  $\beta_1 = 1$  and  $\beta_0 > 13$  such that we have  $\beta_0 - \beta_1 X_n > 0$  and  $Y_n \sim \text{Poisson}(\exp(\beta_0 - \beta_1 X_n))$  are more likely to have a large set of values. For  $\gamma = 1$ , we decided to simulate 299 observations per dataset. For  $\gamma = 0.1$ , since we were not able to get good RMSE with 299 observations per dataset, we increased to 999 per dataset. We then perform inference simulations with one of the the  $\beta(\beta_0 \text{ or } \beta_1)$  fixed and the other gradually decreasing and observe how RMSE changes.

```
n_obs = 999
X <- ou_sim(gamma=1, mu=10, sigma=sqrt(2), dt=1, n_obs=999)
plot(1:n_obs, X, main="Xt under OU process")
```



## Simulation for Challenge 2

To start, we simulate from OU process. We perform the experiment 20 times for a fixed set of  $\beta_0, \beta_1, \gamma, \mu$ . 200  $X'_t$ s and  $Y'_t$ s without NA are generated in each experiment, and parameter inferences are performed. Parameter inferences are performed using “Nelder-Mead” optimization. Similarly, we simulate from BM process. We perform the experiment 20 times for a fixed set of  $\beta_0, \beta_1$ . 400  $X'_t$ s and  $Y'_t$ s without NA are generated in each experiment. Parameter inferences are performed using “BFGS” optimization. Our setups are as follows: OU setup, optimization method is “Nelder Mead”, number of observation is 100.

- setup 1:  $\beta_0 = 10, \beta_1 = 0.5, \tau = 1, \mu = 0$
- setup 2:  $\beta_0 = 5, \beta_1 = 0.5, \tau = 2, \mu = 0$
- setup 3:  $\beta_0 = 5, \beta_1 = 1, \tau = 2, \mu = 0$
- setup 4:  $\beta_0 = 10, \beta_1 = 0.5, \tau = 2, \mu = 1$

And BM setup, optimization method is "BFGS, number of observation is 400.

- setup 1:  $\beta_0 = 10, \beta_1 = 0.5, \mu = 0, \gamma = 5$
- setup 2:  $\beta_0 = 5, \beta_1 = 0.5, \mu = 0, \gamma = 5$
- setup 3:  $\beta_0 = 5, \beta_1 = 1, \mu = 0, \gamma = 5$

## Results

We present and analyze our results with respect to two challenges as follows.

### Challenge 1

The result of our simulations are as follows:

```
sim1_result <- read.csv("sim1_result.csv")
sim1_result <- sim1_result[sim1_result["num_multistart"]==10,]
show_col <- c("beta0", "beta1", "n_obs", "theta.mu", "theta.t", "theta.tau",
             "rmse.mu", "rmse.t", "rmse.tau")
sim1_result <- sim1_result[with(sim1_result, order(-beta1, -beta0)),]
```

### gamma=1

Below is the simulation result with  $\gamma = 1$ . We were able to get  $\text{RMSE}_\theta < 1\%$  and  $\text{RMSE}_\tau < 10\%$  but  $\text{RMSE}_t < 20\%$  only. The result shows that when  $\beta_1 = 1$  the minimum value of  $\beta_0$  to maintain a good RMSE is between 11 and 12. The result shows that when  $\beta_0 = 15$  the minimum value of  $\beta_1$  to maintain a good RMSE is between 0.8 and 1.

```
signif(sim1_result[sim1_result[, "gamma"]==1, show_col], 3)
```

##	beta0	beta1	n_obs	theta.mu	theta.t	theta.tau	rmse.mu	rmse.t	rmse.tau
## 2	15	1.0	299	10	1	1	0.0095	1.5e-01	0.049
## 19	14	1.0	299	10	1	1	0.0088	1.7e-01	0.051
## 20	13	1.0	299	10	1	1	0.0081	1.5e-01	0.049
## 21	12	1.0	299	10	1	1	0.0083	1.7e-01	0.057
## 3	11	1.0	299	10	1	1	0.3300	1.8e+02	3.200
## 4	9	1.0	299	10	1	1	0.8800	1.8e+11	3.500
## 5	15	0.8	299	10	1	1	0.0230	4.2e+00	0.570
## 6	15	0.6	299	10	1	1	0.0680	1.4e+01	1.600

### gamma=0.1

The simulation result with  $\gamma = 1$  is shown below. In this setup, we were able to get  $\text{RMSE}_\theta < 5\%$  and  $\text{RMSE}_\tau < 10\%$  but  $\text{RMSE}_t < 20\%$  only. The result shows that our model were not able to achieve good accuracy when there are only 299 observations in the dataset. However, with the number of observations being 999, the model were able to get much better *RMSE*. In this setup, the minimum value of  $\beta_0$  for good accuracy is between 9 and 11 when  $\beta_1 = 1$ , and the minimum value of  $\beta_1$  for good accuracy is between 0.9 and 1.

```
signif(sim1_result[sim1_result[, "gamma"]==0.1, show_col], 3)
```

##	beta0	beta1	n_obs	theta.mu	theta.t	theta.tau	rmse.mu	rmse.t	rmse.tau
## 7	15	1.0	299	10	10	1	0.048	3.00	0.700
## 8	15	1.0	999	10	10	1	0.016	0.15	0.070
## 22	14	1.0	999	10	10	1	0.014	0.14	0.072
## 23	13	1.0	999	10	10	1	0.014	0.14	0.065
## 24	12	1.0	999	10	10	1	0.013	0.16	0.074
## 9	11	1.0	999	10	10	1	0.015	0.18	0.084
## 10	9	1.0	999	10	10	1	0.094	19.00	1.100
## 25	15	0.9	999	10	10	1	0.015	0.29	0.110
## 11	15	0.8	999	10	10	1	0.023	1.20	0.370
## 12	15	0.6	999	10	10	1	0.022	1.70	0.580

### gamma=10

For  $\gamma = 10$ , as shown below, we were able to get  $\text{RMSE}_\theta < 1\%$  and  $\text{RMSE}_\tau < 10\%$  but not able to estimate parameter  $t$  nicely even with 999 observations per dataset. The lower bound for  $\beta_0$  for a good  $\text{RMSE}_\tau$  is between 11 and 12 and the lower bound for  $\beta_1$  for good  $\text{RMSE}_\tau$  is between 0.6 and 0.8.

```
signif(sim1_result[sim1_result[, "gamma"]==10, show_col], 3)
```

	beta0	beta1	n_obs	theta.mu	theta.t	theta.tau	rmse.mu	rmse.t	rmse.tau
## 13	15	1.0	299	10	0.1	1	0.0049	1.6e+00	0.042
## 14	15	1.0	999	10	0.1	1	0.0032	1.4e+00	0.024
## 26	14	1.0	999	10	0.1	1	0.0036	1.3e+00	0.024
## 27	13	1.0	999	10	0.1	1	0.0032	1.0e+00	0.027
## 28	12	1.0	999	10	0.1	1	0.0036	1.1e+00	0.030
## 15	11	1.0	999	10	0.1	1	1.0000	5.8e+03	9.400
## 16	9	1.0	999	10	0.1	1	0.7000	1.9e+12	4.300
## 29	15	0.9	999	10	0.1	1	0.0034	1.1e+00	0.023
## 17	15	0.8	999	10	0.1	1	0.0032	1.2e+00	0.022
## 18	15	0.6	999	10	0.1	1	0.0330	4.5e+01	0.930

## Challenge 2

For  $X_t$  simulated from OU model, as shown below, for 4 setups, the probabilities to pick OU model by AIC increase as  $\gamma$ 's go up from 0.01 to 4.

```
sim2_OU_result <- read.csv("sim2_table1.csv", check.names=FALSE)

## Warning in read.table(file = file, header = header, sep = sep, quote = quote, :
## incomplete final line found by readTableHeader on 'sim2_table1.csv'

print(sim2_OU_result,row.names = FALSE)

##      gamma 0.01 0.05 1 2 3 4
## ou setup 1  0.40 0.80 1 1 1 1
## ou setup 2  0.45 0.65 1 1 1 1
## ou setup 3  0.45 0.55 1 1 1 1
## ou setup 4  0.30 0.65 1 1 1 1
```

For  $X_t$  simulated from BM model, as shown below, for 4 setups, the probabilities to pick BM model by AIC increase as  $\sigma$ 's go up from 1e-6 to 2. However, when  $\sigma = 3$ , the probabilities goes down.

```
sim2_BM_result <- read.csv("sim2_table2.csv", check.names=FALSE)

## Warning in read.table(file = file, header = header, sep = sep, quote = quote, :
## incomplete final line found by readTableHeader on 'sim2_table2.csv'

print(sim2_BM_result,row.names = FALSE)

##      sigma 1E-6 1E-5 1E-4 1E-3 1E-2 1E-1 1 2 3
## bm setup 1  0.00 0.05 0.00 0.90 0.80 0.9 0.95 1.00 0.9
## bm setup 2  0.05 0.00 0.00 0.05 0.55 0.7 0.10 0.85 0.6
## bm setup 3  0.05 0.05 0.05 0.20 0.75 0.7 0.70 0.90 0.6
```

## Discussion

### Other Potential Models for Molecular Dynamics Study

#### Morse Interaction Model

Morse Interaction model(Schelstraete, Schepens, and Verschelde 1999) is described as a combination of BM and improved OU process. Under this model,  $X_t$  is a Markov process satisfying the stochastic differential equation (SDE):

$$dX_t = -U'(X_t)dt + \sigma dB_t$$

And  $U'(x)$  is the derivative of the Morse potential energy function:

$$U'(x) = \gamma \cdot (1 - e^{-\alpha \cdot (x-u)})$$

The distance  $X_t$  is set to be strictly positive,  $X_t > 0$ . When  $X_t$  is too large, repulsive forces allow bond breaking to occur - the molecules again resemble BM. When the donor and acceptor again get close to each other, the bond is reformed and exchange of photons takes place. The dynamics then resembles OU process.

Morse interaction model is a comprehensive approach towards dynamics study of molecular interaction.

## Lennard-Jones Potential

The interaction between two non-bonded and un-charged atoms, known as Van der Waals interaction, has been expressed in terms of potential energy. Lennard-Jones potential (Libretexts 2019) is probably the most famous pair potential describing interatomic Van der Waals forces. It consists of two parts:

- 1) A steep repulsive term; and
- 2) A smooth attractive term

$$V(r) = 4\epsilon \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^6 \right]$$

In particular,  $V(x)$  is the intermolecular potential between the two molecules, and  $r$  is the distance of separation between both molecules.

Apart from being a widely used model itself, Lennard-Jones potential also sometimes forms one of ‘building blocks’ of many force fields, due to its computational expediency.

However, Lennard-Jones potential is not designated to model donor-acceptor interactions.

## Challenge 1

### Lower Bounds of beta0

Under  $\mu = 10, \tau = 1$ , we simulated 1000 datasets and find the mean and range of the simulated  $X_n$  datasets as below:

```
sapply(c(0.1,1,10), function(gamma) {
  ou_X <- ou_sim(gamma=gamma, mu=10, sigma=sqrt(2*gamma), dt=1, n_obs=999)
  c(gamma=gamma, min=min(ou_X), mean=mean(ou_X), max=max(ou_X))
})
```

```
##           [,1]      [,2]      [,3]
## gamma  0.100000  1.000000 10.000000
## min    7.138558  6.729285  6.477269
## mean   9.899166 10.025087 10.027157
## max   13.160638 14.005026 13.122742
```

With  $\beta_1 = 1$ , when  $\beta_0$  drops under 13,  $\beta_0 - \beta_1 X_n$  has a probability of becoming negative. As  $\beta_0$  continue to decrease, the probability to become negative increases. When  $\beta_0 - \beta_1 X_n$  is negative and  $\exp(\beta_0 - \beta_1 X_n) < 1$ , the  $Y_n \sim \text{Poisson}(\exp(\beta_0 - \beta_1 X_n))$  generate only few values, with large probability being either 0 or 1. This means that the characteristics of  $X_n$  datasets are not reflected to  $Y_n$  and inference on  $X_n$  parameters using  $Y_n$  data will perform poorly.

### Lower Bounds of beta1

To understand why the value of  $\beta_1$  affect inference result, we simulated 100  $X_n$  and  $Y_n$  datasets of size 299 and computed their corresponding  $p(Y|X)$ .

```
vec_beta1 <- seq(0.1,1,by=0.1)
loglik <- sapply(vec_beta1, function(beta1) {
  beta0 <- 15
  X <- ou_sim(gamma=1, mu=10, sigma=sqrt(2), dt=1, n_obs=299)
```

```

Y <- y_sim(X, beta0=beta0, beta1=beta1)
sum(dpois(Y,exp(beta0-beta1*X), log=TRUE))
})
rbind(beta1=vec_beta1, loglik=loglik)

```

```

##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## beta1      0.100      0.200      0.300      0.400      0.500      0.600      0.700
## loglik -2539.005 -2339.053 -2233.243 -2102.789 -1897.939 -1742.476 -1636.847
##           [,8]      [,9]     [,10]
## beta1      0.800      0.900      1.000
## loglik -1444.945 -1314.853 -1189.851

```

This shows that under the same parameter and  $\beta_0$ , different  $\beta_1$  can have different loglikelihood. From this observation, we further investigated into  $l(\theta|X, Y) = \log\{p(Y|X) \times p(X|\theta)\}$  and found that changes in  $\beta_1$  affect the term  $\sum_{n=0}^N Y_n(\beta_0 - \beta_1 X_n) = \sum_{n=0}^N Y_n \beta_0 - \beta_1 X_n Y_n$ . Specifically, if we consider  $\beta_1$  as a weight for the term  $X_n Y_n$ , then decreasing  $\beta$  is reducing the impact of the term  $X_n Y_n$  on  $l(\theta|X, Y)$ . We notice that this is the only term  $Y_n$ -related term in finding  $\hat{X}_\theta$  in laplace approximation. Therefore, in finding  $\hat{X}_\theta$ , decreasing  $\beta_1$  acts as decreasing the importance of the value of  $Y_n$ . We suspect that this is the reason why  $\beta_1$  has a lower bound. However, even if this is the case, it is difficult to determine what the lower bound is without performing simulation.

### Why n\_obs has Impact on RMSE when gamma=0.1

When  $\gamma = 0.1$ , equation 2 is approximately  $X_{t+\Delta t}|X_t \sim N(0.1\mu + 0.9X_t, 0.2\tau)$ , which means that the value of  $X_t$  has plays important factor in the likelihood of  $X_{t+\Delta t}$ . Since  $X_t$  are not given and only approximated using laplace, inference would depend more on Laplace approximation. Since  $X_t$  are approximated using MLE and MLE methods are generally affected number of observations, we suspect that the performance of Laplace approximation in this model would also depend on dataset size.

## Challenge 2

### OU Simulation

For experiments simulated from OU, we can see a consistent improvement in picking the correct model when  $\gamma$  goes up. As  $\gamma \rightarrow 0$ ,  $\omega_{\Delta t} \rightarrow 1$  and  $\tau$  is still a constant, equation 4 becomes  $X_{t+\Delta t}|X_t \sim N(X_t, 0)$ , which is in the same form as equation 2. Thus, it is difficult to distinguish the OU simulation from the BM model. As  $\gamma$  goes up,  $\omega_{\Delta t} \rightarrow 0$ , 2 and 4 differ in mean and variance, and hence become easier to distinguish.

### BM Simulation

For experiments simulated from BM, we can see a consistent improvement in picking the correct model when  $\sigma_{BM}$  goes up, however, the improvement stops at  $\sigma_{BM} = 3$  for most of the experiments. We analyze in the following two aspects. First of all, when  $\sigma_{BM} \rightarrow 0$ , equation 2 becomes  $X_{t+\Delta t}|X_t \sim N(X_t, 0)$ . From the discussion in the previous paragraph, when  $\gamma \rightarrow 0$ , equation 4 also becomes  $X_{t+\Delta t}|X_t \sim N(X_t, 0)$ . Thus, if we do parameter inferences are done correctly, we can also find OU parameters that fit the data well. It will be difficult to tell the BM simulated data from OU models. Furthermore, when  $\sigma$  becomes too large with respect to  $\beta_0$  and  $\beta_1$ , for example, when  $\sigma = 3$  with respect to  $\beta_0 = 5$  and  $\beta_1 = 1$ , the probability of generating large  $X_t$  becomes higher. In this case,  $\exp(\beta_0 - \beta_1 X_t)$  becomes close to 0, and the generated  $Y_t$ 's from Poission equation 1 are more likely to be 0's. Consequently,  $Y_t$  does not provide enough information to the underlying generating process, and hence difficult to tell apart BM model from OU model.

Another observation to note is that as a general rule,  $|\beta_0 - \beta_1|$  cannot be too large; otherwise, many NAs from equation 1 will be generated because R cannot handle values  $> e^{22}$ . Referring to how we handle NAs in the Methodology section, repetively generating values until no NAs results will result in huge computation cost.



## **Potential Improvements in Implementation**

### **Other Model Selection Criterion**

Only AIC is used in model selection for Challenge 2. We can consider other model selection criteria for Challenge 2. For example, AIC and Mallow CP. Using a weighted of different criteria for model selection makes the result less biased.

### **Random Start for BM model**

As discussed in the methodology section, random. We figured out how to implement random start for OU model by reparameterization. Consequently, the accuracies of OU inferences have been improved significantly. It will be useful to figure out similar approaches for BM model.

### **R Drawback**

One of the challenges encountered in the project is the efficiency of running simulation results. Each simulation with different parameters took over 1 hour to obtain the final results. This exposes one of the disadvantages of R, which is that R does not support running test cases in parallel.

## Appendix

### 1. csv for multistart in Methodology

```
multistart_result <-read.csv("multi_start_result.csv")
multistart_result
```

```
##      beta0 beta1 gamma n_dataset n_obs num_multistart theta.omega theta.mu
## 1      15      1      1        100   299              1  0.3678794      10
## 2      15      1      1        100   299             10  0.3678794      10
##      theta.tau theta.t theta.gamma theta.sigma rmse.omega rmse.mu rmse.tau
## 1              1          1          1    1.414214      1.50  1.2000    11.000
## 2              1          1          1    1.414214      0.15  0.0095     0.049
##      rmse.gamma rmse.t rmse.sigma
## 1              0.88 330.00      0.190
## 2              0.15  0.15      0.065
```

### 2. Code to generate results for Challenge 1

#### (1) test cases

```
sim1_cases <- read.csv("sim1_cases_new.csv")
sim1_cases
```

```
##      beta0 beta1 gamma n_dataset n_obs num_multistart
## 1      15      1.0  1.0        100   299              1
## 2      15      1.0  1.0        100   299             10
## 3      11      1.0  1.0        100   299             10
## 4       9      1.0  1.0        100   299             10
## 5      15      0.8  1.0        100   299             10
## 6      15      0.6  1.0        100   299             10
## 7      15      1.0  0.1        100   299             10
## 8      15      1.0  0.1        100   999             10
## 9      11      1.0  0.1        100   999             10
## 10     9      1.0  0.1        100   999             10
## 11     15      0.8  0.1        100   999             10
## 12     15      0.6  0.1        100   999             10
## 13     15      1.0 10.0        100   299             10
## 14     15      1.0 10.0        100   999             10
## 15     11      1.0 10.0        100   999             10
## 16      9      1.0 10.0        100   999             10
## 17     15      0.8 10.0        100   999             10
## 18     15      0.6 10.0        100   999             10
## 19     14      1.0  1.0        100   299             10
## 20     13      1.0  1.0        100   299             10
## 21     12      1.0  1.0        100   299             10
## 22     14      1.0  0.1        100   999             10
## 23     13      1.0  0.1        100   999             10
## 24     12      1.0  0.1        100   999             10
## 25     15      0.9  0.1        100   999             10
## 26     14      1.0 10.0        100   999             10
## 27     13      1.0 10.0        100   999             10
## 28     12      1.0 10.0        100   999             10
## 29     15      0.9 10.0        100   999             10
## 30     15      0.8  1.0        100   299             10
## 31     15      0.8  1.0        100  1999             10
```

## (2) Script

```
require("smfret")

#' Run simulation 1 on all given cases and save RMSE in a csv
#'
#' @param all_cases dataframe of cases to teston. Each row is a test test case.(see **Details**)
#' @param my_ci the indices of cases to test on
#' @param out_file the name of the output file
#' @return
#' @details
#' 1. parameter `all_cases` must has column
#'   `c("beta0","beta1","gamma","n_dataset","n_obs","num_multistart")`.
#' 2. This function will not overwrite existing files. Instead, a warning will show
#'   and output will be saved in a temporary file.
#' 3. To save time, it is recommended that you run cases in different processes and
#'   merge result afterwards.
#' @examples
#' sim_1_script(out_file="testing.csv", my_ci=1:2)
sim_1_script <- function(all_cases=read.csv("sim1_cases_new.csv"),
  my_ci=1:nrow(all_cases), out_file=paste0("sim1_result_",my_ci, ".csv")) {
  my_cases <- all_cases[my_ci,]
  result <- apply(my_cases, 1, function(tc) {
    print(tc)
    sim_1(
      beta0=as.numeric(tc[["beta0"]]), beta1=as.numeric(tc[["beta1"]]),
      omega=exp(-as.numeric(tc[["gamma"]])),
      n_dataset=as.numeric(tc[["n_dataset"]]), n_obs=as.numeric(tc[["n_obs"]]),
      num_multistart=as.numeric(tc[["num_multistart"]]))
  })
  rmse <- cbind(my_cases, t(sapply(1:nrow(my_cases), function(i) {
    c(theta=result[[i]]$true_param, rmse=result[[i]]$rmse)
  })))
  rmse <- apply(rmse, 2, as.numeric)
  if (file.exists(out_file)) {
    tmp_out_file = tempfile()
    warning(paste("ERROR IN SIMULATION 1:", out_file, "exists... Saving to",
      tmp_out_file, "Instead. Please manually retriive file if desired.))
    write.csv(rmse, file=tmp_out_file, row.names=FALSE)
  } else{
    write.csv(rmse, file=out_file, row.names=FALSE)
  }
}

hi <- sim_1_script(out_file="testing.csv", my_ci=1:2)
```

## (3) simulated result

```
read.csv('sim1_result.csv')
```

##	beta0	beta1	gamma	n_dataset	n_obs	num_multistart	theta.omega	theta.mu
## 1	15	1.0	1.0	100	299	1	3.678794e-01	10
## 2	15	1.0	1.0	100	299	10	3.678794e-01	10
## 3	11	1.0	1.0	100	299	10	3.678794e-01	10
## 4	9	1.0	1.0	100	299	10	3.678794e-01	10

## 5	15	0.8	1.0	100	299	10	3.678794e-01	10
## 6	15	0.6	1.0	100	299	10	3.678794e-01	10
## 7	15	1.0	0.1	100	299	10	9.048374e-01	10
## 8	15	1.0	0.1	100	999	10	9.048374e-01	10
## 9	11	1.0	0.1	100	999	10	9.048374e-01	10
## 10	9	1.0	0.1	100	999	10	9.048374e-01	10
## 11	15	0.8	0.1	100	999	10	9.048374e-01	10
## 12	15	0.6	0.1	100	999	10	9.048374e-01	10
## 13	15	1.0	10.0	100	299	10	4.539993e-05	10
## 14	15	1.0	10.0	100	999	10	4.539993e-05	10
## 15	11	1.0	10.0	100	999	10	4.539993e-05	10
## 16	9	1.0	10.0	100	999	10	4.539993e-05	10
## 17	15	0.8	10.0	100	999	10	4.539993e-05	10
## 18	15	0.6	10.0	100	999	10	4.539993e-05	10
## 19	14	1.0	1.0	100	299	10	3.678794e-01	10
## 20	13	1.0	1.0	100	299	10	3.678794e-01	10
## 21	12	1.0	1.0	100	299	10	3.678794e-01	10
## 22	14	1.0	0.1	100	999	10	9.048374e-01	10
## 23	13	1.0	0.1	100	999	10	9.048374e-01	10
## 24	12	1.0	0.1	100	999	10	9.048374e-01	10
## 25	15	0.9	0.1	100	999	10	9.048374e-01	10
## 26	14	1.0	10.0	100	999	10	4.539993e-05	10
## 27	13	1.0	10.0	100	999	10	4.539993e-05	10
## 28	12	1.0	10.0	100	999	10	4.539993e-05	10
## 29	15	0.9	10.0	100	999	10	4.539993e-05	10
##	theta.tau	theta.t	theta.gamma	theta.sigma	rmse.omega	rmse.mu	rmse.tau	
## 1	1	1.0	1.0	1.4142136	1.5e+00	1.2000	11.000	
## 2	1	1.0	1.0	1.4142136	1.5e-01	0.0095	0.049	
## 3	1	1.0	1.0	1.4142136	3.4e-01	0.3300	3.200	
## 4	1	1.0	1.0	1.4142136	1.2e+00	0.8800	3.500	
## 5	1	1.0	1.0	1.4142136	3.6e-01	0.0230	0.570	
## 6	1	1.0	1.0	1.4142136	9.4e-01	0.0680	1.600	
## 7	1	10.0	0.1	0.4472136	4.4e-02	0.0480	0.700	
## 8	1	10.0	0.1	0.4472136	1.6e-02	0.0160	0.070	
## 9	1	10.0	0.1	0.4472136	2.1e-02	0.0150	0.084	
## 10	1	10.0	0.1	0.4472136	4.2e-02	0.0940	1.100	
## 11	1	10.0	0.1	0.4472136	3.3e-02	0.0230	0.370	
## 12	1	10.0	0.1	0.4472136	5.7e-02	0.0220	0.580	
## 13	1	0.1	10.0	4.4721360	8.6e+02	0.0049	0.042	
## 14	1	0.1	10.0	4.4721360	6.1e+02	0.0032	0.024	
## 15	1	0.1	10.0	4.4721360	5.8e+03	1.0000	9.400	
## 16	1	0.1	10.0	4.4721360	1.8e+04	0.7000	4.300	
## 17	1	0.1	10.0	4.4721360	4.5e+02	0.0032	0.022	
## 18	1	0.1	10.0	4.4721360	7.3e+03	0.0330	0.930	
## 19	1	1.0	1.0	1.4142136	1.7e-01	0.0088	0.051	
## 20	1	1.0	1.0	1.4142136	1.5e-01	0.0081	0.049	
## 21	1	1.0	1.0	1.4142136	1.8e-01	0.0083	0.057	
## 22	1	10.0	0.1	0.4472136	1.6e-02	0.0140	0.072	
## 23	1	10.0	0.1	0.4472136	1.5e-02	0.0140	0.065	
## 24	1	10.0	0.1	0.4472136	1.7e-02	0.0130	0.074	
## 25	1	10.0	0.1	0.4472136	1.8e-02	0.0150	0.110	
## 26	1	0.1	10.0	4.4721360	5.4e+02	0.0036	0.024	
## 27	1	0.1	10.0	4.4721360	3.9e+02	0.0032	0.027	
## 28	1	0.1	10.0	4.4721360	4.4e+02	0.0036	0.030	

```
## 29      1      0.1      10.0  4.4721360    4.7e+02  0.0034    0.023
##      rmse.gamma  rmse.t  rmse.sigma
## 1      0.88 3.3e+02    0.190
## 2      0.15 1.5e-01    0.065
## 3      0.31 1.8e+02    0.130
## 4     11.00 1.8e+11    2.100
## 5      0.26 4.2e+00    0.078
## 6      0.58 1.4e+01    0.130
## 7      0.43 3.0e+00    0.043
## 8      0.16 1.5e-01    0.027
## 9      0.21 1.8e-01    0.053
## 10     0.43 1.9e+01    0.150
## 11     0.32 1.2e+00    0.022
## 12     0.55 1.7e+00    0.029
## 13     0.50 1.6e+00    0.310
## 14     0.49 1.4e+00    0.310
## 15     0.71 5.8e+03    0.390
## 16     0.94 1.9e+12    0.840
## 17     0.43 1.2e+00    0.270
## 18     0.58 4.5e+01    0.370
## 19     0.20 1.7e-01    0.072
## 20     0.16 1.5e-01    0.068
## 21     0.20 1.7e-01    0.091
## 22     0.16 1.4e-01    0.029
## 23     0.15 1.4e-01    0.035
## 24     0.17 1.6e-01    0.047
## 25     0.18 2.9e-01    0.027
## 26     0.45 1.3e+00    0.280
## 27     0.42 1.0e+00    0.250
## 28     0.50 1.1e+00    0.270
## 29     0.42 1.1e+00    0.260
```

### 3. Code to generate results for Challenge 2 (simulation from OU)

#### (1) Input file

```
sim2_OU_cases <- read.csv('sim2_OU_cases.csv')
sim2_OU_cases
```

```
##      X beta0 beta1 tau mu
## 1 1      10   0.5   1  0
## 2 2       5   0.5   2  0
## 3 3       5   1.0   2  0
## 4 4      10   0.5   2  1
```

#### (2) Script

```
require("smfret")

all_cases_default <- read.csv("sim2_OU_cases.csv", row.names=1)
test_cases_default <- data.frame(gamma=c(0.01, 0.05, 1, 2, 3, 4))
output_file_default <- "sim2_ou_results.RData"

#' Generate observations from a Morse SDE.
#'
#' @param all_cases dataframe of cases to test. Each row is a test case.(see *Details*)
```

```

#' @param test_cases dataframe of sigmas to test
#' @param output_file the name of the output file
#' @return
#' @details
#' 1. parameter `all_cases` must has column
#' `c(beta0, "beta1", "tau", "mu")`.
#' 2. parameter `test_cases` must have column "gamma"
sim2_OU_script <- function(all_cases=all_cases_default, test_cases=test_cases_default,
  output_file=output_file_default) {
  for(i in 1:nrow(all_cases)) {
    row <- all_cases[i,]
    cur_beta0 <- row[['beta0']]
    cur_beta1 <- row[['beta1']]
    cur_tau <- row[['tau']]
    cur_mu <- row[['mu']]
    ou_result <- apply(test_cases, 1, function(info) {
      gamma <- info[["gamma"]]
      sigma <- sqrt(2*gamma)*cur_tau

      sim2(from="ou", n_dataset = 20, n_obs = 400,
        beta0=cur_beta0, beta1=cur_beta1, gamma=gamma, mu=cur_mu, sigma=sigma)
    })
    capture.output(print(
      paste("Simulate from OU with beta0=", cur_beta0, "beta1=", cur_beta1, "tau=",
        cur_tau, "mu=", cur_mu)), file=output_file, append=TRUE)
    capture.output(print(cbind(test_cases, t(ou_result))), file=output_file, append=TRUE)
  }
}

sim2_OU_script()

```

#### 4. Code to generate results for Challenge 2 (simulation from BM)

##### (1) Input file

```

sim2_BM_cases <- read.csv('sim2_BM_cases.csv')
sim2_BM_cases

```

```

##   X beta0 beta1 mu gamma      method
## 1 1    10  0.5  0     5 Nelder-Mead
## 2 2     5  0.5  0     5 Nelder-Mead
## 3 3     5  1.0  0     5 Nelder-Mead
## 4 4    10  0.5  0     5      BFGS
## 5 5     5  0.5  0     5      BFGS
## 6 6     5  1.0  0     5      BFGS

```

##### (2) Script

```

require("smfret")

all_cases_default <- read.csv("sim2_BM_cases.csv", row.names=1)
test_cases_default <- data.frame(sigma=c(0.000001, 0.00001, 0.0001, 0.001, 0.01, 0.1))

output_file_default <- "sim2_bm_results.RData"

#' Generate observations from a Morse SDE.

```

```

#'
#' @param all_cases dataframe of cases to test. Each row is a test case.(see *Details*)
#' @param test_cases dataframe of sigmas to test
#' @param output_file the name of the output file
#' @return
#' @details
#' 1. parameter `all_cases` must has column
#'    `c("beta0","beta1","gamma","n_dataset","n_obs","num_multistart")`.
#' 2. parameter `test_cases` must have column "sigma"
sim_2_BM_script <- function(all_cases=all_cases_default, test_cases=test_cases_default,
                             output_file=output_file_default) {
  for(i in 1:nrow(all_cases)) {
    row <- all_cases[i,]
    cur_beta0 <- row[['beta0']]
    cur_beta1 <- row[['beta1']]
    cur_mu <- row[['mu']]
    cur_gamma <- row[['gamma']]
    cur_method <- toString(row[['method']])
    print(paste("cur_method=", cur_method))

    bm_result <- apply(test_cases, 1, function(info) {
      sim_2(from="bm", sigma=info[["sigma"]], n_dataset = 20, n_obs = 400,
            beta0=cur_beta0, beta1=cur_beta1, mu=cur_mu, gamma=cur_gamma, method=cur_method)
    })
    capture.output(
      print(paste("Simulate from BM with beta0=", cur_beta0, "beta1=", cur_beta1, "mu=",
                  cur_mu, "gamma=", cur_gamma, "with method=", cur_method)), file=output_file, append=TRUE)
    capture.output(cbind(test_cases, t(bm_result)), file=output_file, append=TRUE)
  }
}

sim_2_BM_script()

```

## References

- G. E., Uhlenbeck, and L. S. Ornstein. 1930. “On the Theory of Brownian motion.” *Physical Review* 36 (1): 823–41.
- Ha, Taekjip. 2001. “Single-Molecule Fluorescence Resonance Energy Transfer.” *Methods* 25 (1): 78–86.
- Hirota, Matsumoto, and Hisanori Kiryu. 2016. “SCOUP: A Probabilistic Model Based on the Ornstein-Uhlenbeck Process to Analyze Single-Cell Expression Data During Differentiation.”
- Ideker, Trey, and Roded Sharan. 2008. “Protein Networks in Disease.” *Genome Research* 18 (4): 644–52.
- Lemke, EA, AA Deniz, and RJ Groarke. 2016. “Förster Resonance Energy Transfer.”
- Libretexts. 2019. “Lennard-Jones Potential.” *Chemistry LibreTexts*. Libretexts. [https://chem.libretexts.org/Bookshelves/Physical\\_and\\_Theoretical\\_Chemistry\\_Textbook\\_Maps/Supplemental\\_Modules\\_\(Physical\\_and\\_Theoretical\\_Chemistry\)/Physical\\_Properties\\_of\\_Matter/Atomic\\_and\\_Molecular\\_Properties/Intermolecular\\_Forces/Specific\\_Interactions/Lennard-Jones\\_Potential](https://chem.libretexts.org/Bookshelves/Physical_and_Theoretical_Chemistry_Textbook_Maps/Supplemental_Modules_(Physical_and_Theoretical_Chemistry)/Physical_Properties_of_Matter/Atomic_and_Molecular_Properties/Intermolecular_Forces/Specific_Interactions/Lennard-Jones_Potential).
- Sasmal, Dibyendu K, Laura E Pulido, Shan Kasal, and Jun Huang. 2016. “Single-Molecule Fluorescence Resonance Energy Transfer in Molecular Biology.” *Nanoscale* 8 (48): 19928–44.
- Schelstraete, Sigurd, Wijnand Schepens, and Henri Verschelde. 1999. “Energy Minimization by Smoothing Techniques: A Survey.” *Molecular Dynamics: From Classical to Quantum Methods*, 129–85.
- Wallace, Bram, and Paul J Atzberger. 2017. “Förster Resonance Energy Transfer: Role of Diffusion of Fluorophore Orientation and Separation in Observed Shifts of FRET Efficiency.” *PloS One* 12 (5): e0177122.
- Yang, Shilong, James B Witkoskie, and Jianshu Cao. 2002. “Single-Molecule Dynamics of Semiflexible Gaussian Chains.” *The Journal of Chemical Physics* 117 (24): 11010–23.