

# S610-Final\_Project-sSINDy

*ChunHsien Lu*

*12/11/2019*

## Introduction

In this project, we mainly follow the paper written by Boninsegna et. al. This paper focuses on the reconstruction of an SDE. It does well on the example they provide. In order to understand the ideas provide in the paper, here, we attempt to apply the same method to Brownian motions(BMs).

Suppose  $X_t$  is an Ito's stochastic process with constant diffusion, i.e.

$$dX_t = b(X_t)dt + \sigma dB_t$$

where  $b$  and  $\sigma$  are functions,  $B_t$  is a BM. The following is two important properties that can help us to reconstruct  $b$  and  $\sigma$ :

$$b(x) = \lim_{s \rightarrow 0} \frac{1}{s} E[X_{t+s} - X_t | X_t = x]$$

and

$$\sigma(x) = \lim_{s \rightarrow 0} \frac{1}{s} E[(X_{t+s} - X_t)(X_{t+s} - X_t) | X_t = x].$$

Since what we want to do is BM,  $b(x) = 0$  and  $\sigma(x) = 1$ . From the two properties, fixed some small  $h > 0$ , we try to fit  $b(x)$  by using  $[X_{(m+1)h} - X_{mh}]_m$  and  $[X_{mh}]_m$ . Furthermore,  $\sigma(x)$  is via  $[(X_{(m+1)h} - X_{mh})(X_{(m+1)h} - X_{mh})]_m$  and  $[X_{mh}]_m$ .

Intuitively, the algorithm proposed in the paper is just by using backward elimination and cross validation. The design matrix is  $[M]_{ij} = [f_j(X_{ih})]_{ij}$  where  $f_j$ 's are possible candidates in  $b(x)$ . Assume  $n$  functions  $f_j$  are chosen. The full model is fitted as

$$(X_{(i+1)h} - X_{ih}) \sim \sum_{j=1}^n \beta_j f_j(X_{ih}).$$

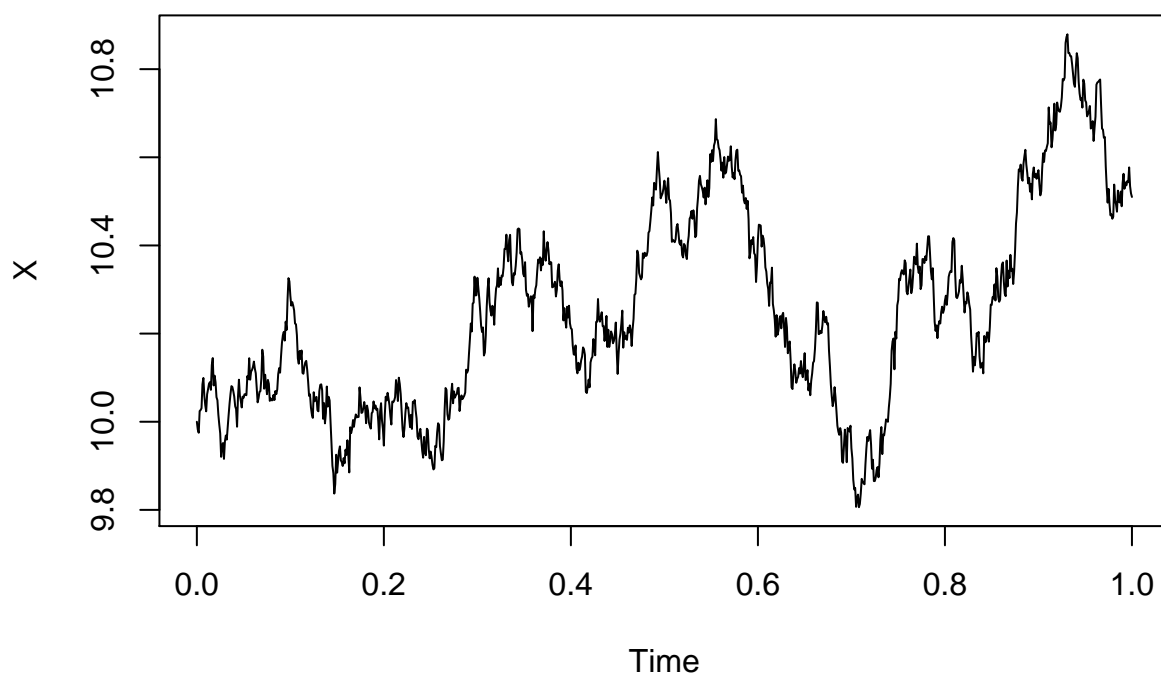
by using all candidates. The smaller model with  $k - 1$  predicted variables is removed the one with lowest absolute value in the model with  $k$  variables. Via this way, we can have  $n + 1$  possible models. ( $y = 0$  is included) Now, CV comes into play to decide the best model.

The codes can be found in the following link:  
[https://github.com/Pielann/sto\\_SINDy.git](https://github.com/Pielann/sto_SINDy.git)

## Simulation Design

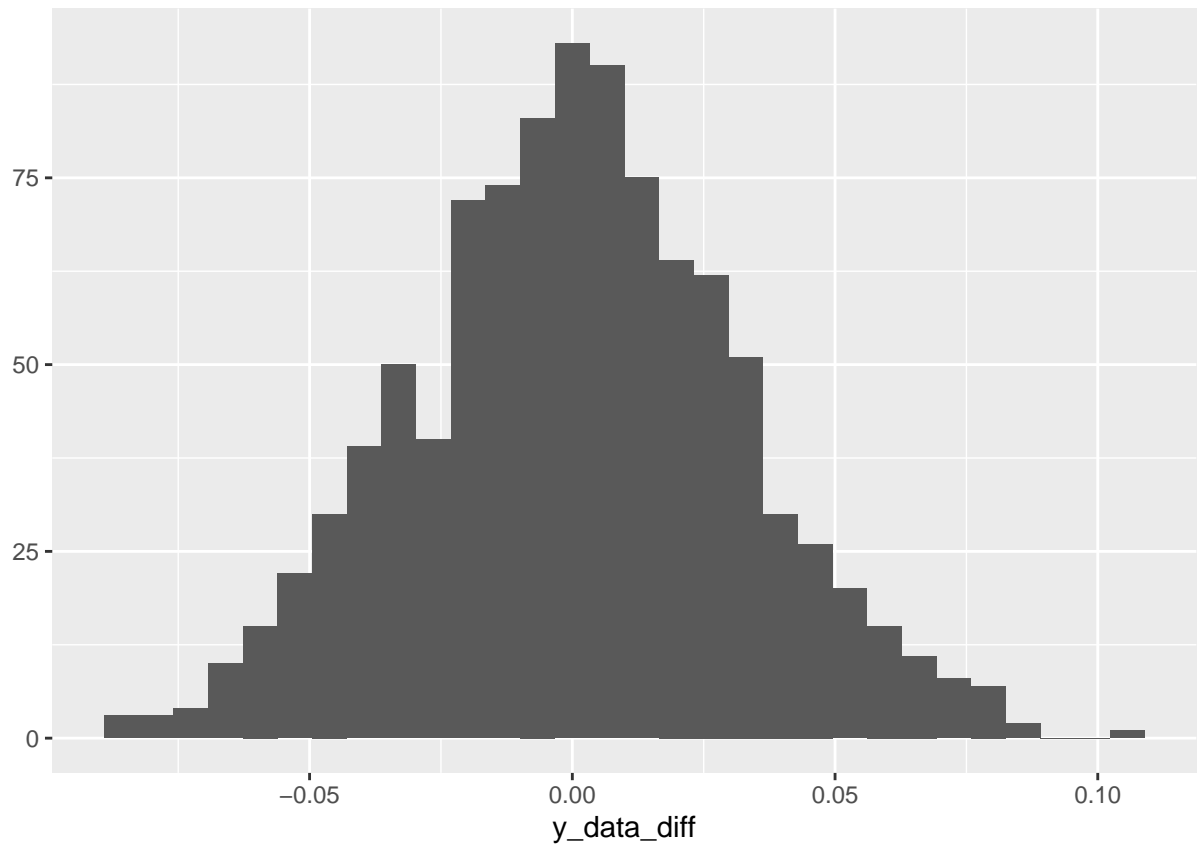
Fit the drift term and explain the output of sSINDy:

**A path of a BM start at 10**



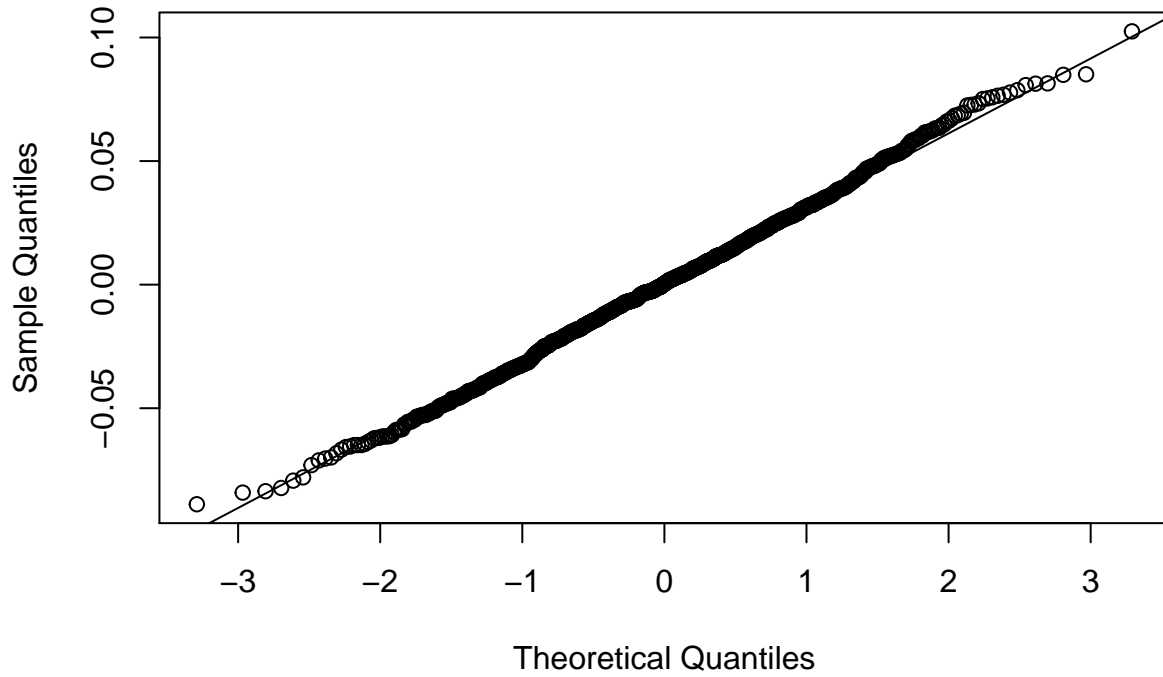
```
# Some facts for y_data  
#This plot should be similar to normal  $N(0, var=t)$   
ggplot2::qplot(y_data_diff)
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
qqnorm(y_data_diff)  
qqline(y_data_diff)
```

## Normal Q-Q Plot



```
cat(" mean:",mean(y_data_diff),"\\n",
    "variance:",var(y_data_diff))
```

```
## mean: 0.0005100079
## variance: 0.0009834589
```

Since  $X_t$  is a BM, then  $[X_{(m+1)h} - X_{mh}]_m$  should follow  $N(0, h)$ . We choose  $h = 0.001$  here. As shown above, the variance is 0.0009835, which is really closed to  $h$ .

```
## Warning: package 'magrittr' was built under R version 3.5.3
```

```
## $coef_order
## [1] 4 3 2 1
##
## $cv_error
## [1] 8.844620 9.102773 9.097858 9.202561 9.332238
##
## $'min#'
## [1] 0
##
## $coef
## [1] 0 0 0 0
```

We did it by 10-fold cross-validation(CV) to compare the models which are up to 3rd order. More precisely, we can compare the models

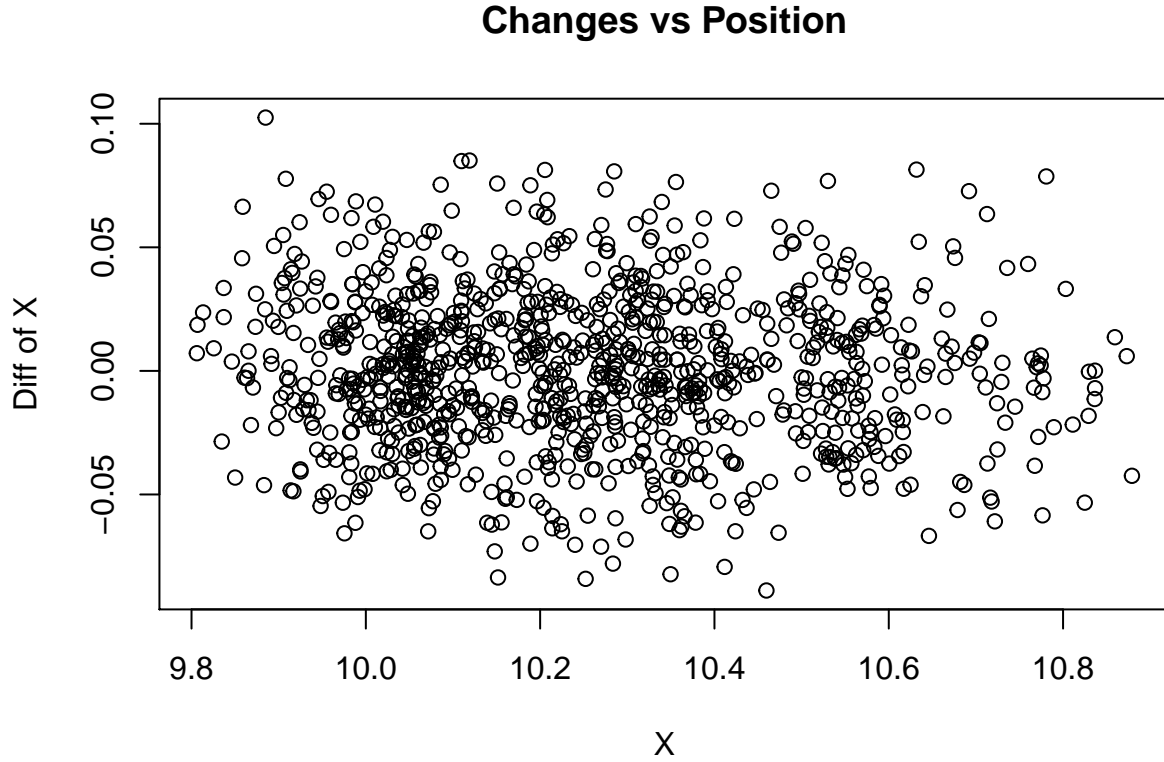
$$y = 0, \quad y = \beta_{i_1}x^{i_1}, \quad y = \beta_{i_1}x^{i_1} + \beta_{i_2}x^{i_2}, \quad y = \beta_{i_1}x^{i_1} + \beta_{i_2}x^{i_2} + \beta_{i_3}x^{i_3}, \quad y = \beta_{i_1}x^{i_1} + \beta_{i_2}x^{i_2} + \beta_{i_3}x^{i_3} + \beta_{i_4}x^{i_4},$$

where  $i_1, i_2, i_3, i_4 \in \{0, 1, 2, 3\}$  are distinct. The output “coef\_order” express the order of the coefficients to be removed. 1 means the first variable to be removed. Hence, the model with three variables is

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

since the fourth term is the first one to be thrown. “cv\_error” denotes the total errors via CV for models with variables from 0 to 3. “min#” means the model with least errors. Moreover, the “coef” is the coefficients from the best choice. Surprisingly, the true function for drift term is 0. It shows that, in this example, we do really rebuild the drift term.

**Compute the diffusion term:**



Via above plot, it tells us that the variance with respect to  $X$  can help us compute  $\sigma(x)$ . We roughly observe the variances don't vary as  $X$  changes. Moreover, since we just have a sample path and the fitting of  $b(x)$  is constant zero, the true variance will be equal to the standard deviation of

$$\left[ \frac{X_{(m+1)h} - X_{mh}}{\sqrt{h}} \right]_m.$$

Below computes the diffusion term from the sample path which is 0.9917. It is also really closed to 1. Hence, the diffusion is rebuilt.

**## Diffusion Term: 0.991695**

## Confidence Intervals

### CI for Drift Terms

```
##          [,1] [,2] [,3] [,4]
## 2.5%         0    0    0    0
## 97.5%         0    0    0    0
```

This simulation shows that the 95% CI are 0 for 1000 samples which has really good performance.

### CI for Diffusion Terms

```
##          2.5%          97.5%
## 0.9840602 1.0120430
```

The true diffusion term for standard BM is 1. The CI contains the true value.

## References

Boninsegna et. al, Sparse learning of stochastic dynamical equations, J. Chem. Phys. 148, 241723 (2018).