

# Tutorial for MATLAB reconstruction package

## Contents

1. The package in brief
2. Different data types and proper data formats
3. Two different modeling approaches
  - 3.1 Parametric models
  - 3.2 Spline models
4. Two optimization solvers
5. A short description about the inputs and outputs of the main codes
6. Check three requirements in real datasets in advance
  - 6.1 Data stationarity
  - 6.2 Data Markovicity
  - 6.3 Data resolution: a key data feature in this study
7. Simulating data from parametric and spline models
8. Data standardization
9. High-resolution data: Euler reconstruction
  - 9.1 Reconstructing a dataset simulated from a linear model
  - 9.2 Reconstructing a dataset simulated from a nonlinear model
  - 9.3 Reconstructing an ecological dataset
10. Sparsely sampled data: Hermite reconstruction
  - 10.1 When does Hermite reconstruction crash? Strategies and recommendations
  - 10.2 Reconstructing a simulated from a linear model
  - 10.3 Reconstructing a dataset simulated from a nonlinear model
  - 10.4 Reconstructing an ice-core dataset
  - 10.5 The concept of effective potential
11. Handling big datasets
12. Handling replicate datasets
13. Assessing the uncertainty of the results

## 1. The package in brief

This package implements a maximum likelihood estimation (MLE) inference in order to fit Langevin models to univariate time series data. The process of fitting a stochastic differential equation, including Langevin models, to data is commonly referred to as ‘system reconstruction’ in the literature (Siegert and Friedrich 2001, Rinn et al. 2016). Based on the correlational structure of the data, the package classifies datasets into three categories: ‘high-resolution’, ‘low-resolution’, and ‘bad-resolution’. If data resolution belongs to bad-resolution category, reconstruction is unlikely to be successful, and it is recommended to use a dataset with a higher resolution. For high-resolution data, the Euler inference technique, which we refer to as ‘Euler reconstruction’, can be implemented while for low-resolution data a more advanced but computationally more intensive inference algorithm, which we refer to as ‘Hermite’ reconstruction, can be recommended. Hermite reconstruction is based on a refinement (Bakshi and Ju 2005) to a reconstruction approach developed by Aït-Sahalia (Aït-Sahalia 2002). For real datasets, it is recommended to follow Hermite reconstruction, even for high-resolution data, unless if the data resolution is exceptionally high.

The package supports two different modeling strategies: ‘parametric models’ and ‘spline models’ (simply put a spline function is constructed by joining piecewise polynomials smoothly to form a flexible and complex functional form). Spline modeling, offers a particularly appealing framework when the appropriate functional form of the underlying system is not straightforward to speculate. Splines offer a flexible structure that can capture nonlinearities inherent in the data-generating system. Furthermore, since splines are linear functions in terms of model parameters, their use often leads to faster and more accurate results.

The can support both ‘typical’ datasets (a single time series) and ‘replicate’ datasets (multiple time series all believed to belong to the same data-generating system). Additionally, the package can analyze a well-mixed fraction of an extremely large typical or replicate dataset sampled randomly across the entire data. The package can also reconstruct datasets consisting missing values. Once model parameters are estimated, the package can assess the corresponding uncertainty of the estimated parameters. We encourage the users to develop an understanding of the ideas and techniques by actively running the codes. That is to say, learning by doing.

The package is compatible with MATLAB 2022 and requires the following toolboxes: Curve Fitting Toolbox, Optimization Toolbox, Symbolic Math Toolbox, Econometrics Toolbox, and Signal Processing Toolbox. Additionally, it relies on the MATLAB package ‘ARMASA’ (Broersen, 2003), available for free download from the following link

<https://nl.mathworks.com/matlabcentral/fileexchange/1330-armasa>

Moreover, the code ‘armasel\_s.m’ from the reference (Erkelens et al., 2013) is necessary. The authors have kindly permitted the inclusion of their code in our package. To ensure proper functionality, we have compiled both the ARMASA package and the ‘armasel\_s.m’ code into a folder named ‘Burg’. Please add the path of this folder to your MATLAB working path.

## 2. Different data types and proper data formats

3. In this package, we distinguish between two types of data: 'typical' and 'replicate' data. Typical data refers to a single uninterrupted time series dataset, whereas replicate data consists of multiple separate time series datasets that are generated from the same underlying data generating system. Typical data can be supplied as an array. On the other hand, replicate data should be supplied as cell arrays. Each cell should contain a single replicate dataset, following the same format as typical data. In other words, a replicate dataset should be organized as a cell array, with each cell containing a typical dataset. Both typical and replicate datasets can include 'missing values', which should be specified using 'NaN' notation. The package in brief

## 4. Two different modeling approaches

### 4.1 Parametric models

In this modeling framework, which we term '*parametric reconstruction*', one specifies a Langevin model along with a parameter vector  $\theta$ . A parametric Langevin model is represented by the stochastic differential equation:

$$dx = \mu(x; \theta)dt + \sigma(x; \theta)dW, \quad (1)$$

where  $\mu(x; \theta)$  denotes the deterministic component of the system, known as the '*drift*' function, and  $\sigma(x; \theta)$  represents the stochastic component, known as the '*diffusion*' function. '*W*' refers to a Wiener process, making the noise source  $dW$  Gaussian distributed and white (uncorrelated). It is worth noting that in some literature focusing on the equivalent Fokker-Planck formulation of Langevin model (*I*), the function  $\frac{1}{2} \sigma^2(x; \theta)$  is referred to as the diffusion function. The role of  $\sigma(x; \theta)$  is to weigh the impact of noise source per state  $x$  measuring the noise intensity.

A Langevin model (*I*) is termed '*additive*' if the diffusion function  $\sigma(x; \theta)$  is constant, meaning it does not vary with the state variable  $x$  (although it may depend on parameters  $\theta$ ). Otherwise, the Langevin model is termed '*multiplicative*'.

### 4.2 Spline models

Spline models are also parametric but they have a flexible form that can adapt to the shape of many nonlinear functions. Splines are an accurate tool for univariate data and fast to compute. To distinguish this modeling approach from parametric reconstruction we call it '*spline reconstruction*'.

Spline models are also parametric, but they offer a flexible form that can adjust to the shape of many nonlinear functions. In this modeling approach, the model parameters correspond to the values of the drift and diffusion functions over a relatively coarse mesh of the state space, known as the '*knot sequence*'. *What makes it convenient to work with splines is the fact that splines are linear functions in terms of parameters, even though they are non-linear in terms of state variables.* Unlike in parametric reconstruction, the user does not need to specify a model him/herself. Splines are particularly useful for univariate data and are computationally efficient. This package only considers spline modeling for univariate data. To differentiate this modeling approach from parametric reconstruction, we refer to it as '*spline reconstruction*'.

## 5. Two optimization solvers

We use two different optimization solvers in this package. The first solver is '*fmincon*' which is a built-in MATLAB solver. The second solver is the 'Grey wolf optimizer' (GWO) (Mirjalili et al. 2014), actually an improved GWO (Nadimi-Shahraki et al. 2021), abbreviated as '*gwo*', in this package. *fmincon* is a local solver but is a fast solver. *gwo*, on the other hand, is a global solver and is slower. We also utilize the MultiStart option to turn *fmincon* into a global solver. *fmincon* is our default solver.

## 6. A short description about the inputs and outputs of the main codes

The main code in the package is called 'euler\_reconstruction.m'. Below, we explain its inputs and outputs here. Have a look here but the best way to learn is to run several examples after this section.

```
Res = euler_reconstruction(data,dt,'name',value,...)
```

**data:** Vector with a fixed time step

**dt:** The fixed time step between consecutive data points

### *name-value pairs*

**'lb':** Vector with the lower bounds of all parameters. For spline models the default is -10 for all knot values of mu and 0 for all knot values of sigma.

**'ub':** Vector with the lower bounds of all parameters for spline models the default is 10 for all the knot values of mu and sigma)

**'L':** Left boundary of the data (default is min(data))

**'R':** Right boundary for the data (default is max(data))

**Note:** When working with small datasets, there may be fewer data points near the borders. This can adversely affect the quality of the fitted model, especially when using spline models (parametric models are not affected). To mitigate this issue, it is recommended to choose a relatively larger lower boundary and a smaller upper boundary for your data.

**'solver':** Optimization solver for the maximum-likelihood estimation problem. The solvers are 'fmincon' and 'gwo' (default is fmincon). 'gwo' is a global solver but is slower.

**'gradient\_fun':** The gradient vector of the objective function. All the optimization solvers can work without this but we recommend to use this option whenever applicable.

**Note:** This option is applicable only to parametric models (but is not applicable to 'Hermite reconstruction'). For nonlinear parametric models, it is recommended to utilize this option. This approach helps prevent the solver from becoming stuck at points that are not even local minima (stagnation).

**'useparallel':** Use parallel computing (default is false)

**'search\_agents':** Number of searching agents (default is 5)

**'maxiter':** Maximum number of iterations (default is 'realmax' which, in practice, means infinity)

**'nknots':** Is a two-element vector where the first (second) element specifies the number of knots you want to allocate for mu (sigma). For additive noise use [n 1] which means you use n knots for mu and a single knot for sigma (default is [8 8]).

**'knots':** The values of the knots (alternative to nknots).

**'spline':** a two-element string which specifies the types of splines for mu and sigma (the default is 'CC'). Spline types for mu and sigma are as follow

**'L'** = linear interpolation (i.e., a straight line)

**'C'** = cubic spline interpolation

**'Q'** = quadratic spline interpolation

**'P'** = pchip spline interpolation (pchip respects the monotonicity in data)

**'SCS'** = cubic smoothing spline

170 'Approximate' = 'SCS' but uses 'L' for fast fitting  
 171  
 172 (For instance, 'LL' means that you want to specify linear interpolation 'L' for both mu and sigma. 'CL' means that  
 173 you want to specify cubic spline 'C' for mu but a linear spline 'L' for sigma. Likewise, 'SCSL' means 'SCS' for mu  
 174 and 'L' for sigma).

175  
 176 *The following name-value pairs are only suitable for parametric models*

177 'mu': parametric function handle for mu.

178  
 179 'sigma': parametric function handle for sigma (if it is empty then an additive sigma (i.e., constant) is considered)  
 180

181 **Note:** mu and sigma must be vectorized. For instance, `mu = @(x,par)par(1)*x^3+par(2)*x` is not suitable,  
 182 rather `mu = @(x,par)par(1).*x.^3+par(2).*x` is appropriate.

183  
 184 'npars': (optional) number of parameters.  
 185

186 'gradient\_fun': handle to gradient function (generated with 'eulergrad'). If you want the code to calculate the  
 187 gradient use this option (then you might get a more accurate result).  
 188

### 189 ***The optimizing dialog screen***

190 While the package is solving the MLE problem, a dialog screen appears, providing updates on the optimization  
 191 progress at each iteration. This feature is particularly valuable when applying Hermite reconstruction, allowing users  
 192 to monitor the outcomes closely. If the results show gradual improvements, users can safely stop the code by pressing  
 193 the 'stop' button.

### 194 ***A simple plotting option***

195 After you are done with the command `'res=euler_reconstruction(data,dt,'name',value,... )'` you  
 196 can enjoy a nice graphical interface using the command `'plot_results(res)'`.

## 197 **7. Check three requirements in real datasets in advance**

### 198 **7.1 Data stationarity**

199 In order to fit a Langevin model ( $I$ ) to data, it is essential for the data to be stationary, at least in a weak sense. In  
 200 simpler terms, stationarity implies that the statistical properties of the system, and hence the dataset under  
 201 consideration, remain invariant over time. Weak stationarity within a fixed time window entails that the mean and  
 202 variance of the data remain constant, and that the autocorrelation function depends solely on the time lag rather than  
 203 on the initial and final times within the specified window. However, if stationarity is violated across the entire dataset,  
 204 it may be possible to identify shorter, (possibly overlapping) time windows during which the data exhibit stationarity  
 205 (see **Example 10** in subsection 10.3). Reconstruction can then be performed separately for each of these windows.  
 206 The stationarity of the data can be assessed using the Augmented Dickey-Fuller test (ADF test) (Dickey and Fuller  
 207 1979). We conducted the ADF test on all real datasets in this tutorial. It is worth noting that simulated data are  
 208 inherently stationary because model ( $I$ ) is stationary.

209 As an example, let's apply the ADF test to the 'OUdata1D.mat' dataset as below

```
210 S = load('OUdata1D.mat');
211 data = S.data;
212 [~,~,~,~,reg] = adftest(data,'model','ARD','lags',0:20); % the input 0:20 is the
213 number of lags we try in fitting an autoregressive model to data.
214 [~,lagndx] = min([reg(:).BIC]); % this tells us how many lags we need (lagndx is 1)
215 [h, Pvalue,~]=adftest(data,'model','ARD','lags',lagndx);h
```

216 1

217 In this test we apply ADF test at least twice. First, we consider an array of lags and apply ADF test to see how many  
 218 autoregressive lags our data needs. In the above example we considered the lags 0:20 and the quantity `lagndx` tells  
 219 us the required number of lags (which corresponds with the lowest BIC) we need which is 1 (if we get 20 then clearly,  
 220 we should repeat the command `adftest(data, 'model', 'ARD', 'lags', 0:20)` with a bigger array of lags).  
 221 Next, we apply the ADF test again and the output `h` is the test result. If `h` = 1 (as is the case above and this was  
 222 expected) then the dataset is stationary, otherwise it is non-stationary. As for another example with real data type the  
 223 following

224 We should apply the ADF test at least twice. Firstly, we consider an array of lags and apply the ADF test to determine  
 225 the optimal number of autoregressive lags required for our data. In this example, we considered lags ranging from 0  
 226 to 20, and the quantity `lagndx` indicates the number of lags corresponding to the lowest Bayesian Information  
 227 Criterion (BIC). For instance, if `lagndx` returns 1, it suggests that only one lag is required. However, if `lagndx` returns  
 228 20, it indicates the need to repeat the ADF test with a larger array of lags. Subsequently, we perform the ADF test  
 229 again, and the output `h` represents the test result. A value of `h` = 1 indicates that the dataset is stationary, as observed  
 230 in the example above. Conversely, if `h` is not equal to 1, it signifies that the dataset is non-stationary. For another  
 231 example, using real data, we try the following command lines

```
232 data = readmatrix('NGRIP20.csv');
233 [~,~,~,~,reg] = adftest(data,'model','ARD','lags',0:20); % the input 0:20 is the
234 number of lags we try in fitting an autoregressive model to data.
235 [~,lagndx] = min([reg(:).BIC]); % this tells us how many lags we need (lagndx is 5)
236 [h, Pvalue,~]=adftest(data,'model','ARD','lags',lagndx);h
237 1
```

238 This confirms the stationarity of the second dataset. All datasets in this tutorial are stationary.

## 239 7.2 Data Markovicity

240 First of all, the MATLAB package called 'ARMASA' (Broersen 2003) is needed for this section. This package can be  
 241 freely downloaded from the following link

242 <https://nl.mathworks.com/matlabcentral/fileexchange/1330-armasa>

243 Reconstructing real datasets presents unique challenges, notably due to the correlation of noise at very small scales, a  
 244 phenomenon highlighted by Einstein in his seminal work on Brownian motion (Einstein 1905). Langevin models in  
 245 (*I*), on the other hand, are Markov models (In short, the Markov property dictates that the future state of a system,  
 246 given its present state, is independent of the entire history of past states). Consequently, the reconstruction process  
 247 must adhere to a coarser time scale, known as the '*Markov-Einstein*' (*ME*) time scale (Friedrich et al. 2011), ensuring  
 248 the fulfillment of the Markov property. This implies that if the ME time scale equals 1, the entire dataset is Markov  
 249 and can be used directly. However, for an ME time scale of 2, a sample of original data with every second data point  
 250 (e.g., `data(1:2:end)`) should be considered in the analysis, and so forth. Data simulated from the Langevin model  
 251 (*I*), are Markov (i.e., ME time scale is 1) by construction. For real datasets ME time scale should be calculated.

252 It is important to note that the Markov property holds at any time scale larger than the ME time scale, allowing for  
 253 reconstructions at these higher scales as well. Determining the ME time scale, however, is far from straightforward.  
 254 Traditional methods for estimating the ME time scale often involve binning data and this requires extensive data  
 255 (Friedrich et al. 2011), leading to results that may be sensitive to the chosen bin size and potentially introduce bias in  
 256 smaller datasets. To address these challenges, we propose a more streamlined and data-efficient method. This approach  
 257 involves fitting an autoregressive (AR) model to the data and examining the order of the fitted AR model. Specifically,  
 258 if an AR(1) model emerges as the optimal fit, this strongly suggests an ME time scale of 1, indicating that the original  
 259 dataset is Markov. Similarly, an AR(*p*) model suggests an ME time scale of *p* meaning that a sample with every *p*-th  
 260 data point (e.g., `data(1:p:end)`) is Markov and this sample should be considered for reconstruction not the original  
 261 dataset. Here, we assess the ME time scale for a few datasets. As first example, type the following commands

```
262 S = load('OUdata1D.mat');
```

```

263 data = S.data;
264 order = 10; % order is the maximum AR order being considered (should be long enough)
265 AR = ME_TimeScale(data,order)

```

266 and you get

```

267 AR = [1.0000 -0.9899]

```

268 This indicates that an AR(1) model provides the best fit for this dataset (note that the number of elements in the  
269 estimated AR vector after the first element, which is always 1, corresponds to the order of the fitted AR process.)  
270 Therefore, the ME time scale is 1, confirming the Markovian nature of this dataset. This was expected, since this  
271 dataset is simulated. Let's now examine a real dataset

```

272
273 data = readmatrix('BGA_stdlevel_2011.csv');
274 order = 10;
275 AR = ME_TimeScale(data,order)

```

276 and you get

```

277 AR = [1.0000 -0.6261 -0.2177 -0.0840 -0.0339 -0.0120 -0.0068 -0.0009
278 -0.0080 0.0005 -0.0089]

```

279 This suggests that the dataset exhibits long-range correlations. However, the magnitudes of the AR coefficients beyond  
280 the fourth element (-0.0840) or the fifth element (-0.0339) are small and we can safely assume that the AR order is  $p$   
281 = 4 or even  $p = 3$ . We consider the ME time scale of 4 (note that any time scale bigger than 4 is also Markov and we  
282 could consider it though this comes at the expense of reduced data resolution). This indicates that the dataset is not  
283 Markov but a data ratification by considering every forth data point `data(1:4:end)` is Markov and reconstruction should  
284 be performed on this smaller and coarser sample. As for another example we examine the Markovicity of a climate  
285 dataset below

286

287

```

288 data = readmatrix('NGRIP20.csv');
289 order = 10;
290 AR = ME_TimeScale(data,order)

```

291 and you get

```

292 AR = [1.0000 -0.4879 -0.2198 -0.1231 -0.0541 -0.0030 -0.0186 -0.0290
293 -0.0159 -0.0061 -0.0262]

```

294 The AR order is  $p = 3$  is reasonable. However, since this dataset is very small we considered  $p = 2$  which is acceptable.

### 295 **7.3 Data resolution: a key data feature in this study**

296 Prior to embarking on system reconstruction, it is crucial to estimate the resolution of the data. Based on its resolution,  
297 data is categorized into three categories: 'high-resolution', 'low-resolution', and 'bad-resolution'. This categorization  
298 is partly theoretical and partly empirical (build on our experience on many real and simulated datasets). But, it is useful  
299 in choosing an appropriate reconstruction procedure. To measure data resolution a quantity called '*relaxation time*  
300 *scale*' (Honisch and Friedrich 2011) should be calculated. Use the command `R = RelaxationTime(data)` to get a  
301 rough estimate of the relaxation time  $R$  ( $R$  expresses the relaxation time as a multiple of data time step. For instance,  
302  $R=176$  means that relaxation times is 176 times the data resolution). A dataset with  $R \leq 1$  falls in the category of bad-  
303 resolution,  $1 < R \leq 100$  falls in the category of low-resolution, and  $100 < R$  falls in the category of high-resolution (see  
304 Figure 1 in the main text). The threshold  $R \leq 1$  has a theoretical basis while the threshold  $100 < R$  is based on experience.  
305 Note that this categorization is a rough categorization since we cannot accurately estimate  $R$ . However, it is not  
306 necessary to accurately estimate the time scales; rather, having a general '*feel*' for them is sufficient. In the previous

307 section on data Markovicity we learned that reconstruction of real data should be conducted on a rarified sample of  
308 data that adheres to Markov property. Therefore, it is important to calculate the relaxation time for that specific sample  
309 not the entire data.

310 If data resolution is of bad-resolution, then reconstruction is not possible since measurement time between consecutive  
311 data points is so big so that the dynamics of the true data-generating system is no longer reflected in such a dataset. In  
312 all other categories a reconstruction scheme called 'Euler reconstruction' should be followed first. If data resolution is  
313 so high, then Euler reconstruction is often sufficient. However, if data resolution is low then a 'Hermite reconstruction'  
314 should be performed right after that which is built upon Euler reconstruction. Hermite reconstruction is more accurate  
315 than Euler reconstruction but it is computationally more expensive. Based on our experience, for most real data Hermite  
316 reconstruction often leads to a considerable improvement of Euler reconstruction even if data resolution is high (only  
317 when resolution is extremely high the two reconstruction procedures have close performance). Therefore, it is  
318 recommended to always follow Hermite reconstruction for real data. We elaborate on Euler and Hermite  
319 reconstructions in section 9 and section 10, respectively. As an example, let's try finding the relaxation time for the  
320 dataset in **Example 1** in Section 9. Type the following commands

321 **8. The package in brief**

322 **9. The package in brief**

323 **10. The package in brief**

324 **11. The package in brief**

```
325 S = load('OUdata1D.mat');  
326 data = S.data;  
327 R = RelaxationTime(data); R  
328 98.1858 (number of time steps)
```

329 As we expected (since this dataset was generated from the OU model  $dx = -x dt + dW$ , whose relaxation time is 1.  
330 The dataset has a simulation time step  $dt$  of 0.01. So, in theory, R should be 100 time units. However, due to the finite  
331 size of the data, we obtained  $R \sim 98.18$ ). Now, we assess the ME time scale for an ecological dataset. Type the following  
332 commands

```
333 data = readmatrix('BGA_stdlevel_2011.csv');  
334 data = data(1:4:end); % Important: We learned from previous section that ME time  
335 scale in this dataset is 3. So, we must consider every third data point  
336 R = RelaxationTime(data); R  
337 219.2118 (number of time steps)
```

338 It is important to pay attention to the fact that this real dataset is not Markov. However, a rarified sample of this dataset,  
339 obtained by considering every forth data point, i.e., `data(1:4:end)`, is Markov. Therefore, when estimating the  
340 relaxation time, it is essential to analyze this rarified sample rather than the original dataset. Despite being highly  
341 rarified, this sample still maintains a high resolution. Applying the Euler reconstruction to this dataset gives relatively  
342 accurate results. However, since this dataset is a real dataset we perform a Hermite reconstruction on this dataset in  
343 later sections as well. As for another real dataset we examine the ME time scale of an ice-core dataset. Type the  
344 following

```
345 data = readmatrix('NGRIP20.csv');  
346 data = data(1:2:end); % Important: We learned from previous section that ME time  
347 scale in this dataset is 2. So, we must consider every other data points  
348 R = RelaxationTime(data); R  
349 63.1028 (number of time steps)
```



350 So, the dataset has a low-resolution and most probably Euler reconstruction is not operational for this dataset. Finally,  
351 we examine the relaxation time of a replicate dataset as below

```
352 S = load('MayData1D_Replicate.mat');
353 data = S.data;
354 R = RelaxationTime(data); R
```

355 219.5342 (number of time steps)

356 And this puts this rarified sample in the category of low-resolution. The computational burden in the second phase of  
357 reconstruction (called Hermite reconstruction, see section 10) increases as we increase either J or K (these are  
358 parameters of the Hermite reconstruction). A small relaxation time indicates the need for choosing large values of K  
359 (for a fixed J) for the estimation procedure to work efficiently. Conversely, a large relaxation time means that highly  
360 accurate result can be obtained with small K (for a fixed J). Note that here ‘small’ relaxation time refers to small values,  
361 typically close to 1 from above, while ‘large’ indicates values significantly larger than 1. **Table 1** summarizes the data  
362 requirements for all the datasets in this tutorial.

Datasets	Real or simulated?	ME time scale	Relaxation time scale (s)	Category
OUdata1D.mat	Simulated	1	98.18	(almost) high-resolution
MayData1D.mat	Simulated	1	3641.8	high-resolution
MayData1D_Replicate.mat	Simulated	1	219.5342	high-resolution
BGA_stdlevel_2011.csv	Real (ecology)	3 or 4	219.2118	high-resolution
OUdata1D.mat (Every 100 <sup>th</sup> data points)	Simulated	1	1.0074	Extremely low-resolution
MayData1D.mat (Every 300 <sup>th</sup> data points)	Simulated	1	12.79	low-resolution
NGRIP20.csv (We analyzed every other data point)	Real (ice-core)	2 or 3	63.1028	low-resolution

363 **Table 1.** A summary of data requirements for all the datasets used in this tutorial.

## 364 12.Simulating data from parametric and spline models

365 The command for the simulations is 'simulate'. In order to generate a dataset from a parametric model, you need  
366 to specify the parameters: the model type `ModelType` (i.e., 'parametric'), a lower bound L for data (if empty, the  
367 code considers it to be  $-\infty$ ), an upper bound R for data (if empty, the code considers it to be  $\infty$ ), the drift function `mu`  
368 (which should be a function handle), the diffusion function `sigma` (which should also be a function handle), the initial  
369 state `x0`, time step `dt` (which should be small, relative to the scale of the problem), and the number of data points T.  
370 As a first example, consider the one-dimensional Ornstein-Uhlenbeck (OU) model  $dx = \mu x dt + \sigma dW$  with  
371 parameters  $\mu = -1$  and  $\sigma = 1$ . The following command lines generate a dataset from OU model, starting from  $x_0 =$   
372 0, with the time step of  $dt = 0.01$  and  $T = 10^5$  data points (see Figure2, left panel, for an illustration)

```
373 ModelType = 'parametric';
374 L = []; % lower bound on data is -infinity
375 R = []; % upper bound on data is infinity
376 a = -1; s = 1;
377 mu = @(x)a.*x; % drift function
378 sigma = @(x)s+0.*x; % diffusion function
379 dt = 0.01; % time step of Euler-Maruyama integration
380 x0 = 0; % initial state
```

```

381 T=10^5;    % simulation length
382 x=simulate(ModelType,L,R,mu,sigma,dt,x0,T);

383 as you see the lower bound L and the upper bound R are empty. Therefore, the package, by default, considers a lower
384 bound of  $-\infty$  and an upper bound of  $\infty$ .

385 Finally, let's simulate a dataset from a spline model. Note that in this package spline modeling is only possible for one-
386 dimensional models. To generate a dataset from a spline model, you need to define the parameters: model type
387 ModelType (i.e., 'spline'), SplineType (with 'CC' being the default and often used), knots (a rather sparse
388 mesh across the state space), a vector of parameters (which are the corresponding values of drift and diffusion functions
389 at knots), a lower bound L, an upper bound R, a time step dt, and the number of data points T. For details on these
390 inputs see section 5. The following command lines generate data for a spline model (which is reconstructed via a spline
391 reconstruction in section 9.2, Example 4)

392 ModelType = 'spline';
393 SplineType = 'CC';
394 L = -4;
395 R = 4;
396 knots = linspace(L, R, 8);
397 par =[4.8855    1.1983    1.1006    0.24633   -0.10136   -0.69903   -1.3182   -6.4891 ...
398 0.84428    1.0242    1.0017    0.98554    1.0031    0.97949    1.0185    0.79658]; % this
399 %is a parameter vector estimated following a spline reconstruction in subsection 9.1, Example 3
400 dt = 0.01; % time step of Euler-Maruyama integration
401 x0 = 0; % initial state
402 T=10^5;    % simulation length
403 x=simulate(ModelType, SplineType, par, L, R, knots, dt, x0, T);

```

### 404 13.Data standardization

405 To help the reconstruction procedure, it is better to standardize the data by subtracting the mean and dividing by the  
406 standard deviation. This step is especially important when the range of data is large. Standardization helps confine the  
407 search region into smaller and more manageable searching spaces, reducing the risk of numerical instabilities. For an  
408 example of data standardization see Section **Error! Reference source not found.** If standardization is applied to data, then the reconstruction procedure follows  
409 these steps: data standardization, performing the reconstruction on the standardized data, and then back-transforming  
410 the results to the original magnitude and scale of the original data.

412 Standardizing the data (by subtracting the mean and dividing by the standard deviation) is recommended to support  
413 the reconstruction process, particularly when the data has a wide range. Standardization reduces the search space of  
414 parameters to smaller, more manageable regions, helping to prevent numerical instabilities. For an example of data  
415 standardization, see Section 10.3 (**Example 10** and **Figure 6**). When standardization is applied, the  
416 reconstruction procedure involves three steps: (1) standardizing the data, (2) performing the reconstruction on the  
417 standardized data, and (3) back-transforming the results to the original scale and magnitude.

418 Consider the Langevin model (**Error! Reference source not found.**), i.e.,  $dx = \mu(x; \theta)dt + \sigma(x; \theta)dW$  and the  
419 standardization  $z = (x - m_d)/s_d$  where  $m_d$  and  $s_d$  are the mean and standard deviation of data. If the Langevin model  
420  $dz = \mu_z(z)dt + \sigma_z(z)dW$  describes the dynamics of the transformed process  $z$ , then the corresponding Langevin  
421 model for the original process  $x$  is as follows:

$$dx = s_d \mu_z((x - m_d)/s_d; \theta)dt + s_d \sigma_z((x - m_d)/s_d; \theta)dW, \quad (2)$$

422 If either of the drift function or diffusion function is linear in the parameters, then all we need to do is to replace the  
423 estimated vector of parameters  $\theta$  in (2) with  $s_d \theta$  and remove the factor  $s_d$ . For instance, if both  $\mu_z$  and  $\sigma_z$  are linear in  
424 terms of the parameter vector  $\theta$ , then (2) will be simplified to

425 
$$dx = \mu_z((x - m_d)/s_d; s_d \theta)dt + \sigma_z((x - m_d)/s_d; s_d \theta)dW.$$

426 Except for pchip splines, all the splines in this package are linear functions of parameters (though not of state  
427 variable). Even pchip splines are close to being linear. This reflects the ease of working with splines.

## 428 14. Euler reconstruction

### 429 14.1 Reconstructing a dataset simulated from a linear model

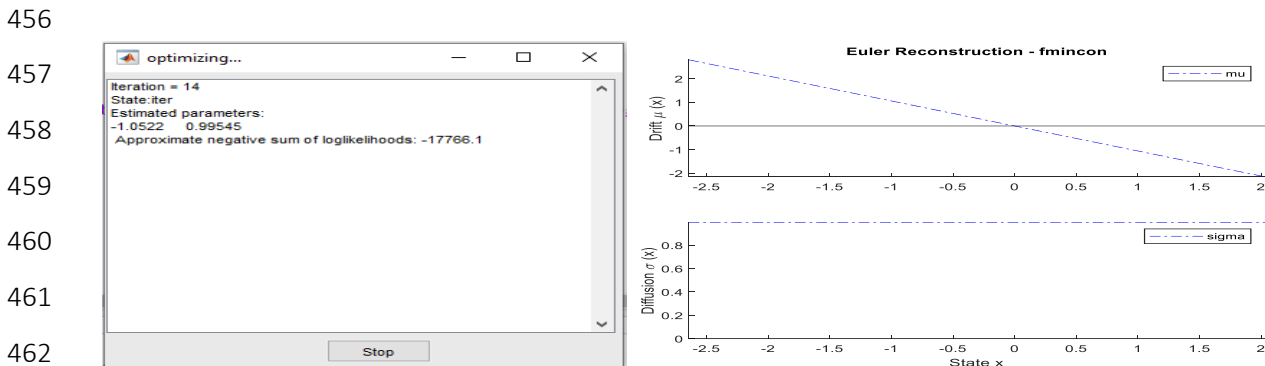
430 **Example 1.** In the first case study we apply the parametric and spline reconstruction techniques to a dataset being  
431 simulated from the OU process  $dx = \mu x dt + \sigma dW$  with parameters  $\mu = -1$  and  $\sigma = 1$ . In this dataset we used the  
432 time step of  $dt = 0.01$  and number of data points are  $T=10^6$  although here we only use the first 20000 data points  
433 (See **Figure 2**, left panel. See also **Example 10** where we use a very sparse sample of the entire dataset). Type the  
434 following commands

```
435 S = load('OUdata1D.mat');
436 data = S.data; %load the data
437 data = data(1:20000); %This is a big timeseries with 10^6 data points. Here, we just
438 %use its first 20000 data points
439 dt = 0.01;
440 mu = @(x,par)par(1).*x;
441 sigma = @(x,par)par(2);
442 result1 = euler_reconstruction(data, dt, 'mu', mu, 'sigma', sigma, ...
443 'gradient_fun', eulergrad(mu, sigma), 'lb', [-200 eps], 'ub', [200 200]);
```

444 And, you get

```
445 Estimated parameters :
446 -1.0586      0.99531
447 - sum of log-likelihoods : -17766.0856
```

448 While the package is solving the problem, a dialog screen appears on your screen (see **Figure 1**, left panel). This dialog  
449 screen is particularly helpful when the reconstruction process is slow or shows gradual improvement, such as when  
450 dealing with large or low-resolution datasets that require Hermite reconstruction. By monitoring the dialog screen, you  
451 can observe any slight improvements in the reconstruction process. If you notice minimal improvement, you can  
452 terminate the process by pressing the 'stop' button (of course, it is not the case with this small dataset). To assess  
453 the progress of the reconstruction, pay attention to the objective value displayed on the dialog screen, which represents  
454 the negative sum of log-likelihoods. A lower objective value indicates a better fit. If the objective value stops declining  
455 or decreases very slowly, you can safely terminate the code.



463 **Figure 1.** The left panel displays a dialog screen that appears during the package's execution, providing parameter updates during  
464 the optimization process. Additionally, the right panel showcases a graphical representation of a parametric model (**Example 1**)  
465 utilized for reconstructing the OU model.

466 We have only used the first 20000 data points of this large dataset with  $10^6$  data points (we will need the entire data in  
 467 the later sections). The model is parametric as we have defined the drift  $\mu(x)$  and diffusion  $\sigma(x)$  functions  
 468 using function handles. The model consists of two parameters: one for the drift function and one for the diffusion  
 469 function. We have specified a vector of lower bounds as `[-200 eps]`, where the lower bound for the first parameter  
 470 is -200 and the lower bound for the second parameter is `eps`, serving as an 'infinitesimal'. It is important to note that  
 471 the diffusion function must remain positive. While the code generally cannot check for the positivity of the diffusion  
 472 function in parametric models, for this simple additive model, the code will alert you if you overlook this requirement.  
 473 However, for more complex multiplicative noise models, you must verify this by yourself (For instance, if  $\sigma(x) =$   
 474  $x.^2 + 1 + a$  then  $a > -1$  should be fulfilled for the  $\sigma(x)$  to remain positive and you can easily check this by  
 475 plotting  $\sigma(x)$  as a function of state  $x$  and parameter  $a$ ). Additionally, we have defined a vector of upper bounds  
 476 as `[-200 200]`, indicating that both parameters are bounded by 200 from above. To visualize the results, use the  
 477 command `plot_results(result1)`.

478 **Example 2.** Let's now try a bit different parametric model here. Type the following commands

```
479 S = load('OUdata1D.mat');
480 data = S.data;
481 data = data(1:20000);
482 data(1:100:end) = nan; %this is to show you that the package works in the presence
483 of NaNs
484 dt = 0.01;
485 mu = @(x,par)par(1).*x+par(2).*x.^2;sigma = @(x,par)par(3); %here, we have 3
486 %parameters
487 result2 = euler_reconstruction(data, dt, 'mu', mu, 'sigma', sigma, ...
488 'gradient_fun', eulergrad(mu, sigma), 'lb', [-200 -200 eps], 'ub', [200 200 200]);
```

489 and you get

```
490 Estimated parameters :
491 -1.0735    -0.036727     0.99536
492 - sum of log-likelihoods) : -17410.5382
493 - sum of log-likelihoods : -17766.2684
```

494 In this example, we have introduced NaN values into every 100th data point to demonstrate the package's capability to  
 495 handle missing data. Additionally, we have augmented the former drift model with an extra quadratic term,  
 496  $\text{par}(2) \cdot x.^2$ , to assess if the package recognizes unnecessary terms. The second parameter is estimated to be -  
 497 0.036727. Notably, the second parameter estimate remains small, as expected. However, when using larger data  
 498 portions, the estimate for the second parameter tends to decrease. Now, which model provides a better fit? `result1`  
 499 or `result2`? To answer this question, compare the objective function values, i.e., the negative sum of log-likelihoods.  
 500 A lower objective value indicates a better fit. Surprisingly, `result2` exhibits a slightly smaller objective value,  
 501 suggesting it is a better fit. This outcome may seem counterintuitive, but it is expected. With a finite dataset length,  
 502 models with more parameters tend to yield better fits. However, as the dataset size increases, the discrepancy between  
 503 the models diminishes, and the estimates of additional parameters converge to zero

504 **Example 3.** Now, let's try a spline model for our data. Type the following commands

```
505 S = load('OUdata1D.mat');
506 data = S.data; %load the data
507 data = data(1:20000);
508 dt = 0.01;
509 L = -2;
510 R = 1.8; %since we have a 'spline' model it is better to shrink the state space
```

```

511 data(data<L | data>R) = nan; %This is VERY important: In spline modeling if you
512 consider a smaller range for your data then you must assign 'nan' to those few data
513 points falling outside this range.
514 mu = 8;
515 sigma = 8; %since mu and sigma are numbers this means that we want to consider
516 %spline modeling with 8 knots for mu and 8 knots for sigma
517 result3 = euler_reconstruction(data, dt, 'nknots', [nmu nsigma], 'spline', 'CC', 'L',
518 ...
519 L, 'R', R, 'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma)
520 + 10, 'solver', 'fmincon'); %we have 8+8 parameters, so, 'lb' and 'ub' should have 16
521 %elements. The vector of lower bounds 'lb' has 8 lower bounds for mu (which are -10) and 8
522 %lower bounds for sigma (which are eps, i.e., infinitesimal). All 16 elements of 'ub' are 10

```

523 and, you get

```

524 Estimated parameters:
525 4.8855    1.1983    1.1006    0.24633   -0.10136   -0.69903   -1.3182   -6.4891
526 0.84428    1.0242    1.0017    0.98554    1.0031    0.97949    1.0185    0.79658
527 - sum of log-likelihoods): -17730.1649

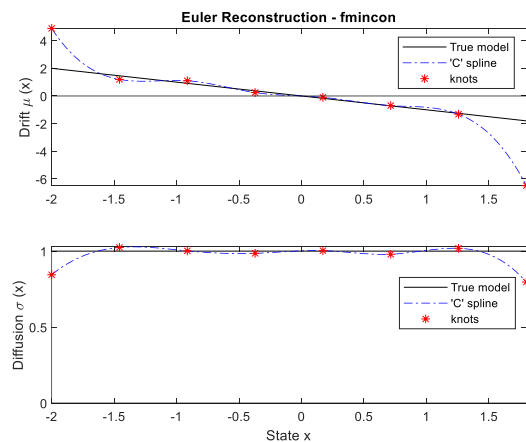
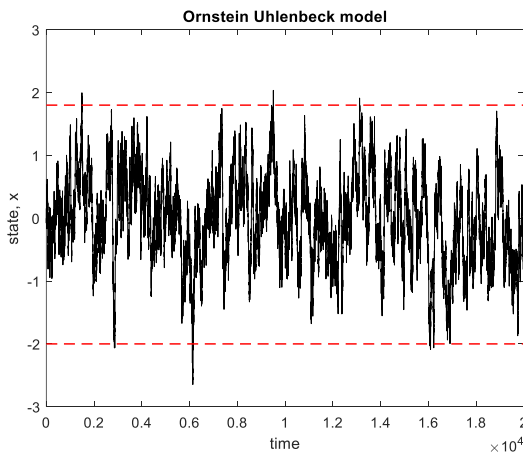
```

528 In this example, we have specified a spline model, unlike the previous examples. When using a spline model, the drift  
529  $\mu$  and  $\sigma$  functions should be numeric values only. For instance, specifying  $\mu=8$  indicates a spline model for the  
530 drift function with 8 equidistant knots across the state space (the same applies to  $\sigma=8$ ). By default, the state space  
531 is set to  $[\min(\text{data}) \ \max(\text{data})]$ , but it is recommended to narrow this range by specifying a larger lower limit  $L$   
532 and a smaller upper limit  $R$  for the data. This adjustment is beneficial because there are typically very few data points  
533 near the data borders, which can adversely affect the accuracy of the  $\mu$  and  $\sigma$  functions near data borders (i.e.,  
534 the first and last knots) (see **Figure 2**, right panel). For this dataset containing 20000 data points, we have chosen the  
535 range  $[-2 \ 1.8]$ , as only 0.000028% of the data fall outside this range (see **Figure 2**, left panel). *It is crucial to note*  
536 *that when utilizing spline modeling and opting for a smaller range for your data, any data point falling outside this*  
537 *reduced range must be assigned NaN.* This ensures that the spline model accurately reflects the specified data range  
538 and 'respects' the order of data in the smaller dataset. To visualize the estimated results and the true model, use the  
539 following commands (see **Figure 2**, right panel).

```

540
541 mu = @(x,par)par(1).*x;
542 sigma = @(x,par)par(2)+0.*x; %this is true model
543 par = zeros(2,1);par(1) = -1;par(2) = 1; %true model parameters
544 xplot = linspace(L,R,1000); % a dense mesh across the considered range
545 plot_results(result3,xplot,mu(xplot,par),sigma(xplot,par));
546

```



547

**Figure 2.** (Left panel) the first 20000 part of a dataset along with the reduced range being considered, highlighting the vast majority of data points within the specified range. The dataset is generated from the OU process  $dx = \mu x dt + \sigma dW$  with parameters  $\mu = -1$  and  $\sigma = 1$ , with a time step of  $dt = 0.01$  and number of data points are  $T=10^6$ . (Right panel) a graphical representation showcases a spline model featuring  $\mu = 8$  knots (indicated by red stars) for the drift function and  $\sigma = 8$  knots for the diffusion function. These functions are represented by dot-dashed blue curves, while the true model is depicted by solid black curves. The data are as in **Figure 1**.

## 14.2 Reconstructing a dataset simulated from a nonlinear model

**Example 4.** In the second case study we apply the parametric and spline reconstruction techniques to a dataset being simulated from the following stochastic version of overgrazed model of May(May 1977) with additive noise

$$dx = \left\{ rx \left( 1 - \frac{x}{K} \right) - \frac{\gamma x^2}{x^2 + a^2} \right\} dt + \sigma dW,$$

where the model parameters are  $r = 1.01, K = 10, \gamma = 2.75, a = 1.6, \sigma = 0.4$ . We have simulated a dataset containing  $3 * 10^5$  data points with time step  $dt = 0.01$  (**Figure 3**, left panel). We fit a parametric model to the first third of this dataset. Type the following

```
S = load('MayData1D.mat');
data = S.data; %load the data
data = data(1:100000); %We only use the first third of the dataset
dt = 0.01;
mu = @(x,par)par(1).*x.*(1-x./par(2))-par(3).*x.^2./(x.^2+par(4).^2);
sigma = @(x,par)par(5);
result4 = euler_reconstruction(data, dt, 'mu', mu, 'sigma', sigma, 'gradient_fun',
eulergrad(mu, sigma), ...
'lb', zeros(1,5)+eps, 'ub', 15.*ones(1,5), 'useparallel', true, 'search_agents', 5);
```

and, you get

```
Estimated parameters :
1.1805      9.8416      3.1242      1.5279      0.39972
- sum of log-likelihoods) : -180062.6377
```

This is a nonlinear model, making it more challenging to solve the underlying optimization problem due to the presence of multiple local minima. When dealing with nonlinear model, we recommend using the name value pairs 'gradient\_fun' to mitigate the risk of stagnation. Stagnation occurs when the optimization process gets stuck in a solution, which may not even be a local minimum (for instance, this can occur when the objective function is very flat). In such situations the code may continue for a long time, requiring manual intervention to stop and restart. To aid in identifying stagnation, dialog box is essential, allowing users to monitor the optimization process and stop it whenever it progresses slowly. In this example, the searching region of parameters is chosen to be a cube in the positive orthant (i.e., all parameters are positive) with sides of 15. Increasing the size of searching region rises the risk of stagnation. How can we determine an appropriate search region a priori then? Here are several considerations:

1. Utilize the `fmincon` solver, which is the default solver and is fast. Even if `fmincon` fails to converge to the true solution, it provides valuable insights into the approximate search region to explore.
2. Consider using the global solver `gwo`, although it may not be as fast. Initially, apply `gwo` to a large search region with a significant number of `search_agents` (e.g., 500) for a brief period to identify a smaller, more appropriate search region. This is not intended to find the true solution, but rather to gain insights into an approximate search region for later exploration using `fmincon`.
3. Opting for spline modeling simplifies and facilitates the optimization problem, as splines are linear functions of parameters, despite being non-linear in state variables. Standardizing the data (subtracting the mean and dividing by the standard deviation) enables the selection of narrower search regions for the parameters of  $\mu$



and `sigma`. We, recommend to use a searching region of `[-10 10]` for parameters of `mu` and `[eps 10]` for parameters of `sigma`, respectively (see **Example 3**). If it turns out that this search region is still small, you can simply choose a larger time step `dt`. It is important not to confuse this with the time step used for simulations, which should indeed be small. When performing the reconstruction, `dt` acts as a 'scale' parameter, so the specific number chosen is not important. For the simulated data, we chose the same time step `dt`, as in the simulations for proof of concept, even though it wasn't necessary to do so. Since splines are linear in terms of parameters, multiplying `dt` by a factor `k` will divide the drift parameters by `k` and the diffusion parameters by the square root of `k`. In technical terms, in a Langevin model (**Error! Reference source not found.**), i.e., the model  $dx = \mu(x)dt + \sigma(x)dW_t$  under the change of time scale  $\tau = kt$  the Langevin model becomes  $dx = 1/k \mu(x)d\tau + 1/\sqrt{k} \sigma(x)dW_\tau$ . For more details, see **Example 5**.

However, in cases where a parametric model is preferred and a small search region is chosen, incorrect solutions may still occur. For instance, in this example you might also get the following wrong result (roughly speaking, with 5 `search_agents` and the considered searching region, you can expect to obtain the correct answer approximately 80% of the time)

```
Estimated parameters :
0.0405253      4.52763      0.0252577      83.2777      0.400786
- sum of log-likelihoods) : -89896.9211
```

But we have a criterion to determine the correct solution: the solution with the smallest objective value is the fittest.

**Example 5.** Now, let's proceed to fit a spline model. Type the following lines

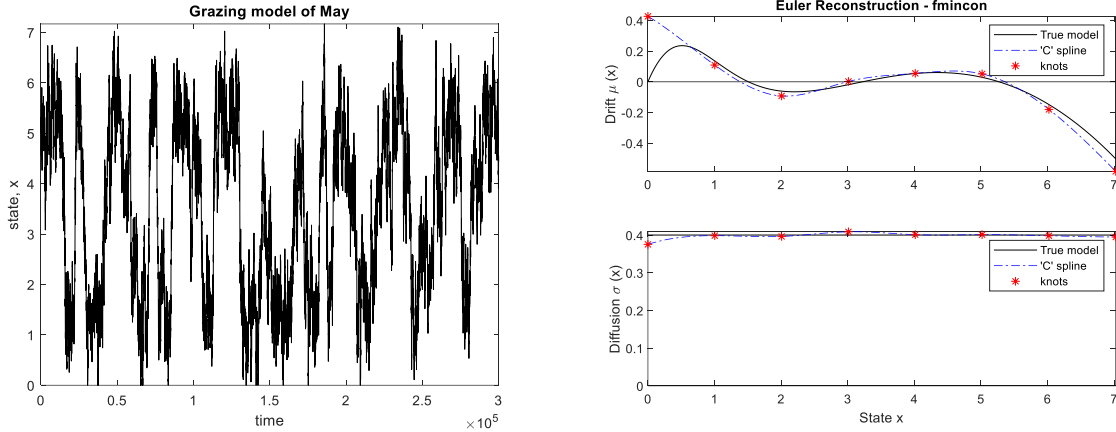
```
S = load('MayData1D.mat');data = S.data; %load the data
data = data(1:100000); %We only use the first third of the data
dt = 0.01;
L = min(data);R = max(data);
mu = 8;sigma = 8; % A spline model with 8 knots for mu and 8 knots for sigma
result5 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'CC', 'L',
...
L, 'R', R, 'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma)
+ 10, 'solver', 'fmincon', 'search_agents', 1);
```

and, you always get the following result

```
Estimated parameters :
0.4263      0.10858      -0.09268      0.0019837      0.054657      0.051057      -0.18033      -0.58468
0.37525      0.39834      0.39633      0.4086      0.40105      0.40098      0.39824      0.3948
- sum of log-likelihoods) : -180071.2916
```

and to see a plot type the following (see **Figure 3**, right panel)

```
r=1.01;K=10;g=2.75;a=1.6;s=0.4; % true parameter values
par = [r K g a s];
mu = @(x,par)r.*x.*(1-x./K)-g.*x.^2./(x.^2+a.^2);sigma=@(x,par)s; % true model
xplot=linspace(L,R,2000);
plot_results(result5,xplot,mu(xplot,par),sigma(xplot,par))
```



**Figure 3.** (Left panel) A dataset simulated from the grazing model of May.  $3 \times 10^5$  data points with a time step of  $dt = 0.01$  are simulated from the overgrazed model of May  $dx = \left\{ rx \left( 1 - \frac{x}{K} \right) - \frac{\gamma x^2}{x^2 + a^2} \right\} dt + \sigma dW$  with parameters  $r = 1.01, K = 10, \gamma = 2.75, a = 1.6, \sigma = 0.4$ . (Right panel) A graphical illustration for a spline model with 8 knots (red stars) considered for both drift and diffusion functions (dot-dashed blue curves) together with the true model (solid black curves).

Some explanations for this spline model. First, we did not specify any lower and upper bound for the dataset since in this case there exist enough data across the state space. In such cases the solver chooses the default option  $[\min(\text{data}), \max(\text{data})]$  for the state space. Second, we only chose one `search_agents`. However, note that this spline model has 16 parameters, yet a single search agent was sufficient to obtain the only solution. Third, this problem has only one solution. This highlights the simplicity of working with spline models, which are recommended. Splines are composed of simple polynomial building blocks and are flexible structures capable of adapting to complex functional forms with unknown nonlinearities. Consequently, spline modeling imposes less pressure on the optimization problem and often yields a unique solution. Forth, in spline modeling we do not use gradient. Calculating a gradient vector computationally requires  $O(n)$  operations where  $n$  is the number of parameters. However, the computational time to estimate them numerically using MATLAB's finite-difference methods requires roughly the same operations. As the underlying MLE procedure for spline models is easier, we opt not to pass a gradient. Fifth, in general, selecting a proper model for parametric modeling can be challenging. Therefore, we emphasize the importance of the spline modeling approach. Even if there is a preference for a parametric model, we recommend starting with a spline model to gain insights into the functional form. Attempting to fit an improper parametric model to the data can result in issues such as longer execution times, decreased accuracy, and stagnation.

In **Figure 3**, the left tail of the May model did not appear in the reconstructed model. This is not related to stuff like numerical inaccuracies, estimation errors, etc. Indeed, this phenomenon occurs with any positive dataset generated by the May model. For an explanation about this refer to section 12.

### 14.3 Reconstructing an ecological dataset

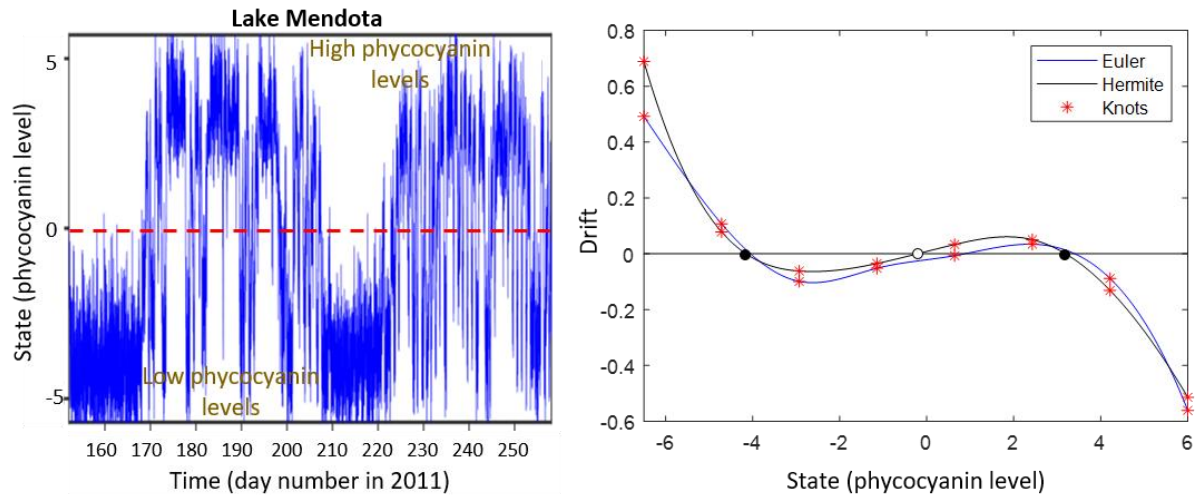
**Example 6.** Here, we apply a spline reconstruction to a univariate time series of phycocyanin concentrations in Lake Mendota (Carpenter et al. 2020). This dataset has a high resolution, with measurements taken at minute intervals. We focus on a period during summer thermal stratification in 2011, a period when Cyanobacterial blooms are common (see **Error! Reference source not found.**, left panel). For further details on this dataset we refer you to the references (Arani et al. 2021, Magnuson et al. 2023). We do not standardize this real dataset since it is already the standardized level of phycocyanin concentrations (for further details on this read the appendix of (Arani et al. 2021)). Furthermore, this dataset does not satisfy one of the data requirements mentioned in section 6.3 since it is not Markov. However, a rarified sample of this dataset with every forth data point (which is still high resolution) is Markov and ME time scale is 4 (see Table 1). Although this dataset has high resolution, we apply both the Euler and Hermite reconstructions. For real datasets, it is recommended to use both reconstruction procedures, unless the dataset has an exceptionally high resolution which is not the case here. Below are the relevant code lines to implement both reconstruction algorithms (see Section 10 for a detailed discussion on Hermite reconstruction)



```

671 data = readmatrix('BGA_stdlevel_2011.csv');
672 data = data(1:4:end); % From Table1, we see that this dataset is not Markov. But, a
673 rarified sample with every third data point is Markov
674 dt = 1; % This is completely arbitrary.
675 L = -6.5; R = 6;
676 data(data<L | data>R) = nan;
677 mu = 8; sigma = 1; % An additive spline model with 8 knots for mu
678 result_euler6 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'QQ',
679 'L', ...
680 L, 'R', R, 'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma)
681 + 10, 'solver', 'fmincon', 'search_agents', 5);
682 legpoints = legitimate_points(data, dt, 'prev', result_euler6, 'prev_range', 0.5,
683 'j', 3, 'k', 12);
684 result_her6 = hermite_reconstruction(data, dt, 'prev', legpoints, 'solver',
685 'fmincon');
686

```



687  
688 **Figure 4.** (Left) High-resolution phycocyanin levels measured every minute in Lake Mendota in 2011. While this dataset is not  
689 Markov (so it does not meet one of the data requirements outlined in section 6.3 (refer to **Table 1**)), a rarified sample of this dataset,  
690 including every forth data point, is Markov and analysis is applied to this ratified sample. (Right) Estimated drift functions using  
691 Euler (blue) and Hermite (black) reconstructions. Given the high resolution of the dataset, the difference between the two drift  
692 function is relatively small. However, this small misestimating results in an unusually small basin of attraction for the high-  
693 phycocyanin state compared to the low-phycocyanin state. The open circle displays the location of the repelor separating the two  
694 attraction basins (alternatively, the red dashed line in the left panel represents the repelor). The estimated noise levels are 0.4396  
695 (Euler) and 0.4567 (Hermite), respectively.

696 which leads us to the following result for the Euler reconstruction

697  
698 Estimated parameters :

```

699 0.49164    0.10753   -0.098849   -0.051798   -0.0066688    0.034236   -0.088322   -
700 0.56085    0.43961
701 - sum of log-likelihoods) : 22751.7848

```

702 And the following result for the Hermite reconstruction

703 Estimated parameters :

```

704 0.68814    0.078596   -0.061772   -0.034591    0.033261    0.050505   -0.13094   -
705 0.51218    0.45664

```

- sum of log-likelihoods : 22697.8381

We considered an additive model. Therefore, we have 9 parameters: 8 drift parameters and a single noise parameter (both noise parameters are highlighted in blue). There are two reasons to additive modeling. First, in this example, we need to compare the performance of the Euler and Hermite reconstructions in terms of their corresponding attraction basins. Calculating attraction basins for multiplicative models is more involved and requires the concept of ‘effective potential’, which we introduce and describe in Section 11. Second, since this dataset is rather large, using a multiplicative model would make Hermite reconstruction significantly more time-consuming.

**Error! Reference source not found.**, right panel depicts the estimated drift functions obtained using Euler (blue) and Hermite (black) reconstructions. Since this dataset has high resolution, the difference between the two drift functions is rather small (The difference between the two reconstruction algorithms becomes more pronounced as the data resolution decreases (for example, see Example 12), and it approaches 0 as the resolution gets higher). However, notable discrepancies still exist. In particular, the Euler reconstruction poorly estimates the locations of the alternative stable states—especially the repeller separating the attraction basins (see the dashed red line in the left panel or the open circle in the right panel). This leads to a pronounced imbalance between the sizes of basins of attraction, which can significantly affect resilience quantification. This example highlights that for real-world datasets, it is generally safer to rely on Hermite reconstruction, unless the resolution is extremely high (which is not the case here). The estimated noise intensities (that are highlighted in blue) are 0.4396 (Euler) and 0.4567 (Hermite), indicating close agreement between the two algorithms.

## 15. Sparsely sampled data: Hermite reconstruction

It is not uncommon to encounter datasets with low-resolution especially in life science, presenting a challenge for reconstruction techniques like Euler reconstruction, which typically require at least a medium resolution for effective results. To address the issue of sparsely sampled data, our package implements a reconstruction technique based on a refinement approach by Bakshi, et.al (Bakshi and Ju 2005), building upon the work of Aït-Sahalia (Aït-Sahalia 2002) for univariate data. Aït-Sahalia's method involves constructing a sequence of converging closed-form expansion of ‘transition density’  $p(x(t + \Delta) | x(t))$  using Hermite polynomials for which there is no closed form in almost all stochastic models. Therefore, we call it ‘Hermite reconstruction’. Hermite reconstruction offers a higher level of accuracy but comes at a higher computational cost compared to the Euler reconstruction. When the data resolution is low Hermite reconstruction should be implemented; otherwise, the use of Euler reconstruction may lead to inaccurate results. Hermite reconstruction has the capacity to enhance the accuracy of Euler reconstruction to some extent.

When implementing Hermite reconstruction, the package needs Euler parameter estimation as a starting guess. Here, the procedure has two phases in which in the first phase the Euler reconstruction is implemented. In the second phase, the package follows Hermite reconstruction and aims to enhance Euler estimation by exploring a ‘small’ region in the parameter space around the Euler estimation. Hermite reconstruction requires two key inputs: J and K which determine the number of terms one includes in the Hermite expansion of transition density (and hence the likelihood function). J represents the number of ‘spatial’ terms in the Hermite expansion of likelihood function using Hermite polynomials, while K represents the number of ‘temporal’ terms in the Taylor expansion, in terms of sampling time  $dt$ , of Hermite coefficients. High J and/or K improve estimation accuracy at the cost of higher computation time. Typically, a small J suffices (for technical details, see the last paragraph on page 2 of Aït-Sahalia’s paper (Aït-Sahalia 2002)), and in all case studies in the tutorial, we used  $J=3$ . However, as data resolution decreases, a bigger K is necessary to enhance estimation accuracy. Based on our experience, for data with high resolution, a value of  $K \leq 6$  is generally sufficient while for low-resolution data higher values of K such as  $9 \leq K \leq 12$  may be necessary.

### 15.1 When does Hermite reconstruction crash? Strategies and recommendations

Hermite reconstruction works by constructing a sequence of closed-form expansions of transition density using Hermite polynomials. However, this density approximation might not converge to a positive density, leading to an undefined objective value, in situations where data resolution is very low and initial parameters are rather far from the optimal parameters. To inspect this, the package attempts to find some starter parameter values in the vicinity of the

Euler solution obtained in the first phase, which we call '*legitimate points*', where the objective function is defined. The function `legitimate_points` is responsible for this task. The challenge lies in finding the first legitimate point. Therefore, whether or not the package is able to tackle Hermite reconstruction boils, primarily, down to finding this first legitimate point (a feasibility problem). If the package can find it, then the problem is often tractable; otherwise, the problem is considered intractable based on the modeling strategy adopted (although the package might take a while to find the first legitimate solution, subsequent ones will be found faster). If the former occurs, the package uses the legitimate points first and implements a surrogate optimization for a very short time. Surrogate optimization is a technique used to optimize complex, computationally expensive, black-box, or undefined objective functions by replacing them with simpler surrogate models that approximate their behavior (Koziel and Leifsson 2013). Note that surrogate optimization is not among the optimization solvers you can use, rather it is an internal optimization which is performed solely by the package. Surrogate optimization then gives us a solution. If surrogate optimization finds a solution better than all legitimate solutions (indicated by a message in the command window), then the package seeks a further improvement starting from the surrogate solution, using the solver chosen by the user which is either `fmincon` (recommended) or `gwo`. Otherwise, the package tries to make progress using the legitimate solutions as starters. However, if the package either fails to find a first legitimate solution or consumes a considerable amount of time to do so (a message in the command window appears when a legitimate solution is found) it suggests that the objective function is severely damaged and the considered model is not reasonable. In such cases, it is advisable to consider changing the modeling strategy by trying simpler models to data. Below, are very useful strategies:

- 1) *A The best advice is to use quadratic spline modeling, using 'QQ' flag (or, simply 'Q' flag (or, simply 'Q' flag if additive modeling is preferred) instead of cubic spline modeling.* This significantly reduces the computational burden and increases the likelihood of the package in finding legitimate solutions. Note that irrespective of the flag considered the output provided by the package is always cubic spline. Thus, if the flag 'QQ' is used the package fits a cubic spline mode to the obtained quadratic spline model. This is because cubic splines are computationally more appealing since they have better smoothing properties.
- 2) Prioritize additive models over multiplicative models despite potential limitations in efficiency.
- 3) Opt for models that are linear in parameters. Spline modeling becomes significant here, as splines are linear in terms of parameters while being nonlinear in terms of the state variable. Also, consider models with fewer parameters such as additive models.
- 4) For low-resolution data, using  $K \geq 10$  is not recommended in situations where 1) the dataset is real or large (e.g., the number of data points  $> 10^5$ ), or 2) the model is complex (e.g., a parametric model that is nonlinear in its parameters, like the May model or the number of parameters is large, or the model is cubic spline), or 3) the data resolution is extremely low. In such situations, using large K values can lead to a complex and computationally intensive optimization problem making it difficult to find legitimate solutions. Often,  $K=9$  is sufficient. Nonetheless, we usually use  $K=12$  whenever either of these conditions are violated.

We illustrate these points through several examples in this section, where Hermite reconstruction is applied to three datasets: one generated by a linear OU model, another by a nonlinear May model, and an ice-core climate dataset.

## 15.2 Reconstructing a dataset simulated from a linear model

**Example 7.** In this example, we apply Hermite reconstruction to the same dataset used in **Example 1** which was generated from the OU model  $dx = \mu x dt + \sigma dW$  with parameters  $\mu = -1$  and  $\sigma = 1$ . Here, we use parametric reconstruction. The dataset has a time step of  $dt=0.01$  and contains  $T=10^6$  data points. However, we select every 100<sup>th</sup> data point, resulting in an extremely low-resolution dataset with a time step of  $dt=1$  and  $T=10^4$  data points. The resolution of this dataset is the minimum resolution required for any reconstruction algorithm to function properly in theory. This is so because the relaxation time (see section 6.3 for details) for this dataset is 1. To see this, type the following commands

```
S = load('OUdata1D.mat');
data = S.data;
data = data(1:100:end);
R = RelaxationTime(data);R
```

```

800
801 1.0073

802 Where R should, in theory, be 1 if we had a longer dataset. This relaxation time signifies the lowest resolution threshold
803 (in theory) below which all reconstruction procedures fail. Now, we perform Euler reconstruction (first phase). Type
804 the following commands

805 S = load('OUdata1D.mat');
806 data = S.data; %load the data
807 data = data(1:100:end); % Only every 100 data points are considered
808 dt = 1; % note that the mother dataset has the time step of dt=0.01 which is multiplied
809 %by 100 to match the time scale of this sample
810 mu = @(x,par)par(1).*x;sigma = @(x,par)par(2);
811 result6 = euler_reconstruction(data, dt, 'mu', mu, 'sigma', sigma, 'gradient_fun',
812 eulergrad(mu, sigma), 'lb', [-200 eps], 'ub', [200 200], 'useparallel',true, 'solver',
813 'fmincon', 'search_agents', 5);

814 and you get

815 Estimated parameters :
816 -0.63555      0.65456
817 - sum of log-likelihoods): 9950.4726

818 This estimate deviates significantly from the true parameter values due to the fact that we rarified the original dataset
819 to create this low-resolution dataset. This (Euler reconstruction) marks the completion of the first phase. Moving on to
820 the second phase, we employ Hermite reconstruction, wherein the package explores in the vicinity of the Euler
821 estimation to improve it. At this stage you have two choices to pick an optimization solver: fmincon or gwo. Opting
822 for fmincon (recommended) involves an initial step of running the function legitimate_points to find at least N
823 = search_agents legitimate starting points before proceeding to solve the underlying optimization problem (MLE).
824 Subsequently, you can execute the main function hermite_reconstruction to estimate the optimal parameter
825 values. This sequential approach is essential because fmincon cannot start optimization with an infeasible solution,
826 although it can fortunately cope with intermediate infeasible solutions, to some extent, all the way to the optimal
827 parameter values. However, for extremely low-resolution data and nonlinear models, there's a higher risk of the
828 fmincon solver crashing, whereas the gwo solver, although resilient to crashes, may be considerably slower in such
829 scenarios. Therefore, our recommendation is to try fmincon first. Nonetheless, if you prefer to utilize the gwo solver,
830 you can directly proceed with the hermite_reconstruction function and skip executing legitimate_points
831 since gwo internally calls it.

832 Another consideration is the selection of two parameters J ( $J \geq 3$ ) and K ( $K \geq 1$ ), which are required for the
833 implementation of Hermite reconstruction. In the second paragraph of this section we elaborated on appropriate choices
834 of these two parameters but here we recall it briefly. The bigger these parameters are the more accurate result we get
835 but this comes at the expense of higher computational time. Therefore, appropriate values should be assigned to these
836 parameters based on the data resolution. Typically,  $J=3$  is sufficient. However, for enhanced parameter estimation, a
837 larger value of K is necessary. For high-resolution data  $K \leq 6$  is sufficient while for low-resolution data  $9 \leq K \leq 12$  may be
838 necessary. For low-resolution data we often use  $K=12$  when 1) the model is structurally simple (e.g., spline models or
839 models that are linear in their parameters as is the case in this example), and 2) the dataset is not large (e.g., the number
840 of data points  $< 10^4$ ). This includes all the examples in this section except Example 8 where we used  $K=9$ .

841 Now, we proceed with finding some legitimate solutions (read subsection 10.1 for details) using the following
842 command

843 legpoints = legitimate_points(data, dt, 'prev', result6, 'prev_range', 0.5, 'j', 3,
844 'k', 12);

845 Finally, to implement Hermite reconstruction type the following command

```

```

846 result_her6 = hermite_reconstruction(data, dt, 'prev', legpoints, 'solver', 'fmincon');
847 and you get the following great result
848 Estimated parameters :
849 -1.0093      0.99867
850 - sum of log-likelihoods : 9950.4912
851 Example 8. Let's now try to fit a spline model to the same dataset used in Example 7. For the first phase (Euler
852 reconstruction) type the following commands
853 S = load('OUdata1D.mat');data = S.data; %load the data
854 data = data(1:100:end); % Only every 100 data points are considered
855 L = -2.5;R = 2.5;%Since we have a spline model it is better to shrink the state space
856 data(data<L | data>R) = nan; %This is VERY important: In spline modeling if you
857 consider a smaller range for your data then you must assign 'nan' to those few data
858 points falling outside this range.
859 dt = 1;
860 mu = 7; sigma = 7; %In spline modeling mu and sigma are numbers
861 result7 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'QQ', 'L',
862 L, 'R', R, ...,
863 'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma) + 10,
864 'solver', 'fmincon', 'search_agents', 5);
865 And you always get the following solution
866 Estimated parameters:
867 1.7415      1.0522      0.52421      -0.0032876      -0.52616      -1.1204      -1.6048
868 0.38686      0.65803      0.65051      0.66209      0.63977      0.63572      0.48534
869 - sum of log-likelihoods): 9920.7613
870 It is important to note that for this example, we have adjusted the range of the data due to the implementation of spline
871 modeling. Typically, the range of state space is, by default, [min(data) max(data)] = [-2.6289 2.9383].
872 However, in this example, we have narrowed it down to [-2.5 2.5]. This decision was made because the dataset
873 contains only 2 data points larger than 2.5 and 2 data points smaller than -2.5 out of a total of 10,000 data points. This
874 situation often arises with small or relatively small datasets that have few data points near the data borders. Failing to
875 address this issue can negatively impact the estimation process in terms of both accuracy and speed, particularly during
876 the second phase when Hermite reconstruction is implemented. Additionally, note that we've assigned nan values to
877 those few data points that fall outside the considered range (i.e., the code line data(data<L | data>R) = nan).
878 Important to note is that we have chosen the spline flag 'QQ' (i.e., quadratic spline models for both drift and diffusion
879 functions) not the typical flag 'CC' (cubic spline modes for both drift and diffusion functions) we used before. If you
880 only wish to apply Euler reconstruction, then this choice does not matter a lot, so you can safely use the flag 'CC'. If,
881 however, the final goal is to apply a Hermite reconstruction to data then it is recommended to use the spline flag 'QQ'.
882 This greatly speeds up the calculations, increases the chance of finding legitimate solutions and, further increases the
883 chance of improving the legitimate solutions later. Note, however, that at the end the package fits cubic spline models
884 to the obtained quadratic splines. So, you always get cubic spline results even if you choose quadratic spline modeling.
885 mu = 7 and sigma = 7 are considered. We elaborate on these choices later. To implement the second phase, first
886 type the following command
887 legpoints = legitimate_points(data, dt, 'prev', result7, 'prev_range', 0.5, 'j', 3,
888 'k', 9);
889 which provides us with at least search_agents = 5 legitimate solutions. Finally, type the following command to
890 implement the Hermite reconstruction

```

```

891 result_her7 = hermite_reconstruction(data, dt, 'prev', legpoints, 'solver', 'fmincon',
892 'search_agents', 5);

893 and you get

894 Estimated parameters:
895 2.2275    1.6342    0.86879    0.020833   -0.8565   -1.6008   -2.1243
896 0.92026    1.0035    1.0253    1.0373    1.0147    0.96308    0.86743
897 - sum of log-likelihoods: 9913.6229

898 Unlike Euler spline reconstruction, Hermite spline reconstruction may yield slightly different results every time you
899 run it. This variability is due to the quality of the legitimate solutions obtained. However, different solutions usually
900 closely approximate each other. This is because the use of splines typically leads to an objective function with simple
901 structure with one local minimum or few local minima which are close to each other.

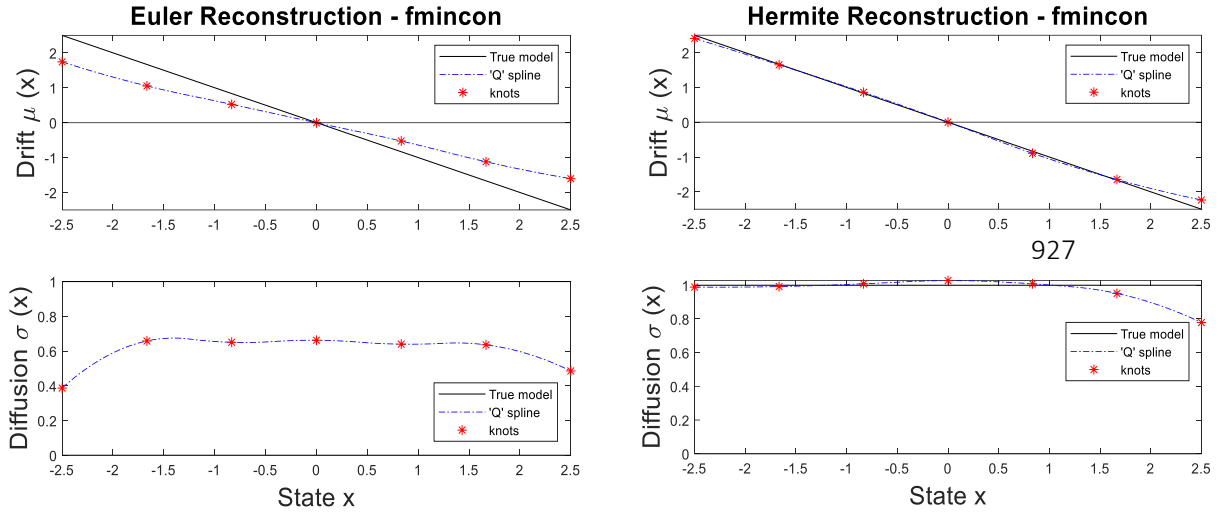
902 Several other points are noteworthy here. First, the lower objective value in this example (i.e., 9913.6229) compared
903 to Example 7 (i.e., 9950.8106) indicates that the model in this example is a better fit. This outcome was expected as
904 the model in this example comprises  $\mu + \sigma = 14$  parameters, whereas the model in Example 7 had only 2
905 parameters. A key point to note is that, although the dataset was generated from the OU model, the spline model
906 outperformed the parametric OU model. This highlights the strength of spline modeling and justifies the superiority
907 and elegance of spline models even in situations where we know the true functional form of the data-generating system.
908 Second, in this example, we fit a multiplicative spline model to data using  $\mu = 7$  knots for the drift function and
909  $\sigma = 7$  knots for the diffusion function. In Hermite reconstruction, we need to be economical with respect to the
910 number of parameters. Attempting  $\mu = 8$  and  $\sigma = 8$  also works but the computational time increases and the
911 package has difficulty in making a progress (while it took around 20 seconds to find 5 legitimate solutions here, it
912 would take a few minutes for the case  $\mu = 8$  and  $\sigma = 8$ ). Furthermore, if the spline flag 'CC' would be used,
913 instead, it would be very hard to find legitimate solutions beyond the cases where  $\mu \geq 5$  and  $\sigma \geq 5$  in a reasonable
914 time. A general advice, therefore, is to use the spline flag 'QQ' and be economical on the number of parameters.
915 Otherwise: 1) it takes a long time for the package to find legitimate solutions, and 2) the package fails to improve upon
916 the Euler reconstruction. Due to the 'curse of dimensionality', as the parameter space expands, finding legitimate
917 solutions becomes increasingly challenging. If you wish to consider a bigger knot sequence, then a good strategy is to
918 consider an additive model (see the next example). Third, since this dataset has the minimal resolution and is a bit
919 large we considered  $K=9$  which is sufficient. To get a plot (see Figure 5), type the following commands

920 mu = @(x,par)par(1).*x;sigma = @(x,par)par(2)+0.*x; %this is true model
921 par = zeros(2,1);par(1) = -1;par(2) = 1; %true model parameters
922 xplot = linspace(L,R,2000); % a dense mesh across the considered range
923 plot_results(result7,xplot,mu(xplot,par),sigma(xplot,par)); % Euler & true models
924 plot_results(result_her7,xplot,mu(xplot,par),sigma(xplot,par));%Hermite & true models
925

```



926



927

928

929 **Figure 5. Illustration of Euler and Hermite multiplicative reconstructions applied to a low-resolution simulated dataset**  
 930 **generated by a linear model.** The left panel depicts the true model alongside the Euler reconstructed model, while the right panel  
 931 depicts the true model alongside the Hermite reconstructed model. The data are simulated from the OU model  $dx = \mu x dt + \sigma dW$   
 932 with parameters  $\mu = -1$  and  $\sigma = 1$ . The original dataset has the time step of  $dt=0.01$  and contains  $T=10^6$  data  
 933 points. However, we select every 100<sup>th</sup> data point, resulting in an extremely low-resolution dataset with a time step of  
 934  $dt=1$  and  $T=10^4$  data points.

935 **Example 9.** In this example we reconstruct the same dataset as in **Example 8**. However, we use an additive spline  
 936 model with more knots. Type the following command lines.

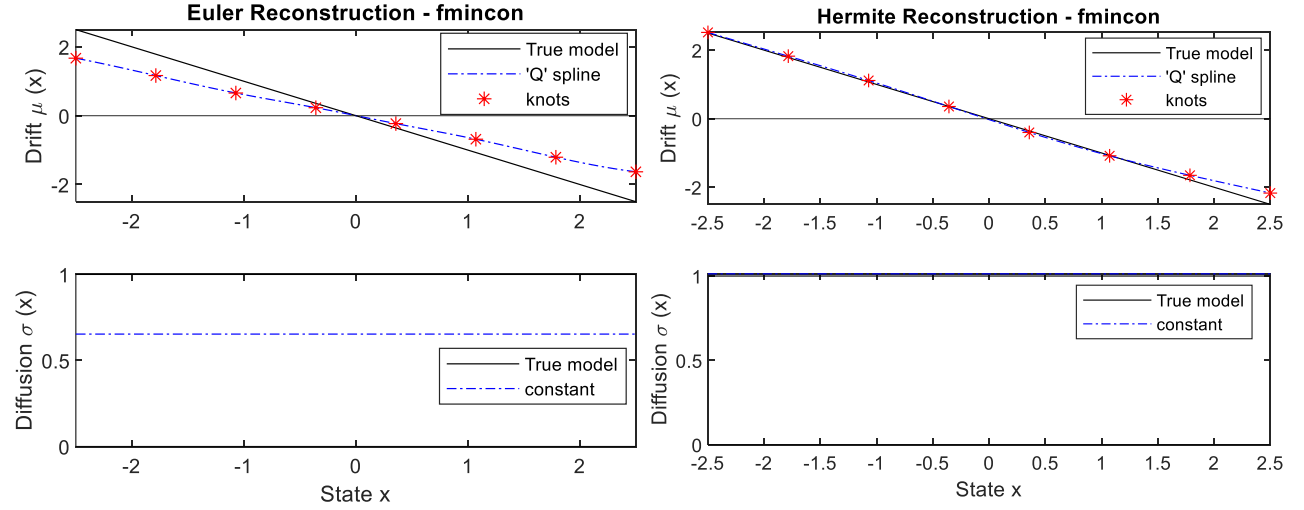
```

937 S = load('OUdata1D.mat');data = S.data; %load the data
938 data = data(1:100:end); % Only every 100 data points are considered
939 L = -2.5;R = 2.5; %Since we have a spline model it is better to shrink the state
940 space
941 data(data<L | data>R) = nan; %This is VERY important: In spline modeling if you
942 consider a smaller range for your data then you must assign 'nan' to those few data
943 points falling outside this range.
944 dt = 1;
945 mu = 8; sigma = 1; %In spline modeling mu and sigma are numbers
946 result8 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'QQ', 'L',
947 L, 'R', R, ...,
948 'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma) + 10,
949 'solver', 'fmincon', 'search_agents', 5);
950
951 legpoints = legitimate_points(data, dt, 'prev', result9, 'prev_range', 0.5, 'j', 3,
952 'k', 9);
953 result_her9 = hermite_reconstruction(data, dt, 'prev', legpoints, 'solver',
954 'fmincon', 'search_agents', 5);
955
956 mu=@(x,par)par(1).*x;sigma=@(x,par)par(2)+0.*x; %this is true model
957 par=zeros(2,1);par(1)=-1;par(2)=1; %true model parameters
958 xplot=linspace(L,R,2000); % a dense mesh across the considered range
959 plot_results(result8,xplot,mu(xplot,par),sigma(xplot,par)); % Euler & true models
960 plot_results(result_her8,xplot,mu(xplot,par),sigma(xplot,par));%Hermite & true models

```

961 **Figure 86** illustrates the outcomes of Euler, Hermite and true models. In this additive model, the Hermite objective  
 962 value (9918.9198) is higher than that in the multiplicative model in **Example 8** (9913.6229) which means that the

multiplicative model is a better fit. However, it was much faster for the package to handle this example since this model has less parameters.



**Figure 6. Illustration of Euler and Hermite additive reconstructions applied to a low-resolution simulated dataset generated by a linear model.** The left panel depicts the true model alongside the Euler reconstructed model, while the right panel depicts the true model alongside the Hermite reconstructed model. The data are simulated from the OU model  $dx = \mu x dt + \sigma dW$  with parameters  $\mu = -1$  and  $\sigma = 1$ . The original dataset has the time step of  $dt=0.01$  and contains  $T=10^6$  data points. However, we select every 100<sup>th</sup> data point, resulting in an extremely low-resolution dataset with a time step of  $dt=1$  and  $T=10^4$  data points.

### 15.3 Reconstructing a dataset simulated from a nonlinear model

Here, we apply Hermite reconstruction to the following nonlinear model which is a stochastic version of the overgrazed model of May(May 1977)

$$dx = \left\{ rx \left( 1 - \frac{x}{K} \right) - \frac{\gamma x^2}{x^2 + a^2} \right\} dt + \sigma dW,$$

where the model parameters are  $r = 1.01, K = 10, \gamma = 2.75, a = 1.6, \sigma = 0.4$ . We have simulated a dataset containing  $3 * 10^5$  data points with time step  $dt = 0.01$ . We consider estimating the parameters of this model by rarifying this data set by considering every 300<sup>th</sup> data points to get a sparse sample with time step  $dt = 3$  and only 1000 data points. Let's first check the resolution of this dataset via its relaxation time. Type the following commands

```
S = load('MayData1D.mat'); data = S.data; %load the data
data=data(1:300:end); %We consider every 300-th data point
RelaxationTime(data)
```

and you get

```
12.7959 (unit of data)
```

which puts this rarified dataset in the category of low-resolution. We recall that a dataset with relaxation time in the interval  $[1 \ 100]$  is considered low-resolution (see subsection 6.3 for more details).

**Example 10.** Here, we try to fit a parametric model to this dataset. If we try a model with drift term  $\mu = @(x,par)par(1).*x.*(1-x./par(2))-par(3).*x.^2./(x.^2+par(4).^2)$  and diffusion term  $\sigma = @(x,par)par(5)$  then Hermite reconstruction cannot improve the Euler reconstruction. The core of difficulty is that this model is nonlinear in terms of two parameters:  $par(2)$  and  $par(4)$ . We, therefore, try to fit the following cubic polynomial model which is linear in terms of parameters. Type the following commands



```

993 S = load('MayData1D.mat');data = S.data; %load the data
994 data=data(1:300:end); %We consider every 300-th data point
995 dt = 3; % Since in the original dataset dt=0.01 and here we considered every 300-th
996 data points the actual time step is 3
997 m = mean(data);s = std(data);
998 data = (data-m)./s;
999 mu = @(x,par)par(1).*x.^3+par(2).*x.^2+par(3).*x+par(4); %this is a standardized
1000 drift model
1001 sigma = @(x,par)par(5);
1002 result9 = euler_reconstruction(data, dt, 'mu', mu, 'sigma', sigma, 'gradient_fun',
1003 eulergrad(mu, sigma), ...
1004 'lb', [-5.*ones(1,4) 0], 'ub', [5.*ones(1,4)
1005 5],'useparallel',true,'search_agents', 5); % Since the model lb for drift parameters
1006 is chosen to be symmetrical about 0, i.e., the interval [-5 5]

1007 we then obtain

1008 Estimated parameters:
1009 -0.038801 -0.0045207 0.041769 0.0017958 0.20353
1010 - sum of log-likelihoods): 375.9106

1011 Here, we provide explanations. First, it is important to note that the original dataset has the time step of  $dt = 0.01$ .
1012 However, since the rarified dataset consists of every 300th data points of the original dataset, its time step is  $dt = 3$ .
1013 As mentioned previously, the choice of time step during the reconstruction process acts as a scale parameter and is
1014 completely arbitrary. Nonetheless, for the sake of validating our approach, we opt to use the actual time steps. Second,
1015 to streamline the optimization process, we standardized the data. Consequently, the corresponding parameters can vary
1016 around 0, and their magnitudes are not significantly different from 0. Therefore, we used a symmetrical lower and
1017 upper bound for each drift parameters, i.e., the interval  $[-5 \ 5]$ . For the noise parameter, we need to consider the
1018 interval  $[0 \ 5]$  to ensure the positivity of the diffusion function. Finally, in order to recover a model that accurately
1019 represents the original data, we must back-transform the parameters to their original scales. This involves subtracting
1020 the state variable by the mean of the data, dividing by the standard deviation of the data, and then multiplying the entire
1021 system by the standard deviation of the data. Consequently, the following drift and diffusion functions model the
1022 original data

1023 mu = @(x,par)s.*par(1).*((x-m)./s).^3+s.*par(2).*((x-m)./s).^2+s.*par(3).*(x-
1024 m)./s+s.*par(4);sigma=@(x,par)s.*par(5);

1025 where m = mean(data) and s = std(data). Next, we get at least 5 legitimate points by the following command
1026 legpoints = legitimate_points(data, dt, 'prev', result9, 'prev_range', 0.5, 'j', 3,
1027 'k', 12);

1028 and, finally we go for Hermite reconstruction as bellow

1029 result_her9 = hermite_reconstruction(data, dt, 'prev', legpoints, 'solver',
1030 'fmincon', 'search_agents', 5);

1031 to obtain

1032 Estimated parameters:
1033 -0.044756 -0.0042904 0.0464 0.0022962 0.22073
1034 - sum of log-likelihoods: 362.5456

1035 The lower objective value for the Hermite reconstruction (362.5456) compared to the Euler reconstruction
1036 (375.9106) indicates an improvement in the parameter estimation. To generate plots for both the Euler and Hermite
1037 outcomes, type the following commands (see Figure 6)

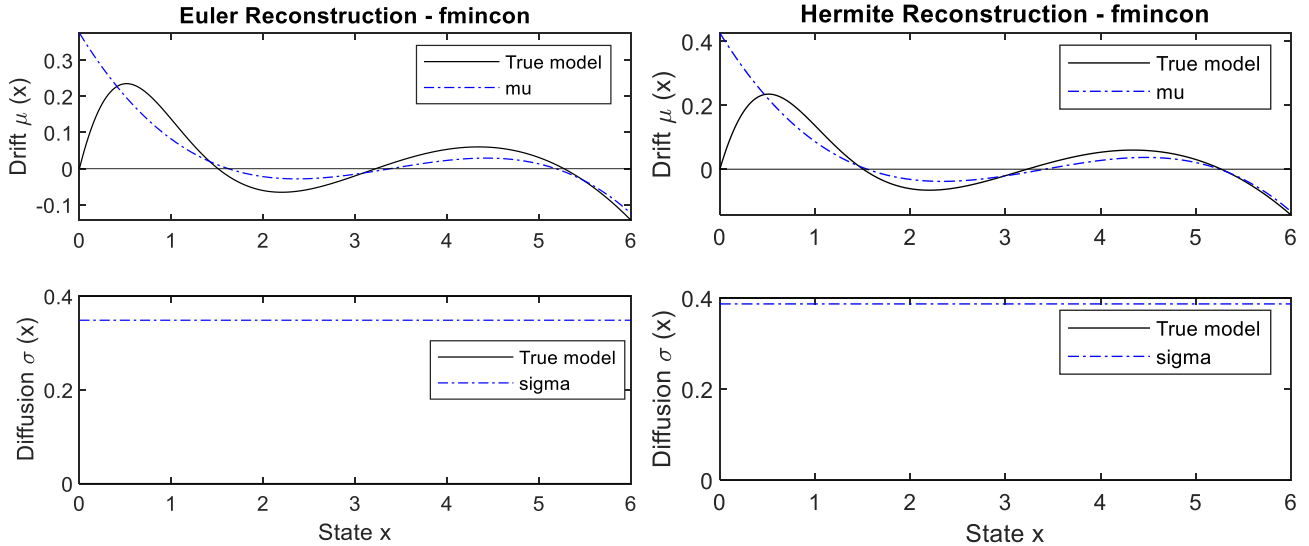
1038 result = result9; %plot for Euler outcomes

```

```

1039 S = load('MayData1D.mat');data = S.data;
1040 m = mean(data);s = std(data);
1041 result.s = s;result.m = m;result.par_est = result.estimated_par;
1042 result.mufun = @(x,par_est)s.*(par_est(1).*((x-m)./s).^3+par_est(2).*((x-
1043 m)./s).^2+par_est(3).*((x-m)./s)+par_est(4));result.sigmapfun =
1044 @(x,par_est)s.*par_est(5)+0.*x; %back-transformed diffusion
1045 mu = @(x,par)par(1).*x.*(1-x./par(2))-par(3).*x.^2./(par(4).^2+x.^2); %true drift
1046 sigma = @(x,par)par(5)+0.*x; %true diffusion
1047 par = zeros(5,1);par(1) = 1.01;par(2) = 10;par(3) = 2.75;par(4) = 1.6;par(5) = 0.4;
1048 %true model parameters
1049 xplot=linspace(0,6,2000);
1050 plot_results(result,xplot,mu(xplot,par),sigma(xplot,par));
1051
1052 result = result_her9; %plot for Hermite outcomes
1053 S = load('MayData1D.mat');data = S.data;
1054 m = mean(data);s = std(data);
1055 result.s = s;result.m = m;result.par_est = result.estimated_par;
1056 result.mufun = @(x,par_est)s.*(par_est(1).*((x-m)./s).^3+par_est(2).*((x-
1057 m)./s).^2+par_est(3).*((x-m)./s)+par_est(4));
1058 result.sigmapfun = @(x,par_est)s.*par_est(5)+0.*x;
1059 mu = @(x,par)par(1).*x.*(1-x./par(2))-par(3).*x.^2./(par(4).^2+x.^2);
1060 sigma = @(x,par)par(5)+0.*x;
1061 par = zeros(5,1);par(1) = 1.01;par(2) = 10;par(3) = 2.75;par(4) = 1.6;par(5) = 0.4;
1062 xplot=linspace(0,6,2000);
1063 plot_results(result,xplot,mu(xplot,par),sigma(xplot,par));
1064

```



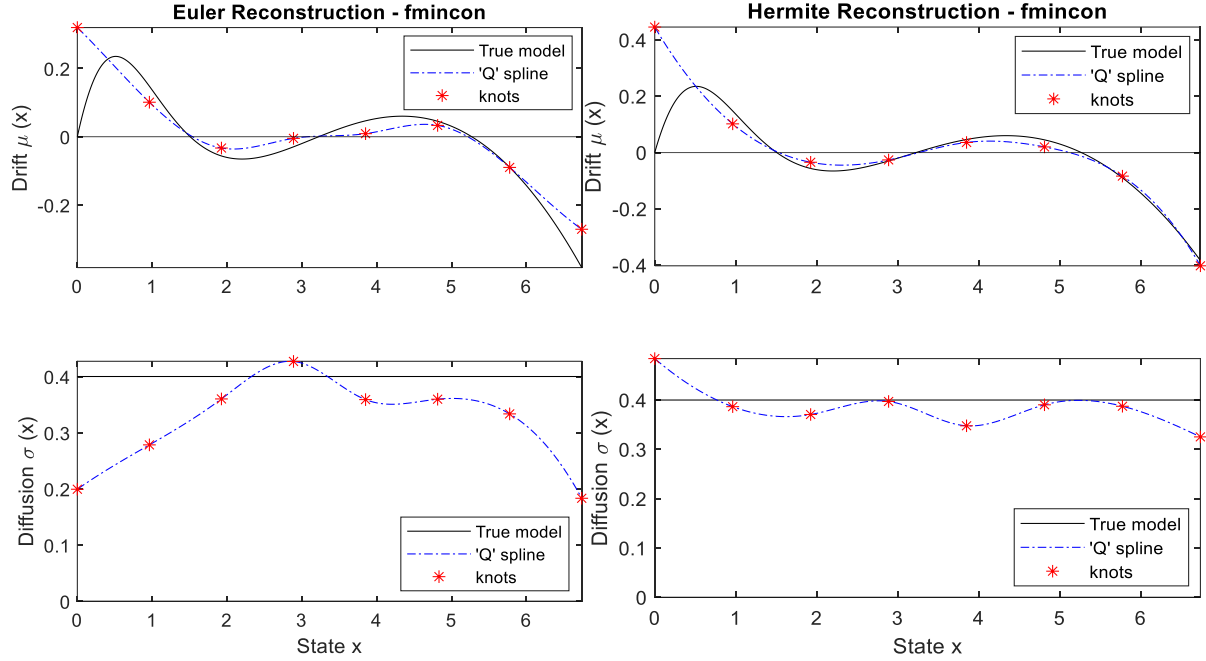
**Figure 6. Illustration of Euler and Hermite reconstructions applied to a low-resolution simulated dataset generated by a nonlinear model.** The left panel depicts the true model alongside the Euler reconstructed model, while the right panel depicts the true model alongside the Hermite reconstructed model. The data are simulated from the grazing model of May  $dx = \left\{ rx \left( 1 - \frac{x}{K} \right) - \frac{\gamma x^2}{x^2 + a^2} \right\} dt + \sigma dW$  with parameters  $r = 1.01, K = 10, \gamma = 2.75, a = 1.6, \sigma = 0.4$ . The original dataset has the time step of  $dt=0.01$  and contains  $T=3 \times 10^3$  data points. However, we select every 300<sup>th</sup> data point, resulting in a low-resolution dataset with a time step of  $dt=3$  and  $T=10^3$  data points.

As is evident from **Figure 6**, there is no a substantial improvement to the Euler reconstruction after applying Hermite reconstruction. The biggest improvement is observed for the diffusion function, while the improvement in the drift function is slight. We improve these outcomes using spline reconstruction in the next example.

```

1076 Example 11. Type the following command lines (we omit the details. See Example 9 for more details)
1077 S = load('MayData1D.mat');data = S.data;
1078 data = data(1:300:end);
1079 L = 0;R = max(data); %we consider the entire range of data
1080 data(data<L | data>R) = nan;
1081 dt = 3;
1082 mu = 8; sigma = 8;
1083 result10 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'QQ', 'L',
1084 L, 'R', R, ...,
1085 'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma) + 10,
1086 'solver', 'fmincon', 'search_agents', 5);
1087 legpoints = legitimate_points(data, dt, 'prev', result10, 'prev_range', 0.5, 'j', 3,
1088 'k', 12);
1089 result_her10 = hermite_reconstruction(data, dt, 'prev', legpoints, 'solver',...,
1090 'fmincon');
1091 these leads us to the following Euler parameter estimation
1092 Estimated parameters:
1093 0.31967 0.10093 -0.03319 -0.0049032 0.0088617 0.033111 -0.089981 -0.27115
1094 0.19942 0.27834 0.36011 0.42682 0.35892 0.35962 0.33349 0.18335
1095 - sum of log-likelihoods): 899.8273
1096 and the following Hermite parameter estimation
1097 Estimated parameters:
1098 0.38364 0.10187 -0.032008 -0.029461 0.036058 0.025406 -0.09752 -0.35229
1099 0.54784 0.38213 0.37651 0.40084 0.34381 0.39658 0.40582 0.29627
1100 - sum of log-likelihoods: 889.6237
1101
1102 To plot the results, type the following command lines (see Figure 10)
1103 mu = @(x,par)par(1).*x.*(1-x./par(2))-
1104 par(3).*x.^2./(x.^2+par(4).^2);sigma=@(x,par)par(5)+0.*x; % true model
1105 par = zeros(5,1);par(1) = 1.01;par(2)= 10;par(3) = 2.75;par(4) = 1.6;par(5) = 0.4;
1106 %true model parameters
1107 xplot = linspace(L,R,2000); % a dense mesh across the considered range
1108 plot_results(result10,xplot,mu(xplot,par),sigma(xplot,par)); % Euler & true models
1109 plot_results(result_her10,xplot,mu(xplot,par),sigma(xplot,par));%Hermite & true models
1110

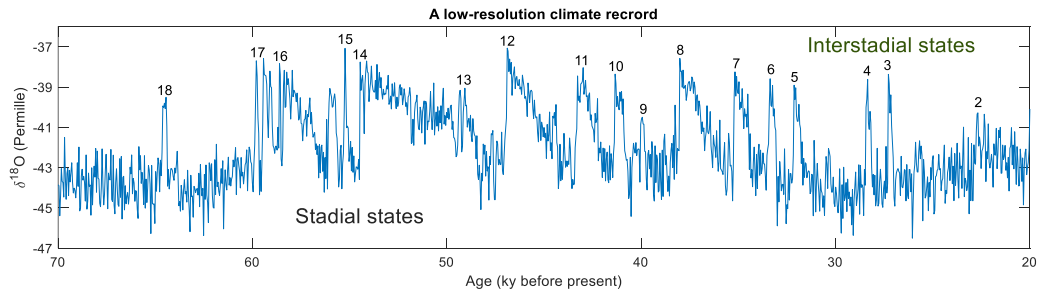
```



**Figure 7. Illustration of Euler and Hermite spline reconstructions applied to a low-resolution simulated dataset generated by a nonlinear model.** The left panel depicts the true model alongside the Euler reconstructed spline model, while the right panel depicts the true model alongside the Hermite reconstructed spline model. The data are simulated from the grazing model of May  $dx = \left\{ rx \left( 1 - \frac{x}{K} \right) - \frac{\gamma x^2}{x^2 + a^2} \right\} dt + \sigma dW$  with parameters  $r = 1.01, K = 10, \gamma = 2.75, a = 1.6, \sigma = 0.4$ . The original dataset has the time step of  $dt=0.01$  and contains  $T=3 \cdot 10^3$  data points. However, we select every 300<sup>th</sup> data point, resulting in a low-resolution dataset with a time step of  $dt=3$  and  $T=10^3$  data points. You can improve the Hermite reconstruction slightly further by opting for larger  $K$  values, like  $K = 10, 11, 12$ , albeit with a slightly increased computational time.

#### 15.4 Reconstructing an ice-core dataset

**Example 12.** In this case study, we reconstruct a  $\delta^{18}O$  record from the North Greenland Ice Core Project (NGRIP) (2004), which serves as a proxy for the temperature of the northern hemisphere. This record spans the last 120 thousand years, encompassing the last glaciation and has a resolution of 20 years. The dataset does not meet two data requirements. First, it is not stationary. Therefore, our analysis is restricted to the time period from 70 to 20 thousand years before the present time (see **Figure 8**) where the data is Markov. For more details on this, refer to (Kwasniok and Lohmann 2009) or see **Table 1**. Throughout the last glaciation, the climate of the northern hemisphere experienced alternating colder (stadial) and warmer (interstadial) states, due to a phenomenon called Dansgaard–Oeschger (DO) events (Dansgaard et al. 1993). Within the considered time window, the majority of DO events, DO2 to DO 18 out of 25 DO events, occurred (2004). Second, the dataset is not Markov but a rarified sample of data with every other data points, i.e., `data(1:2:end)`, is Markov (see subsection 6.3 for details). The relaxation time of this sample is  $63.1028 < 100$  (see **Table 1**) and therefore it falls in the category of low-resolution datasets. This sample comprises a total of 1217 data points and is a small dataset.



**Figure 8.** A  $\delta^{18}\text{O}$  climate record, with a resolution of 20 years, extending from 70 to 20 thousand years before the present time from NGRIP. This is used as a proxy for the temperature of the northern hemisphere which shows that the northern hemisphere climate alternated between cold stadial and warmer interstadial alternative climate states. In this time period majority of Dansgaard-Oeschger (DO) events, DO2 to DO18 out of 25 DO events, occurred (see the numbers).

Due to its low resolution, Hermite reconstruction is necessary. To reconstruct this dataset, use the following commands (for details take a look at previous examples).

```
data = readmatrix('NGRIP20.csv');
data = data(2649:5081);
data = data(1:2:end); % The original dataset is not Markov but a sample with every
other data point is
L = -45.5; R = -38.2;
data(data < L | data > R) = nan;
dt = 1; % Arbitrary
mu = 7; sigma = 7;
result11 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'QQ', 'L',
L, 'R', R, ...,
'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma) + 10,
'solver', 'fmincon', 'search_agents', 10); % we used 20 search agents
legpoints11 = legitimate_points(data, dt, 'prev', result11, 'prev_range', 0.5, 'j',
3, 'k', 9);
result_her11 = hermite_reconstruction(data, dt, 'prev', legpoints11, 'solver',
'fmincon');
```

you get the following Euler results

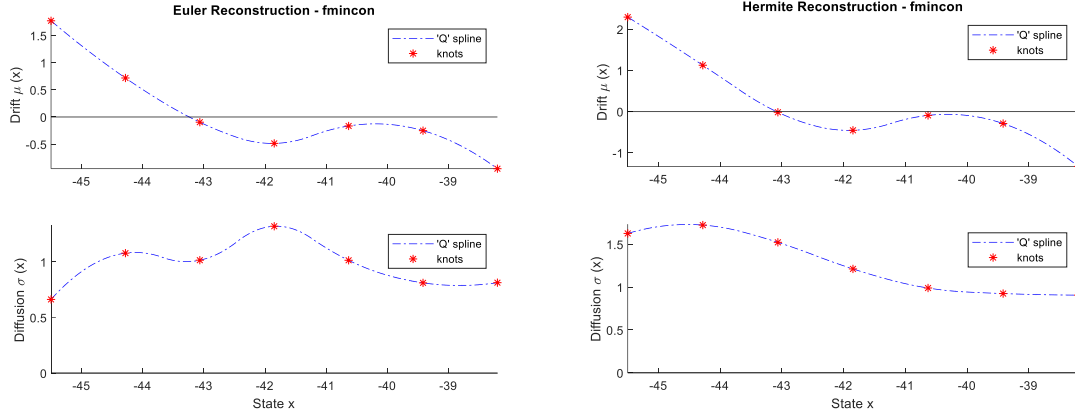
```
Estimated parameters :
2.5801    0.97206    0.0523    -0.50281    -0.18251    -0.23949    -1.1289
0.52863    1.079    1.0395    1.2871    1.0323    0.82494    0.78785
- sum of log-likelihoods) : 1721.4668
```

and the following Hermite results

```
Estimated parameters :
2.2923    1.1239    -0.017088    -0.45515    -0.093535    -0.29108    -1.3325
1.626    1.7236    1.5222    1.212    0.98925    0.92439    0.90594
- sum of log-likelihoods : 1651.6858
```

To depict the results type the following command lines (see **Figure** )

```
xplot = linspace(L,R,2000);
plot_results(result11,xplot);
plot_results(result_her11,xplot);
```



**Figure 9. Illustration of Euler and Hermite spline reconstructions applied to a low-resolution ice-core climate dataset.** The left panel depicts the Euler reconstructed spline model, while the right panel depicts the Hermite reconstructed spline model. The data are illustrated in **Figure 8**.

Some explanations are needed here. First, as mentioned, the entire dataset is not stationary, which is a data requirement in subsection 7.1. However, we have analyzed a portion of data that is stationary. If the analysis of the entire dataset is the goal, then the analysis performed on this portion should be repeated using a moving window approach. In this approach, one analyzes a short segment of data over a short window of time (which is often stationary), then shifts this time window slightly to the right and analyzes the second window, then shifts the second window slightly to the right and analyzes the third segment, so on. Eventually, you should get time-varying drift  $\mu(x, t)$  and diffusion  $\sigma(x, t)$  functions, which are calculated by interpolating the outcomes of all the segments. This is a simple scheme, and in practice, the size of segments does not need to be equal. Second, this dataset is not Markov and this is another data requirement we elaborated on in subsection 7.2. However, every other data point is nearly Markov, i.e., the ME time scale is 2 (see **Table 1**). Third, similar to **Example 8** we considered  $\mu = 7$ ,  $\sigma = 7$ . Since, this is a very small dataset with 1217 data points if you consider 8 knots it also works but it takes a bit more time for the package to find the legitimate points. But, we believe 7 knots should be enough.

## 15.5 The concept of effective potential

Does the climate dataset in the previous subsection have ‘alternative stable states’? To address this question, we often attempt to find the roots of the drift function, i.e., solve for  $\mu(x) = 0$ . Following this approach, we only identify a single equilibrium near -43, as shown in **Figure 10** (see the vertical dashed line in orange). However, this approach is incorrect. Such an approach is suitable for deterministic systems and stochastic systems with additive noise.

Here, we need to take a deeper look at this question. An equilibrium, by definition, is a state where the system settle into forever. Therefore, a stochastic system does not have any equilibrium since its trajectories never rest into any state. Instead, it can have ‘most visited states’ and ‘least visited states’, which are analogous to the concepts of ‘stable equilibria’ and ‘unstable equilibria’ (or, ‘repellers’) in deterministic systems. For simplicity, we henceforth apply the terminology of deterministic systems to stochastic systems we well. To this end, the following quantity called ‘effective potential’

$$U_{\text{eff}}(x) = -2 \left( \int^x \frac{\mu(u)}{\sigma^2(u)} du + \log \sigma(x) \right), \quad (3)$$

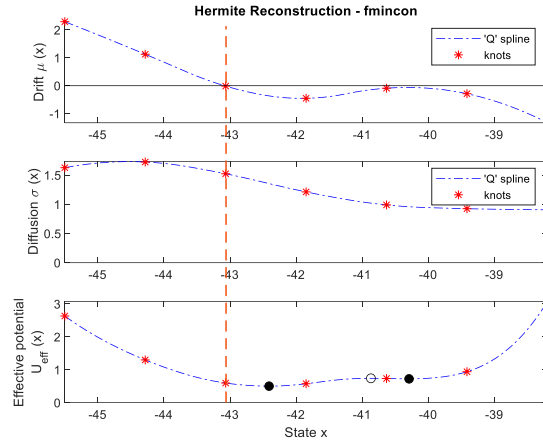
Should be calculated.  $U_{\text{eff}}(x)$  plays a role analogous to the concept of ‘stability landscape’ or ‘potential function’ in deterministic systems (refer to (Arani 2019, MS Arani et al. 2024) for further details).

Unlike the concept of stability landscape which is calculated solely based on the deterministic forces, i.e., the drift function  $\mu(x)$ , the concept of effective potential in **Error! Reference source not found.** incorporates information about the stochastic forces, i.e., the diffusion function  $\sigma(x)$ , too. To determine the equilibria, we need to identify the

minima and maxima of  $U_{\text{eff}}(x)$ , which correspond to stable and unstable equilibria, respectively. We will not delve into extra details here, as the package can perform these calculations. For instance, to plot the estimated drift, diffusion, and effective potential (**Figure 10**, bottom panel) for the Hermite reconstruction, type the following commands

```
xplot = linspace(L,R,2000);
plot_results(result_her11,xplot,'eff_potential');
```

and you get the following plot



**Figure 10. Illustration of the concept of effective and its significance in identifying alternative stable states in a stochastic system.** The top and middle panels are as in **Example 12**, right panel. In the bottom panel the effective potential is depicted, with its minima and maxima corresponding to the stable (solid black dots) and unstable (open circle) equilibria. Notably, the dashed vertical dashed orange line intersects the only root of the drift function (i.e., where  $\mu(x) = 0$ ) yet it does not coincide with any of the minima of the effective potential. This discrepancy suggests that relying solely on the drift function to calculate the equilibria is incorrect.

**Figure 10**, bottom panel, reveals that the climate dataset in **Example 12** possess alternative stable states (the solid circle in the bottom panel) separated by a repellor in between (open circle in the bottom panel). This example clarifies that the only way to determine the equilibria of a stochastic system is to identify the minima and maxima of the effective potential instead of finding the zeros crossings the drift function followed for deterministic systems.

## 16. Handling big datasets

When working with datasets containing millions of data points, the computational burden can be significant, leading us to consider using only a portion of the dataset. However, selecting the appropriate portion is crucial, as opting for the first 10%, last 10%, or middle portion can notably influence the final results, potentially introducing bias into the estimated parameters. Since Langevin models are Markovian, we can employ mini-batch optimization, where we sample a fraction of ‘data pairs’ and solve the underlying optimization problem based on that fraction alone. Here, a ‘data pair’ refers to any consecutive pair  $(x_t, x_{t+1})$  across the data. By randomly selecting a sample comprising just 10% of all data pairs, we can conduct the analysis on this subset. This fraction is well-mixed across the entire dataset and provides a representative sample. To ensure an even more random selection compared to simple random sampling, we recommend and implement a ‘stratified’ random sampling of data pairs. This method offers an excellent representation of the entire dataset. After a random sample of data is obtained, we can follow either of Euler or Hermite reconstruction as explained in previous sections.

**Example 13.** Consider the first dataset in **Example 1** which is simulated from the OU model. The length of this dataset is  $10^6$ . Imagine that we just wish to perform parametric reconstruction using 1% of this dataset. Type the following commands

```
S = load('OUdata1D.mat');data = S.data;
```



```

1247 dt = 0.01; % the time step remains unchanged (this should not be confused with data
1248 rarification)
1249 mu = @(x,par)par(1).*x;sigma = @(x,par)par(2);
1250 result16 = euler_reconstruction(data, dt, 'mu', mu, 'sigma', sigma, 'gradient_fun',
1251 eulergrad(mu, sigma), ...
1252 'reconst_fraction', [10 0.01], 'lb', [-200 eps], 'ub', [200
1253 200],'useparallel',true, 'solver', 'fmincon', 'search_agents', 5);
1254
1255 and you get an answer close to the following (depending on the sample you get)
1256
1257 Estimated parameters:
1258 -0.98133      1.0056
1259 - sum of log-likelihoods): -8780.5937

```

Some explanations here. In the name-value pair `'reconst_fraction', [10 0.01]`, we specify our intention to reconstruct a stratified random sample using 10 strata, representing just 1% of the entire data. Note that, here we consider stratification of time points (1,2, 3, ...) rather than that of data values. The accuracy of the estimated parameters is good (bearing in mind that the true solution is  $[-1,1]$ ). *It is crucial to realize that the resolutions of the mother dataset and its random sample.* Both datasets maintain the same resolution; only their lengths differ (which is why we set `dt = 0.01`, matching the resolution of the OU mother dataset). Given the high resolution of our data, further Hermite reconstruction may not significantly improve the results. Below are the complete code lines for implementing both Euler and Hermite reconstructions.

```

1267 S = load('OUdata1D.mat');data = S.data;
1268 dt = 0.01
1269 mu = @(x,par)par(1).*x;sigma = @(x,par)par(2);
1270 result16 = euler_reconstruction(data, dt, 'mu', mu, 'sigma', sigma, 'gradient_fun',
1271 eulergrad(mu, sigma), ...
1272 'reconst_fraction', [10 0.01], 'lb', [-200 eps], 'ub', [200
1273 200],'useparallel',true, 'solver', 'fmincon', 'search_agents', 5);
1274 legpoints16 = legitimate_points(data, dt, 'prev', result16, 'prev_range', 0.5, 'j',
1275 3, 'k', 4);
1276 result_her16 = hermite_reconstruction(data, dt, 'prev', legpoints16, 'solver',
1277 'fmincon');

```

## 1278 17. Handling replicate datasets

Before proceeding, ensure that you have added the path of the 'Burg' folder to your MATLAB working directory (refer to subsection 6.2 for detailed instructions).

In some cases, we may not have access to a single long dataset but rather to many shorter samples, known as *'replicate data'*. Reconstructing such data is not challenging as long as there is sufficient evidence or theoretical justification to believe that all the data share a common generating system. It's important to note that Langevin models are 'Markov' models, meaning that the future state, given the present state, is independent of the entire past history of states. Therefore, any damage to data at a single time point will only affect the adjacent data points, allowing us to effectively treat damaged values as missing values (NaN). Consequently, we can safely append a NaN at the end of each replicate and then concatenate all the replicates (the order of concatenation does not matter) to create a long dataset. The code `'prepare_replicateData.m'` automates this process. Once the replicate data is prepared, the subsequent calculations are straightforward. You can simply apply the same codes developed for 'typical' datasets (i.e., single time series datasets) to the replicate data. Note that the replicate data must be supplied as a cell array.

**Example 1.** In this analysis, we examine a dataset comprising three high-resolution replicates simulated from the grazing model of May, with parameters matching those in **Example 1**. Each replicate begins from the initial state  $x_0 = 8$  and continues until perturbations drive the system towards 0 biomass. To ensure the removal of transient effects, the



1294 first 5% of each replicate is discarded. Subsequently, we reconstruct this high-resolution replicate dataset using cubic  
 1295 splines as below

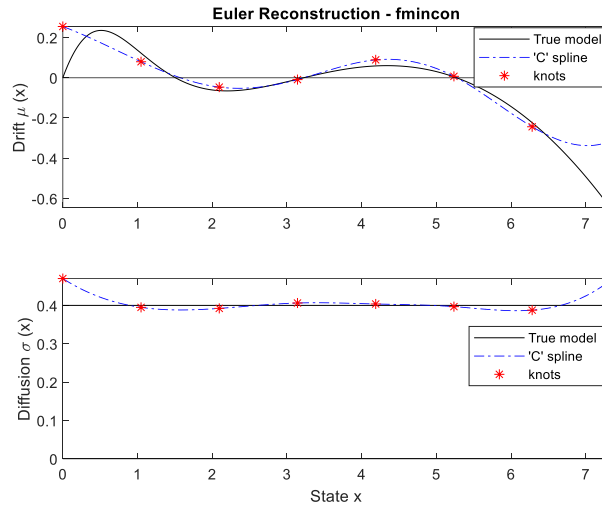
```
1296 S = load('MayData1D_Replicate.mat');
1297 data = S.data; % replicate data should be supplied as a cell array
1298 data = prepare_replicateData(data); % to reconstruct replicate data we first need to
1299 use this function. The rest of calculations are similar to those for typical datasets
1300 L = 0;
1301 R = max(data);
1302 dt = 0.1; % this is the true resolution of replicates
1303 mu = 8; sigma = 8;
1304 result17 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'CC', 'L',
1305 L, 'R', R, ...,
1306 'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma) + 10,
1307 'solver', 'fmincon', 'search_agents', 5);
```

1308 Which leads us to the following solution

```
1309 Estimated parameters:
1310 0.25467  0.079432 -0.046699 -0.0092472  0.088873  0.0069646 -0.24373 -0.31535
1311 0.47003  0.39488  0.3922  0.40622  0.40369  0.39681  0.38765  0.46541
1312 - sum of log-likelihoods): -9192.4963
```

1313 and for a plot type (see **Figure 14**)

```
1314 r=1.01;K=10;g=2.75;a=1.6;s=0.4; % true parameter values
1315 par = [r K g a s];
1316 mu = @(x,par)r.*x.*(1-x./K)-g.*x.^2./(x.^2+a.^2);sigma=@(x,par)s; % true model
1317 xplot=linspace(L,R,2000);
1318 plot_results(result17,xplot,mu(xplot,par),sigma(xplot,par));
```



1319  
 1320 **Figure 14. Reconstructing a replicate dataset.** The top and bottom panels depict the true drift and diffusion functions (solid black  
 1321 curves) alongside the estimated drift and diffusion functions (dot-dashed blue curves) obtained through spline modeling. The  
 1322 dataset consists of three replicates, all simulated from the May model and initiated at position  $x_0 = 8$  until reaching 0  
 1323 biomass. The true model parameters align with those described in **Example 1**.

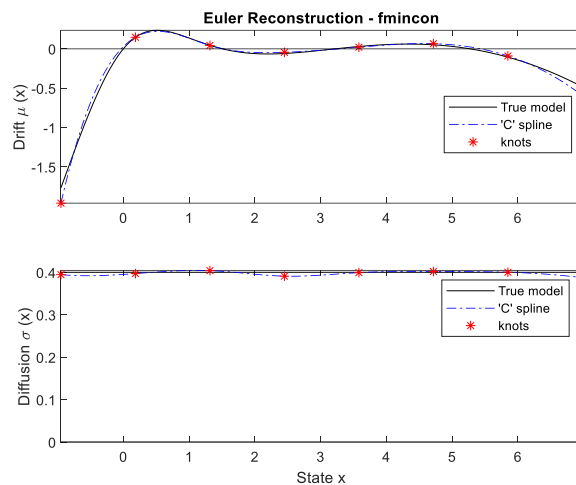
1324 *An explanation for the absence of the left tail in the reconstructed May model*

1325 As is evident in **Figure 3** and **Figure 7**, and **Figure 14** the left tail of the grazing model of May did not manifest in the  
 1326 reconstructed model. This discrepancy between the reconstructed and true models is not attributable to estimation

inaccuracies but rather to a deliberate modeling choice in the ecological context. In this model, our objective was to simulate a dataset with positive state values (biomass), despite the stochastic force potentially pushing trajectories into negative states. To address this, we implemented a ‘reflecting’ boundary at 0 biomass, effectively pushing trajectories back to positive values upon crossing 0. Consequently, the reconstructed model exhibits a steep positive rate of change at 0, in contrast to May’s deterministic model where trajectories slow down near 0 (note that 0 is an equilibrium in the determinist May model). As a result, no positive dataset can reflect this behavior, causing the left tail of May’s model to be omitted—a feature that holds limited ecological significance. To reconstruct the left tail, trajectories need to be allowed to cross 0 and fluctuate around it. To illustrate this concept, we have generated a replicate dataset ‘MayData\_LeftTail.mat’ which can reveal the left tail of May’s model. This dataset consists of 15 replicates, all initially placed at  $x_0 = 0.2$  and terminate once they escape the interval  $[-1, 7]$  via either of the left or right borders (or, they reach the chosen maximum length of  $2 \times 10^4$ ). To exclude transient effects the first 5% of all the replicates are discarded. This dataset has a high-resolution. Therefore, we apply Euler reconstruction and fit a cubic spline model. Type the following commands to recover the left tail in the May model (see **Figure 15**)

```
S = load('MayData_LeftTail.mat');
data = S.data; % replicate data should be supplied as a cell array
data = prepare_replicateData(data);
L = min(data); % note that here the min data value is negative (-0.9493)
R = max(data);
dt = 0.1; % This is the actual time step used to generate this dataset
mu = 8; sigma = 8;
result18 = euler_reconstruction(data, dt, 'nKnots', [mu sigma], 'spline', 'CC', 'L',
L, 'R', R, ...,
'lb', [zeros(1, mu) - 10, zeros(1, sigma)+eps], 'ub', zeros(1, mu + sigma) + 10,
'solver', 'fmincon', 'search_agents', 5);

r = 1.01; K = 10; g = 2.75; a = 1.6; s = 0.4; % true parameter values
par = [r K g a s];
mu = @(x,par)r.*x.*(1-x./K)-g.*x.^2./(x.^2+a.^2); sigma=@(x,par)s; % true model
xplot = linspace(L,R,2000);
plot_results(result18,xplot,mu(xplot,par),sigma(xplot,par));
```



**Figure 15. Revealing the left tail in the May model.** Estimated drift (top panel) and diffusion (bottom panel) functions using a dataset with negative values. In order to recover the left tail of the May model a dataset with negative biomass is needed (which is ecologically unrealistic). The dataset consists of three replicates all initialized at  $x_0 = 0.2$  and terminate once they escape the interval  $[-1, 7]$  via either of the left or right borders (or, they reach the chosen maximum length of  $2 \times 10^4$ ). The first 5% of all replicates are discarded to remove the transient effects.

## 18. Assessing the uncertainty of the results

After estimating the model parameters, we further need to have an estimate about the uncertainty of the estimated parameters. To this end, we need to use the code 'Uncertainty.m'. Here, we have made one single code which is responsible to calculate the uncertainty of the estimated model parameters for all different types of models (parametric, spline), different reconstruction schemes (Euler, Hermite), different data types (typical, replicate, big), and, with or without missing values. Here, we estimate the uncertainty of the parameters for several examples in this tutorial.

Consider **Example 1**. To assess the uncertainty of the parameters, type the following command lines

```
syms mu(x) sigma(x)
par = sym('par%d', [1 2]);
mu(x) = par(1)*x;sigma(x) = par(2);
S=load('OUdata1D.mat');data=S.data;dt=0.01;
ModelType=["Parametric" "Euler"];
h = 10^(-3);
estimated_par = [-1.0586 0.99531];
[hess,err_hess,err_par]=Uncertainty(ModelType,estimated_par,h,data,[],[],dt,par,mu,sigma);
```

and you get

Uncertainty of the parameters (in terms of standard deviation)

```
0.014136    0.00069961
```

## References

2004. High-resolution record of Northern Hemisphere climate extending into the last interglacial period. *Nature* **431**:147-151.
- Aït-Sahalia, Y. J. E. 2002. Maximum likelihood estimation of discretely sampled diffusions: a closed-form approximation approach. **70**:223-262.
- Arani, B. M., S. R. Carpenter, L. Lahti, E. H. Van Nes, and M. Scheffer. 2021. Exit time as a measure of ecological resilience. *Science* **372**:eaay4895.
- Arani, B. M. S. 2019. Inferring ecosystem states and quantifying their resilience: linking theories to ecological data. Wageningen University
- Bakshi, G., and N. J. T. J. o. B. Ju. 2005. A Refinement to Aït-Sahalia's (2002) "Maximum Likelihood Estimation of Discretely Sampled Diffusions: A Closed-Form Approximation Approach". **78**:2037-2052.
- Broersen, P. M. 2003. Automatic Time Series Identification Spectral Analysis with MATLAB Toolbox ARMASA. IFAC Proceedings Volumes **36**:1435-1440.
- Carpenter, S. R., B. M. Arani, P. C. Hanson, M. Scheffer, E. H. Stanley, and E. Van Nes. 2020. Stochastic dynamics of Cyanobacteria in long-term high-frequency observations of a eutrophic lake. *Limnology and Oceanography Letters* **5**:331-336.
- Dansgaard, W., S. J. Johnsen, H. B. Clausen, D. Dahl-Jensen, N. S. Gundestrup, C. U. Hammer, C. S. Hvidberg, J. P. Steffensen, A. Sveinbjörnsdottir, and J. Jouzel. 1993. Evidence for general instability of past climate from a 250-kyr ice-core record. *Nature* **364**:218-220.
- Dickey, D. A., and W. A. J. J. o. t. A. s. a. Fuller. 1979. Distribution of the estimators for autoregressive time series with a unit root. **74**:427-431.
- Einstein, A. J. I. o. t. t. o. t. B. m. 1905. On the movement of small particles suspended in a stationary liquid demanded by the molecular-kinetic theory of heat (English translation, 1956).
- Friedrich, R., J. Peinke, M. Sahimi, and M. R. R. J. P. R. Tabar. 2011. Approaching complexity by stochastic methods: From biological systems to turbulence. **506**:87-162.

1408 Honisch, C., and R. J. P. R. E. Friedrich. 2011. Estimation of Kramers-Moyal coefficients at low sampling  
 1409 rates. **83**:066701.  
 1410 Koziel, S., and L. Leifsson. 2013. Surrogate-based modeling and optimization. Springer.  
 1411 Kwasniok, F., and G. Lohmann. 2009. Deriving dynamical models from paleoclimatic records: Application  
 1412 to glacial millennial-scale climate variability. *Physical Review E* **80**:066104.  
 1413 Magnuson, J. J., S. R. Carpenter, and E. H. Stanley. 2023. North Temperate Lakes LTER: High Frequency  
 1414 Data: Meteorological, Dissolved Oxygen, Chlorophyll, Phycocyanin-Lake Mendota Buoy 2006-  
 1415 current.  
 1416 May, R. M. 1977. Thresholds and breakpoints in ecosystems with a multiplicity of stable states. *Nature*  
 1417 **269**:471-477.  
 1418 Mirjalili, S., S. M. Mirjalili, and A. J. A. i. e. s. Lewis. 2014. Grey wolf optimizer. **69**:46-61.  
 1419 MS Arani, B., S. R. Carpenter, E. H. van Nes, I. A. van de Leemput, C. Xu, P. G. Lind, and M. Scheffer. 2024.  
 1420 Stochastic regimes can hide the attractors in data, reconstruction algorithms can reveal them.  
 1421 bioRxiv:2024.2002. 2017.580797.  
 1422 Nadimi-Shahraki, M. H., S. Taghian, and S. Mirjalili. 2021. An improved grey wolf optimizer for solving  
 1423 engineering problems. *Expert Systems with Applications* **166**:113917.  
 1424 Rinn, P., P. G. Lind, M. Wächter, and J. J. a. p. a. Peinke. 2016. The Langevin Approach: An R Package for  
 1425 Modeling Markov Processes.  
 1426 Siegert, S., and R. J. P. R. E. Friedrich. 2001. Modeling of nonlinear Lévy processes by data analysis.  
 1427 **64**:041107.

1428  
 1429  
 1430  
 1431  
 1432