
Online Change Point Detection of Stochastic Differential Equations

Anonymous Authors¹

Abstract

Online detection of the change points for evolving stochastic differential equations (SDEs) from streaming observations is of paramount importance yet highly challenging. This is attributed to the inherent stochasticity in the data, compounded by the limitations imposed in an online setting, wherein training is conducted at each time step utilizing solely limited collected data. In this study, we introduce a method that integrates the Sparse Identification of Nonlinear Stochastic Dynamics (SINSDy) technique with a modified two-phase Follow-The-Regularized-Leader (FTRL)-proximal algorithm. This framework enables extracting low-dimensional features that enhance the accurate detection of change points in evolving SDEs by online identifying the governing equations. Our proposed methodology, termed Online Sparse Identification of Nonlinear Stochastic Dynamics (Online-SINSDy), not only demonstrates competitive performance in the identification of SDEs but also surpasses existing online change point detection methods and exhibits broader applicability. We show the efficacy of our framework using the dataset produced by representative real-world systems, such as the collective dynamics of the cichlid fish, the rolling bearing vibration data, and the electroencephalogram dataset.

1. Introduction

Data-driven change point detection of complex dynamical systems has drawn increasing attention in various disciplines, including neuroscience (Cribben et al., 2013), epidemiology (Yu et al., 2013), economics (Bodenham & Adams, 2017), and ecology (Francesco-Ficetola & Denoël, 2009). The ability to accurately detect the critical transitions in data patterns contributes significantly to facilitating

informed decision-making. The current mainstream change point detection methods can be broadly categorized into two types. One involves direct analysis of the data, while the other is based on identifying the underlying differential equations from data. Indeed, the data-driven identification of differential equations has emerged as a prominent topic across various domains (Brunton et al., 2016; Rudy et al., 2017; Boninsegna et al., 2018; Champion et al., 2019). Successfully identifying dynamical equations contributes to acquiring the mechanisms of underlying systems and assists in pinpointing change point locations.

Over the past few decades, various techniques have been proposed for the discovery of deterministic differential equations from data (Voss et al., 1999; Daniels & Nemenman, 2015; Lusch et al., 2018), with sparse identification of nonlinear dynamics (SINDy) (Brunton et al., 2016) emerging as one of the prominent methods. Given a set of observation data and a pre-selected library of basis functions, SINDy can automatically select the terms that best align with the underlying deterministic differential equations. However, in many real-world scenarios, complex systems are influenced by random noise, necessitating the use of stochastic differential equations (SDEs) for accurate modeling and analysis of these systems (Carpenter et al., 2020; Jhawar & Guttal, 2020; Jhawar et al., 2020). Therefore, several stochastic variants of SINDy-based methods have been introduced to unveil the underlying SDEs from data (Boninsegna et al., 2018; Wang et al., 2022; Huang et al., 2022), which utilize the Kramers-Moyal (KM) formulae (Risken & Risken, 1996) to estimate the drift and diffusion values and identify the drift and diffusion terms respectively. While effective in capturing the underlying SDEs, these methods are offline algorithms that require substantial amounts of data, rendering them unsuitable for streaming data. Recently, the development of O-SINDy addressed this limitation (Li et al., 2023b). It constructs a sparse regression problem at each time with only one sample and utilizes the follow-the-regularized-leader (FTRL)-proximal method (McMahan, 2011) to solve these sparse regression problems in an online fashion, allowing for sequential processing of samples in streaming data.

Besides, many methods and theories have been proposed to detect abrupt variations in time series data, which can be roughly divided into offline and online methods. The

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

offline algorithms analyze the complete dataset collectively, retrospectively examining the data to recognize when and where the change occurred (Aminikhahgahi & Cook, 2017; Truong et al., 2020; Yu et al., 2023). However, in real-world scenarios, what we often acquire is sequentially arriving streaming data, rendering these offline algorithms no longer applicable. Online change point detection algorithms can be further broadly categorized into methods based on statistical features and prediction errors. Bayesian Online Change-point Detection (BOCD) is a representative online statistical method that performs the detection by computing the probability distribution of the running length since the last change point (Adams & MacKay, 2007). Another typical category of statistical methods detects changes utilizing kernel statistics (Flynn & Yoo, 2019; Arlot et al., 2019; Wei & Xie, 2022; Ferrari et al., 2023; Li et al., 2023a) or the likelihood ratio of the probability distribution (Kawahara & Sugiyama, 2009; Liu et al., 2013) in sliding windows before and after the change point. On the other hand, O-SINDy and Temporal Change-point Detection (TCD) detect the occurrences of change points by measuring the sudden increase of the prediction loss (Hou et al., 2022). Specifically, O-SINDy calculates the predicted values by integrating the identified governing equations, while TCD generates prediction series by the Randomly Distributed Embedding framework (Ma et al., 2018). However, these prediction error-based methods are often not applicable to SDEs because we cannot precisely predict the specific state at the next moment. Additionally, many dynamical phenomena in real-world systems are **noise-induced** (Carpenter et al., 2020; Jhawar et al., 2020), rendering these methods designed for deterministic differential equations often ineffective.

Despite the methods mentioned above, to the best of our knowledge, there are currently limited approaches for simultaneous online identification and change point detection of SDEs. Therefore, in this article, we extend O-SINDy to the stochastic version by combining the stochastic-SINDy with a modified two-phase FTRL-proximal methodology. Additionally, we propose a novel method for detecting the change points in SDEs using the weight dynamics obtained during the identification process. Our proposal framework, named Online Sparse Identification of Nonlinear Stochastic Dynamics (Online-SINSDy), can simultaneously discover the governing SDEs and extract the low-dimensional (d -d) features that can be used to detect the change points from evolving SDEs in an online fashion. The advantages of the proposed framework are summarized as follows.

- Online-SINSDy can precisely identify the underlying SDEs using **limited streaming data** available during each training iteration. In contrast, most existing studies necessitate substantial data sets for concurrent training.

- The variations in weights acquired by Online-SINSDy during training serve as a natural indicator, enabling pinpointing the change points of the underlying SDEs. Its detection performance surpasses that of the baseline methods across numerous **real-world systems**.
- Online-SINSDy exhibits strong robustness to noise, achieving small detection errors and high detection accuracy in identifying system change points across various noise levels.
- Online-SINSDy can concurrently capture the changes in both the **deterministic** and **stochastic** parts of the underlying SDEs. Hence, the applicability of our method is more extensive.

2. Methods

Suppose we receive a few numbers of streaming time series from the same underlying system via a finite number of sensors. Let $\mathbf{x}^k(t_l) \in \mathbb{R}^d$ denote the d -d state of the k th time series at time t_l , where $k = 1, 2, \dots, M$ represent different trajectories collected by different sensors, $t_1 < t_2 < \dots < t_N$ are discretely sampled times. Since the time series arrive continuously as streams, it is assumed that due to the limited memory capacity, only data within a narrow time window of duration T_{tr} preceding the current time can be stored and utilized for training. Let $\mathbf{X}(t_l) \in \mathbb{R}^{MT_{\text{tr}} \times d}$ denote the state matrix, encompassing all available data applicable for training purposes at time t_l .

In this work, we consider dynamics driven by an Ito SDE:

$$d\mathbf{x}(t) = \mathbf{f}[\mathbf{x}(t)]dt + \mathbf{H}[\mathbf{x}(t)]d\mathbf{W}(t),$$

where $\mathbf{f}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the drift function, $\mathbf{H}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ is the diffusion function, and $\mathbf{W}(t)$ is a d -d standard Brown motion. We further define the positive-definite covariance matrix of the diffusion as $\mathbf{G}(\mathbf{x}) = \mathbf{H}^\top(\mathbf{x})\mathbf{H}(\mathbf{x})$.

2.1. Regression problems based on the KM formulae

In most scenarios, we cannot know the specific forms of the drift and diffusion terms. However, based on the KM formulae, the drift and diffusion values can be estimated as

$$\begin{aligned} f_i(\mathbf{x}) &= \lim_{\Delta t \rightarrow 0} \mathbb{E} \left\{ \frac{1}{\Delta t} [x_i(t + \Delta t) - x_i(t)] | \mathbf{x}(t) = \mathbf{x} \right\}, \\ G_{ij}(\mathbf{x}) &= \lim_{\Delta t \rightarrow 0} \mathbb{E} \left\{ \frac{1}{\Delta t} [x_i(t + \Delta t) - x_i(t)][x_j(t + \Delta t) \right. \\ &\quad \left. - x_j(t)] | \mathbf{x}(t) = \mathbf{x} \right\}, \end{aligned} \quad (1)$$

where x_i represents the i -th component of the d -d state of the time series, $f_i(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes the i -th component of the d -d vector-valued drift function $\mathbf{f}(\cdot)$ and $G_{ij}(\cdot) :$

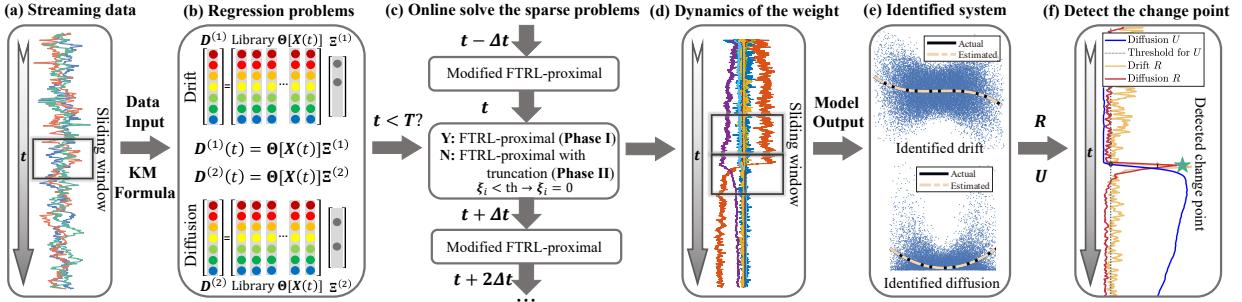


Figure 1. Schematic diagram for Online-SINSDy.

$\mathbb{R}^d \rightarrow \mathbb{R}$ refers to the component located in the i -th row and j -th column of the matrix-valued diffusion function $G(\cdot)$.

A conventional technique for estimating KM coefficients is the binning method (Boninsegna et al., 2018). We first partition the observed states into Q equally spaced intervals. Then the states within each bin are used to estimate the drift and diffusion values at the midpoint of each interval. The details of the binning method can be found in Appendix A. However, in scenarios where there is a limited amount of training data, we can directly utilize all available data to calculate the linear and quadratic variation at each point as

$$\begin{aligned} \mathbf{D}_i^{(1)}(t_l) &= \left\{ \frac{1}{\Delta t} [\mathbf{x}_i^q(t_{l-p+1}) - \mathbf{x}_i^q(t_{l-p})] \right\}_{\substack{q=1,2,\dots,M \\ p=1,2,\dots,T_{\text{tr}}}}, \\ \mathbf{D}_{ij}^{(2)}(t_l) &= \left\{ \frac{1}{\Delta t} [\mathbf{x}_i^q(t_{l-p+1}) - \mathbf{x}_i^q(t_{l-p})] \right. \\ &\quad \left. [\mathbf{x}_j^q(t_{l-p+1}) - \mathbf{x}_j^q(t_{l-p})] \right\}_{\substack{q=1,2,\dots,M \\ p=1,2,\dots,T_{\text{tr}}}}, \end{aligned} \quad (2)$$

where $\mathbf{D}_i^{(1)}$ represents the estimated values of the i -th component of the d -d drift vectors at different states, $\mathbf{D}_{ij}^{(2)}$ refers to the estimated values of the component located in the i -th row and j -th column of the $(d \times d)$ -d matrix-valued diffusion matrices at different states, and $x_i^q(t_l)$ denotes the i -th component of the d -d state at time t_l in the q -th trajectory. In this paper, the time stamps are assumed to be uniformly sampled and the time interval is denoted as Δt , i.e. $\forall l : t_{l+1} - t_l = \Delta t$. We construct the regression problems without the need for bin averaging (Huang et al., 2022).

We then construct a candidate library $\Theta[\mathbf{X}(t_l)] \in \mathbb{R}^{MT_{\text{tr}} \times P}$ consisting of P pre-selected basis functions. For example, the P -order polynomial basis

$$\Theta[\mathbf{X}(t_l)] = \begin{bmatrix} | & | & & \cdots & | \\ 1 & \mathbf{X}(t_l) & \cdots & \mathbf{X}^P(t_l) & | \end{bmatrix}^\top,$$

where the high-order terms of \mathbf{X} include all mixed-terms, for example, when $P = 2$ and $d = 2$, the library is

$\Theta(\mathbf{X}) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2]$. Without loss of generality, we focus on the case where the observed states contain only one space dimension ($d = 1$), allowing us to omit all subscripts in Eq. (2). The method we propose below can be easily generalized to the multivariate case.

Suppose the basis functions in the library encompass all the function terms present in the underlying SDE. We can determine the weight of basis functions corresponding to the drift $\Xi^{(1)}$ and the diffusion $\Xi^{(2)}$ at t_l by solving the following regression problems:

$$\begin{aligned} \mathbf{D}^{(1)}(t_l) &= \Theta[\mathbf{X}(t_l)]\Xi^{(1)}, \\ \mathbf{D}^{(2)}(t_l) &= \Theta[\mathbf{X}(t_l)]\Xi^{(2)}. \end{aligned} \quad (3)$$

2.2. Modified two-phase FTRL-proximal algorithm

An essential observation lies in the fact that numerous practical biophysical SDEs exhibit sparsity in their drift and diffusion expressions, characterized by the presence of only a few terms within the space spanned by basis functions. Additionally, when dealing with streaming data, it is imperative to iteratively optimize the weight vectors over time. Here, we propose a modified two-phase FTRL-proximal algorithm, capable of solving the regression problems online with high accuracy while ensuring the sparsity of the solution.

We define the regression loss of Eq. (3) at time t_l as

$$\begin{aligned} L_{t_l}^{(1)}(\Xi^{(1)}) &= \frac{1}{MT_w} \|\mathbf{D}^{(1)}(t_l) - \Theta[\mathbf{X}(t_l)]\Xi^{(1)}\|_2^2, \\ L_{t_l}^{(2)}(\Xi^{(2)}) &= \frac{1}{MT_w} \|\mathbf{D}^{(2)}(t_l) - \Theta[\mathbf{X}(t_l)]\Xi^{(2)}\|_2^2, \end{aligned}$$

where $\mathbf{D}^{(1)}(t_l)$ and $\mathbf{D}^{(2)}(t_l)$ are the MT_{tr} -d vectors of the estimated target values of the drift and diffusion terms at different states calculated by Eq. (2).

The FTRL-proximal algorithm constructs the following unconstrained optimization problems to address the online

165 sparse regression tasks presented in Eq. (3):

$$\begin{aligned} \Xi(t_{l+1}) &= \arg \min_{\Xi} \sum_{i=1}^l \nabla L_{t_i}[\Xi(t_i)]^\top \Xi + \lambda_1 \|\Xi\|_1 \\ &\quad + \frac{1}{2} \lambda_2 \|\Xi\|_2^2 + \frac{1}{2} \sum_{i=1}^l \sigma_i \|\Xi - \Xi(t_i)\|_2^2, \end{aligned} \quad (4)$$

173 where $\Xi(t_l)$ can denotes either the weight vector of drift
 174 $\Xi^{(1)}$ or the weight vector of diffusion $\Xi^{(2)}$ at time t_l of the
 175 training phase, $\lambda_1 > 0$ and $\lambda_2 > 0$ are the coefficients of the
 176 L_1 and L_2 regularization, whereas $\|\Xi - \Xi(t_i)\|_2^2$ denotes
 177 the proximal term that center additional regularization at the
 178 current weight, which can prevent dramatic change of the
 179 weight vector during the training process, $\sigma_i > 0$ serves as
 180 the learning rate. In fact, the joint action of the first term
 181 and the proximal term is equivalent to training the weight
 182 vectors using online gradient descent method (McMahan,
 183 2011). We provide a detailed explanation of this part in
 184 Appendix E.

185 Note that we can obtain the following closed-form solution
 186 of Eq. (4):

$$\xi_j(t_{l+1}) = \begin{cases} 0, & \text{if } |\zeta_j(t_l)| < \lambda_1, \\ \frac{-\zeta_j(t_l) - \text{sgn}[\zeta_j(t_l)]\lambda_1}{\lambda_2 + \sum_{i=1}^l \sigma_i}, & \text{if } |\zeta_j(t_l)| \geq \lambda_1, \end{cases}$$

193 where $\xi_j(t_l)$ denotes the j th component of $\Xi(t_l)$, $\zeta_j(t_l)$
 194 represents the j th component of $\zeta(t_l) = \sum_{i=1}^l \{\nabla L_{t_i}[\Xi(t_i)] -$
 195 $\sigma_i \Xi(t_i)\}$, and $\text{sgn}(\cdot)$ is the sign function. The detailed
 196 derivation process is provided in Appendix B.

197 While the above algorithm incorporates regularization terms
 198 to promote sparsity, the learned function expressions often
 199 still contain some minor non-zero erroneous terms. Previous
 200 approaches such as O-SINDy address this issue by applying
 201 a truncation operation after the entire training process, forcing
 202 the parameters below a threshold to zero. However, for
 203 the task of identifying SDEs, we observe that weak sparsity
 204 induced by the regularization terms in Eq. (4) can lead to
 205 convergence towards a false surrogate function expression,
 206 resulting in the identification task failure.

207 To overcome this challenge, we propose a two-phase training
 208 method. In Phase I, we update the weight vectors based
 209 on the original FTRL-proximal algorithm. In Phase II, a
 210 truncation operation is performed after each weight update,
 211 any scalar component below a pre-selected threshold is set
 212 to zero. These components are then excluded from further
 213 updates during the subsequent training process. This
 214 approach ensures the sparsity of the parameters and disrupts
 215 the local stable state of the weight vectors induced by
 216 interference terms, facilitating convergence to the authentic
 217 function expressions for the drift and diffusion terms.

2.3. Change point detection based on the weight vectors

During the identification process, we dynamically adjust the weight vectors using streaming data, thereby obtaining the optimal solution to the optimization problem Eq. (4) at each moment. During the training process, if the underlying SDE remains unchanged, the weight vectors will converge after a transition period and show minor fluctuations around their true values. However, when the dynamics experience specific changes at a particular time step, the weight vectors will undergo more substantial variations and gradually converge toward a new stable state corresponding to the altered SDE. Thus, detecting change points can be achieved by identifying abrupt alterations in the weight vectors.

Nevertheless, the occurrence of rare events in the evolution of SDEs can significantly impact the single-step changes of the weight vector, leading to erroneous change point detection. To address these problems, a sliding window of length T_{cpd} is introduced. Change points are detected by quantifying the variation of the average weight vectors between two adjacent time windows. Specifically, we first calculate the average weight vector in $[t_l - T_{\text{cpd}}, t_l]$:

$$\bar{\Xi}(t_l - T_{\text{cpd}}, t_l) = \frac{1}{T_{\text{cpd}}} \sum_{i=l-T_{\text{cpd}}}^{l-1} \Xi(t_i).$$

Then, we define the relative change intensity between two consecutive average weight vectors at the time t_l as

$$R(t_l) = \frac{\|\bar{\Xi}(t_l - T_{\text{cpd}}, t_l) - \bar{\Xi}(t_l, t_l + T_{\text{cpd}})\|_2}{\|\bar{\Xi}(t_l - T_{\text{cpd}}, t_l)\|_2 + \|\bar{\Xi}(t_l, t_l + T_{\text{cpd}})\|_2},$$

which serves as an indicator for identifying the change points. If the SDE remains unchanged, R will fluctuate around a relatively low value. However, when a change point falls within the sliding window, R will experience a significant increase. As the sliding window moves past the change point, R will initially ascend and subsequently descend, exhibiting an unimodal pattern. When the intersection of the two adjacent sliding windows is close to the change point, the change intensity of the average weight vector between the preceding and subsequent windows will be maximized and the value of R will reach its peak. Consequently, we can pinpoint the most likely change point locations by identifying the maximum value of R in the unimodal pattern.

During the online detection process, we still need a method to automatically detect when the indicator starts to exhibit significant changes, alerting us to pay attention to the potential occurrence of a change point during that period. Here, we employ the cumulative sum (CUSUM) control chart to detect the moment when there is an incremental change in the mean of the indicator (Page, 1954). We denote this moment as t_{CU} , and the indicator value at this moment as

R_{CU}. Starting from *t_{CU}*, we record the process of the indicator values rising and then falling until it returns to *R_{CU}* after a short time. We identify the time corresponding to the maximum value of the indicator during this period as the determined change point moment. Given a sequence of the indicator $R(t_1), R(t_2), \dots, R(t_n)$, we first estimate the initial average m_R and standard deviation σ_R from the first k samples of the indicator. Then we can recursively define the upper cumulative process sums using:

$$U(t_l) = \begin{cases} 0, & \text{if } l = 1, \\ \max[0, U(t_{l-1}) + R(t_l)] - m_R - \frac{1}{2}h\sigma_R, & \text{if } l > 1, \end{cases} \quad (5)$$

where h denotes the number of standard deviations from the target mean that makes a shift detectable. The indicator violates the CUSUM criterion for the first time at *t_{CU}* if it obeys $U(t_{CU}) > \gamma\sigma_R$ and $U(t_l) \leq \gamma\sigma_R (\forall t_l < t_{CU})$, where the control limit γ is an artificially selected threshold.

Fig. 1 illustrates the complete schematic diagram of Online-SINSDy. The pseudocode of our method is provided in Appendix C for detailed reference.

3. Results

In this section, we validate the effectiveness of our Online-SINSDy algorithm for both identification and change point detection tasks using various synthetic data and real-world complex systems with randomness.

For the identification task, we validate the outstanding performance of Online-SINSDy on a synthetic double-well system, a recorded dataset of collective dynamics of the cichlid fish, and a rolling bearing vibration dataset. The experimental results demonstrate that our method can accurately identify the underlying SDEs even in real-world scenarios with significant observation noise and missing data. Due to limitations in space, we present the system identification results in Appendix D.

For the change point detection task, to validate the effectiveness of our method, we extensively test it to detect the splicing points of two synthetic systems with identical steady-state distribution, the moment when the group size of the cichlid fish experiences a change, the occurrence of the bearing fault, and the occurrence of epilepsy. To evaluate the robustness of our method, we also introduce varying levels of Gaussian random noise into the observational data.

To further evaluate the effectiveness and robustness of our method, we compare it with several baseline methods. Specifically, we compare our model with the following online change point detection methods: (a) Bayesian Online Changepoint Detection (BOCD); (b) Online Kernel Cumulative Sum (KCUSUM) (Flynn & Yoo, 2019); (c) Online Density Ratio Estimation (DRE) (Kawahara & Sugiyama,

2009); (d) No-prior-knowledge Exponential Weighted Moving Average algorithm (NEWMA) (Keriven et al., 2020). Table 1 demonstrates the diverse metrics that different methods utilize for accomplishing the detection task.

3.1. Detect the splicing points of synthetic systems

First, we consider splicing trajectories of two different synthetic systems whose dynamics are in the form of the following two SDEs:

$$dx = (x - x^3)dt + \sqrt{(2 + 2x^2)}dW, \quad (6)$$

$$dx = -xdt + \sqrt{2}dW. \quad (7)$$

Although these two SDEs are completely different, the steady-state distribution of x in Eq. (6) and Eq. (7) is identical (Nabeel et al., 2022), as is shown by the black dashed line in Fig. 2(c). Thus, it is quite challenging to detect the time points when the underlying SDE translates from one to another according solely to the statistical features.

Here we generate $M = 100$ trajectories, each containing $N_{tr} = 1 \times 10^5$ recorded data points with the time interval $\Delta t = 0.01$. The first half of the time series is generated by Eq. (6), whereas the second half follows Eq. (7). The library is set as $\Theta(x) = [1, x, \dots, x^9]$. During the training process, we take the sliding window length $T_{tr} = 10$ with stride step $\Delta T = 10$. Phase II begins when $t > 200$ with the threshold for the drift and diffusion takes $th^{(1)} = 0.1$ and $th^{(2)} = 0.01$ respectively. The sliding window for change point detection is set as $T_{cpd} = 100$. To calculate the upper CUSUM, in all following experiments, we utilize the first 1000 recorded indicators in Phase II to estimate the average m_R and standard deviation σ_R and set $h = 1$ in Eq. (5).

For the sake of convenience, we utilize detection accuracy and detection error to assess the performance of our method in change-point detection tasks. The detection error ε refers to the time difference between the actual occurrence time of a change point T_{true} and the time at which we detect the occurrence of that change point T_{detect} , i.e. $\varepsilon = |T_{true} - T_{detect}|$. The detection accuracy pertains to the probability of the algorithm identifying change points with the detection error below an acceptable threshold ε_{th} . In real-world online change point detection tasks, the objective is to detect the change points promptly following their occurrence. If the detection error ε is too large, it often loses its significance. Hence, we establish a maximum allowable threshold for the error and the detection with error below this threshold is deemed accurate, i.e. $\varepsilon < \varepsilon_{th}$. We note that many articles directly employ detection delay (the positive part of the stopping time minus the true change-point time) to measure the effectiveness of the detection. In Appendix F, we also include statistical figures corresponding to the detection delay under different noise levels.

Table 1. Comparison of different change point detection methods with diverse metrics they rely on for accomplishing the detection task.

Method	Statistics	Prediction	Identification (Deterministic)	Identification (Stochastic)
O-SINDy (Li et al., 2023b)	✗	✓	✓	✗
TCD (Hou et al., 2022)	✗	✓	✗	✗
RC-TCD (Li et al., 2023a)	✗	✓	✗	✗
BOCD (Adams & MacKay, 2007)	✓	✗	✗	✗
DRE (Kawahara & Sugiyama, 2009)	✓	✗	✗	✗
NEWMA (Keriven et al., 2020)	✓	✗	✗	✗
KCUSUM (Flynn & Yoo, 2019)	✓	✗	✗	✗
Online-SINSDy (ours)	✗	✗	✓	✓

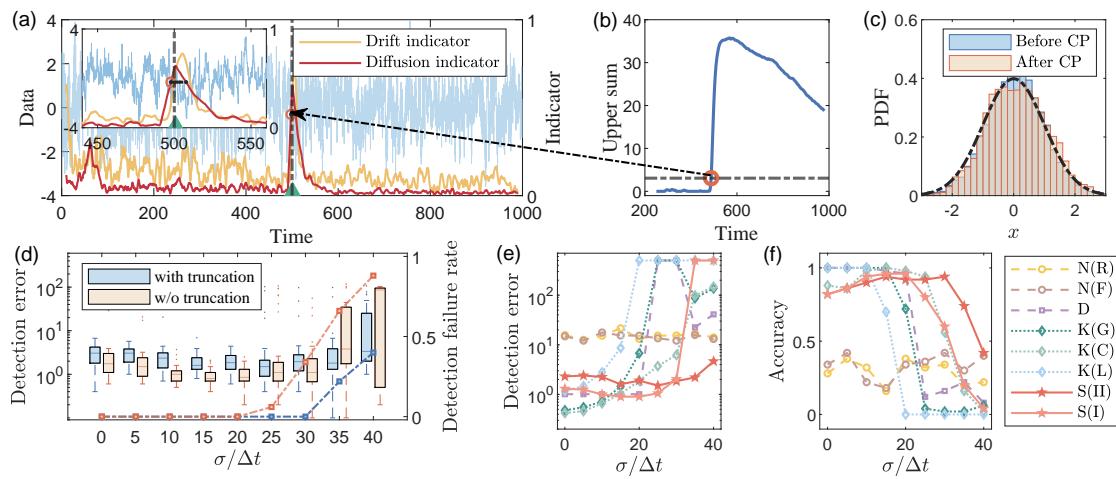


Figure 2. Change point detection results of the synthetic systems with identical steady-state distribution. (a) The streaming data and Online-SINSDy indicators. (b) The upper CUSUM indicator with the control limit $\gamma = 500$. (c) Two separate statistical histograms plotted using data before and after the splicing point, with a black dashed line representing their shared probability density. (d) Box plot of detection errors of Online-SINSDy under different levels of noise. (e-f) The detection errors and accuracy under different noise levels obtained by Online-SINSDy with (S(II)) and without (S(I)) truncation in Phase II, NEWMA with Random Fourier Features (N(R)) and Fastfood approximation (N(F)), DRE (D), KCUSUM with Gaussian (K(G)), Cauchy (K(C)) and Laplace (K(L)) kernels.

Fig. 2(a-b) illustrates the relative change intensity indicator R and the CUSUM indicator U during the training process. As shown in Fig. 2(d), our framework can sensitively detect the change points under varying levels of the noise standard deviation σ . For relatively low noise intensities, Online-SINSDy can accurately detect the existence of the change points and determine their positions with small detection errors, regardless of whether truncation is performed in Phase II. As the noise intensity increases, Online-SINSDy with truncation demonstrates remarkable robustness to the noise, while the algorithm without truncation rapidly loses its detection capability, particularly when $\sigma > 20\Delta t$.

To comprehensively assess the efficacy of our change point detection framework, we conduct a comparative analysis, pitting Online-SINSDy against various baseline methods across diverse noise levels. Each method is subjected to 50 repetitions under each noise intensity, and the median within each group is selected to represent its detection error. Additionally, we define detection accuracy as the proportion of experiments where the detection error was less than 10.

Due to the absence of notable changes in the probability distribution of the data before and after the change point, BOCD fails to detect the presence of the change points. Illustrated in Fig. 2(e-d), KCUSUM with distinct kernels and DRE exhibit smaller detection errors and higher detection accuracy under low noise intensities. However, their robustness to noise is poor, swiftly diminishing their detection capability as noise levels rise. Conversely, NEWMA maintains relatively stable detection errors and accuracy as noise intensity increases. Nonetheless, their detection errors consistently remain at a high level, and the detection accuracy persistently hovers at a low level. In comparison, Online-SINSDy showcases robustness to noise in terms of both detection errors and accuracy, surpassing other baseline algorithms, particularly in scenarios characterized by heightened noise intensity.

3.2. Detect the change of the fish populations

Our framework can also detect the change points of **real-world systems** with randomness, which is quite challenging because of the strong observation noise, the inherent dynamical noise, and the existence of missing data in real scenarios. Here we focus on the dynamics of collective alignment in groups of the cichlid fish *Etroplus suratensis* with different group sizes. Previous work has confirmed that the schooling of fish, which means highly polarized and coherent motion, is a **noise-induced** effect, and the mesoscopic dynamics can be fitted as an SDE (Jhawar et al., 2020):

$$d\mathbf{m} = -\gamma_1 \mathbf{m} dt + \left[\frac{\gamma_2(1 - |\mathbf{m}|^2) + \gamma_1}{N_f} \right]^{1/2} \mathbb{1} \cdot d\mathbf{W}, \quad (8)$$

where $\gamma_1 \approx -0.1$ and $\gamma_2 \approx 4.0$ are constant parameters, N_f denotes the group size, $\mathbb{1}$ is the 2-d identity matrix, and $\mathbf{m} = [m_x, m_y]^\top$ is the 2-d group polarization calculated by

$$\mathbf{m}(t) = \frac{1}{N_f} \sum_{i=1}^{N_f} \frac{\mathbf{v}_i(t)}{|\mathbf{v}_i(t)|}, \quad (9)$$

which represents the collective fish motion direction and the degree of alignment. $\mathbf{v}_i(t)$ in Eq. (9) is the 2-d velocity of the i th fish at time t .

Eq. (8) indicates that the diverse dynamical behaviors exhibited by fish with different group sizes are induced by varying intensities of intrinsic noise. Consequently, we can utilize the data of $\mathbf{m}(t)$ to discern instances when there are changes in the group size. The dataset was constructed by concatenating the trajectories of $\mathbf{m}(t)$ corresponding to $N_f = 15$ and $N_f = 60$, as is shown in Fig. 3(a). The total length of the concatenated data is 4800s with the sampling interval $\Delta t = 0.12s$. The library is set as $\Theta(x) = [1, x, \dots, x^9]$. During the training process, we take the window length $T_{\text{tr}} = 100$ with the stride step $\Delta T = 10$. Phase II begins at $t = 960s$ with the truncation threshold $\text{th}^{(1)} = \text{th}^{(2)} = 0.01$. Notably, all data associated with the splicing points is excluded during the training process. The sliding window for change point detection is set as $T_{\text{cpd}} = 40$.

We find that the dynamics associated with the drift indicators obtained through Online-SINSDy do not manifest a discernible alteration after the change point. This observation suggests that the change in the population does not influence the drift term, aligning with Eq. (8). Consequently, we turn to the indicators corresponding to the two diagonal elements of the diffusion matrix, denoted as G_{xx} and G_{yy} , to identify the change points. Illustrated in Fig. 3(d-e), both indicators promptly identify the moment when changes occur in the population, with the detection error and accuracy displaying robustness across varying noise levels. However, eliminating the truncation in Phase II leads to a notable increase in both detection error and failure rate when the noise intensity $\sigma > 0.5\Delta t$.

According to the results in Fig. 3(b-c), it is challenging for conventional methods to detect the changes in the population. We find that BOCD, DRE, and KCUSUM with Cauchy kernel are unable to identify the changes. NEWMA and KCUSUM with Gaussian and Laplace kernels take a considerable amount of time after the changes have occurred to recognize the variations in the system. In contrast, our approach exhibits significantly lower detection errors and substantially higher detection accuracy compared to the baseline models across diverse noise intensities. Here, we define the detection accuracy as the proportion with a detection error less than 120s in 50 repetitions.

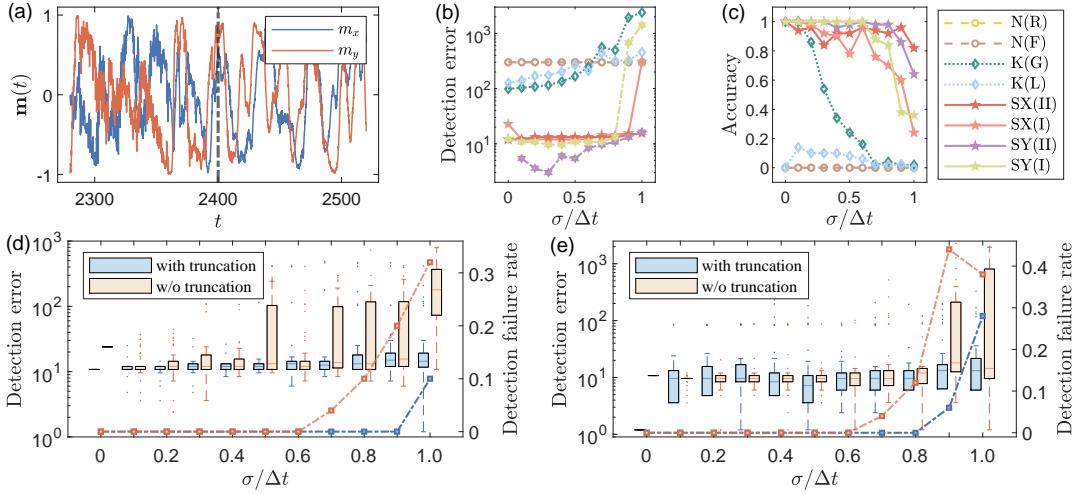


Figure 3. Change point detection results of the fish schooling dataset. (a) Concatenated trajectories of $N_f = 15$ and $N_f = 60$, with 2400s as the splicing point. (b-c) The detection errors and accuracy under different noise levels obtained by different methods. SX and SY represent the Online-SINSDy detection results corresponding to G_{xx} and G_{yy} respectively. (d-e) Box plot of detection errors using G_{xx} indicators (d) and G_{yy} indicators (e) respectively.

3.3. Detect the occurrence of the bearing fault

In manufacturing systems, a bearing fault can lead to a reduction in machine life, a decline in the quality of workpieces, and pose safety risks. Consequently, the diagnosis of bearing faults has emerged as a prominent subject within the engineering community (Safizadeh & Latifi, 2014; Shao et al., 2015; Yuan et al., 2020). Here, we employ Online-SINSDy to identify the occurrence of bearing faults. We generate the dataset by concatenating the trajectory of the normal bearing and that of the fault bearing with a fault diameter of 14 mils at five randomly selected time points to simulate scenarios where the bearings experience faults at different operating conditions. Each trajectory contains around 7×10^5 recorded data with time interval $\Delta t = 1/48,000$ s. The data comes from the Case Western Reserve University bearing dataset.

Here, the library is set as $\Theta(x) = [1, x, x^2, x^3, e^x, e^{2x}, e^{3x}]$. We take the training window length $T_{\text{tr}} = 1,000$ with the stride step $\Delta T = 100$. Phase II begins after observing 20% of the total data with the truncation threshold $\text{th}^{(1)} = 0.1$ and $\text{th}^{(2)} = 0.01$. The sliding window for change point detection is set as $T_{\text{cpd}} = 200$.

During the detection process, we select the point that is detected earlier between the drift and the diffusion indicators as the identified change point. The results are shown in Fig. 4, our method achieves the minimum detection errors compared to the baseline models in various scenarios, except for Case 3, where the detection error of NEWMA with Fourier random features (0.029s) slightly outperforms our approach (0.033s). This indicates that our method can more accurately pinpoint the occurrence of bearing faults,

enabling corresponding countermeasures.

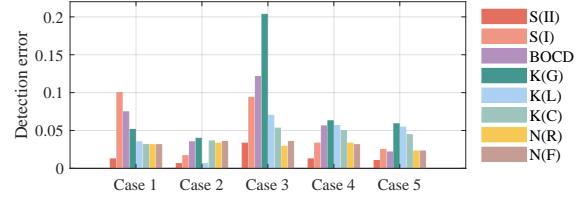


Figure 4. Detection performances in the bearing data.

3.4. Detect the occurrence of the epilepsy

Individuals afflicted with epilepsy experience recurrent seizures that occur at unpredictable times. Detecting the epileptic seizures utilizing the streaming electroencephalogram (EEG) data enables timely intervention, thereby mitigating potentially severe consequences for the patients. Here, we detect the occurrence of epilepsy utilizing the pediatric EEG data (Shoeb, 2009). Most of the recorded EEG signals contain 23 channels with the sampling frequency $\frac{1}{\Delta t} = 256\text{Hz}$.

To assess the efficacy of our approach in this task, we randomly select five patients and employ their EEG data to pinpoint the specific times of epileptic seizures. To make full use of data, we apply our method to detect change points for each of the 23 channels and select the time corresponding to the earliest detected change point as the predicted time for the onset of epileptic seizures. The library is set as $\Theta(x) = [1, x, \dots, x^9]$. During training, we take the time window $T_{\text{tr}} = 1,000$ with the stride step $\Delta T = 100$. After

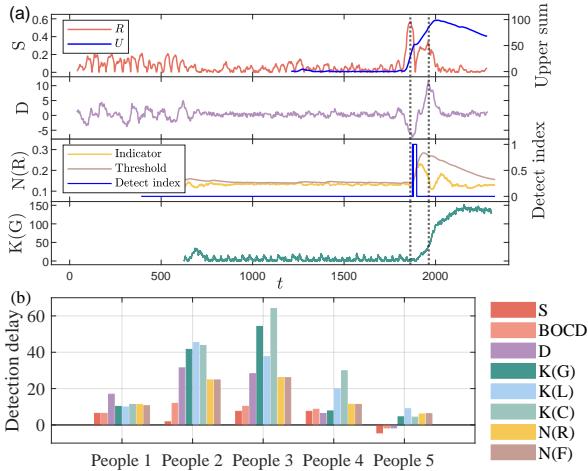


Figure 5. Detection performances in EEG data. (a) Different detection indicators for People 1, including Online-SINSDy, DRE, NEWMA, and KCUSUM. The onset and offset positions of epilepsy are indicated by the black dashed lines. (b) Response time delays using different methods in the data of 5 patients.

observing 20% of the total recorded data, we perform the truncation with $\text{th}^{(1)} = \text{th}^{(2)} = 0.001$. The sliding window for detection is set as $T_{\text{cpd}} = 100$.

Indicators of different methods for People 1 are depicted in Fig. 5(a). The first figure illustrates that our indicators consistently maintain a low level when epileptic seizures are not occurring. However, at the onset and conclusion of epileptic seizures, the indicators experience a significant increase, resulting in a distinctive bimodal pattern. The two peaks precisely align with the onset and conclusion times of epileptic seizures. In comparison, DRE can also identify both the onset and conclusion time points, but with a greater detection delay. BOCD, KCUSUM, and NEWMA can only detect the onset point of epileptic seizures.

As illustrated in Fig. 5(b), our framework shows better detection performance. For People 1-4, the detection delay of our method is lower than that of the baseline models. For People 5, we observe that Online-SINSDy, BOCD, and DRE can provide early warnings before the epileptic seizures, while other methods can only detect changes after the onset of seizures. Furthermore, our approach provides earlier warnings of epilepsy compared with BOCD and DRE.

4. Discussion

In summary, this paper introduces a novel framework, Online-SINSDy, designed for the online change point detection of SDEs. Specifically, Online-SINSDy identifies the drift and diffusion terms of SDEs in an online fashion and utilizes the variations in the identified weight vectors

at each moment as indicators to achieve precise online detection of change points. In comparison to conventional methods relying on the identification of deterministic differential equations, our approach marks a pioneering focus on the inherent dynamical noise within complex systems. We model the system as an SDE, thereby facilitating the concurrent detection of variations in both the deterministic and stochastic parts. The experimental results demonstrate that our framework has an exceptional performance in both identification and change point detection tasks. Notably, Online-SINSDy consistently outperforms benchmark models, exhibiting lower detection errors and higher accuracy across various real-world systems.

It is noteworthy that our method does not necessitate a substantial alteration in any statistic of the system. Hence, it exhibits broader applicability compared to statistically-based methods. Additionally, we acknowledge that certain changes observed in real systems can be induced by noise, as exemplified by the collective dynamics of the cichlid fish. In such scenarios, prediction error-based methods assuming deterministic differential equations, such as O-SINDy and TCD, fail to detect the occurrence of these changes. In contrast, our method can concurrently capture the changes in both the deterministic and stochastic parts of the underlying SDE. This broader capability enhances its applicability compared to prediction error-based methods.

In conclusion, the Online-SINSDy framework demonstrates significant potential in the online identification and change point detection of SDEs. Experimental validations on a diverse set of real-world datasets substantiate the effectiveness, wide applicability, and excellent robustness of our approach.

5. Impact Statements

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

References

- Adams, R. P. and MacKay, D. J. Bayesian online change-point detection. *arXiv preprint arXiv:0710.3742*, 2007.
- Aminikhanghahi, S. and Cook, D. J. A survey of methods for time series change point detection. *Knowledge and Information Systems*, 51(2):339–367, 2017.
- Arlot, S., Celisse, A., and Harchaoui, Z. A kernel multiple change-point algorithm via model selection. *Journal of Machine Learning Research*, 20(162), 2019.
- Bodenham, D. A. and Adams, N. M. Continuous moni-

- 495 toring for changepoints in data streams using adaptive
 496 estimation. *Statistics and Computing*, 27:1257–1270,
 497 2017.
- 498 Boninsegna, L., Nüske, F., and Clementi, C. Sparse learn-
 499 ing of stochastic dynamical equations. *The Journal of
 500 Chemical Physics*, 148(24), 2018.
- 501 Brunton, S. L., Proctor, J. L., and Kutz, J. N. Discovering
 502 governing equations from data by sparse identification
 503 of nonlinear dynamical systems. *Proceedings of the Na-
 504 tional Academy of Sciences*, 113(15):3932–3937, 2016.
- 505 Carpenter, S. R., Arani, B. M., Hanson, P. C., Scheffer, M.,
 506 Stanley, E. H., and Van-Nes, E. Stochastic dynamics of
 507 cyanobacteria in long-term high-frequency observations
 508 of a eutrophic lake. *Limnology and Oceanography Letters*,
 509 5(5):331–336, 2020.
- 510 Champion, K., Lusch, B., Kutz, J. N., and Brunton, S. L.
 511 Data-driven discovery of coordinates and governing equa-
 512 tions. *Proceedings of the National Academy of Sciences*,
 513 116(45):22445–22451, 2019.
- 514 Cribben, I., Wager, T. D., and Lindquist, M. A. Detecting
 515 functional connectivity change points for single-subject
 516 fmri data. *Frontiers in Computational Neuroscience*, 7:
 517 143, 2013.
- 518 Daniels, B. C. and Nemenman, I. Automated adaptive in-
 519 ference of phenomenological dynamical models. *Nature
 520 Communications*, 6(1):8133, 2015.
- 521 Ferrari, A., Richard, C., Bourrier, A., and Bouchikhi, I. On-
 522 line change-point detection with kernels. *Pattern Recog-
 523 nition*, 133:109022, 2023.
- 524 Flynn, T. and Yoo, S. Change detection with the kernel cu-
 525 mulative sum algorithm. In *2019 IEEE 58th Conference
 526 on Decision and Control (CDC)*, pp. 6092–6099, 2019.
- 527 Francesco-Ficetola, G. and Denoël, M. Ecological thresh-
 528 olds: an assessment of methods to identify abrupt changes
 529 in species–habitat relationships. *Ecography*, 32(6):1075–
 530 1084, 2009.
- 531 Hou, J. W., Ma, H. F., He, D., Sun, J., Nie, Q., and Lin,
 532 W. Harvesting random embedding for high-frequency
 533 change-point detection in temporal complex systems. *Na-
 534 tional Science Review*, 9(4):nwab228, 2022.
- 535 Huang, Y., Mabrouk, Y., Gompper, G., and Sabass, B.
 536 Sparse inference and active learning of stochastic dif-
 537 ferential equations from data. *Scientific Reports*, 12(1):
 538 21691, 2022.
- 539 Jhawar, J. and Guttal, V. Noise-induced effects in collec-
 540 tive dynamics and inferring local interactions from data.
 541
- 542 Philosophical Transactions of the Royal Society B, 375
 543 (1807):20190381, 2020.
- 544 Jhawar, J., Morris, R. G., Amith-Kumar, U. R., Danny-Raj,
 545 M., Rogers, T., Rajendran, H., and Guttal, V. Noise-
 546 induced schooling of fish. *Nature Physics*, 16(4):488–493,
 547 2020.
- 548 Kawahara, Y. and Sugiyama, M. Change-point detection
 549 in time-series data by direct density-ratio estimation. In
 550 *Proceedings of the 2009 SIAM International Conference
 551 on Data Mining*, pp. 389–400. Society for Industrial and
 552 Applied Mathematics, 2009.
- 553 Keriven, N., Garreau, D., and Poli, I. Newma: a new method
 554 for scalable model-free online change-point detection.
 555 *IEEE Transactions on Signal Processing*, 68:3515–3528,
 556 2020.
- 557 Li, X., Zhu, Q., Zhao, C., Qian, X., Zhang, X., Duan, X.,
 558 and Lin, W. Tipping point detection using reservoir com-
 559 puting. *Research*, 6:0174, 2023a.
- 560 Li, Y., Wu, K., and Liu, J. Discover governing differen-
 561 tial equations from evolving systems. *Physical Review
 562 Research*, 5(2):023126, 2023b.
- 563 Liu, S., Yamada, M., Collier, N., and Sugiyama, M. Change-
 564 point detection in time-series data by relative density-ratio
 565 estimation. *Neural Networks*, 43:72–8, 2013.
- 566 Lusch, B., Kutz, J. N., and Brunton, S. L. Deep learning
 567 for universal linear embeddings of nonlinear dynamics.
 568 *Nature Communications*, 9(1):4950, 2018.
- 569 Ma, H., Leng, S., Aihara, K., Lin, W., and Chen, L. Ran-
 570 domly distributed embedding making short-term high-
 571 dimensional data predictable. *Proceedings of the National
 572 Academy of Sciences*, 115(43):E9994–E10002, 2018.
- 573 McMahan, B. Follow-the-regularized-leader and mirror de-
 574 scent: Equivalence theorems and l1 regularization. In *Pro-
 575 ceedings of the Fourteenth International Conference on
 576 Artificial Intelligence and Statistics*, pp. 525–533, Cam-
 577 bridge, MA, 2011.
- 578 McMahan, H., Holt, G., Sculley, D., Young, M., Ebner, D.,
 579 Grady, J., Nie, L., Phillips, T., Davydov, E., Golovin, D.,
 580 Chikkerur, S., Liu, D., Wattenberg, M., Hrafinkelsson,
 581 A., Boulos, T., and Kubica, J. Ad click prediction: a
 582 view from the trenches. In Langley, P. (ed.), *Proceedings
 583 of the 19th ACM SIGKDD International Conference on
 584 Knowledge Discovery and Data Mining*, pp. 1222–1230,
 585 Chicago, Illinois, USA, 2013.
- 586 Nabeel, A., A.Karichannavar, Palathingal, S., Jhawar, J.,
 587 and Guttal, V. Pydaddy: A python package for discover-
 588 ing stochastic dynamical equations from timeseries data.
 589 *arXiv preprint arXiv:2205.02645*, 2022.

- 550 Page, E. S. Continuous inspection schemes. *Biometrika*, 41
551 (1/2):100–115, 1954.
- 552 Risken, H. and Risken, H. (eds.). *Fokker-Planck Equation*.
553 Springer Berlin Heidelberg, 1996.
- 554
- 555 Rudy, S. H., Brunton, S. L., Proctor, J. L., and Kutz, J. N.
556 Data-driven discovery of partial differential equations.
557 *Science Advances*, 3(4):e1602614, 2017.
- 558
- 559 Safizadeh, M. and Latifi, S. K. Using multi-sensor data
560 fusion for vibration fault diagnosis of rolling element
561 bearings by accelerometer and load cell. *Information*
562 *Fusion*, 18:1–8, 2014.
- 563
- 564 Shao, H., Jiang, H., Zhang, X., and Niu, M. Rolling bearing
565 fault diagnosis using an optimization deep belief network.
566 *Measurement Science and Technology*, 26(11):115002,
567 2015.
- 568
- 569 Shoeb, A. H. *Application of machine learning to epileptic*
570 *seizure onset detection and treatment*. PhD thesis,
571 Massachusetts Institute of Technology, 2009.
- 572
- 573 Truong, C., Oudre, L., and Vayatis, N. Selective review of
574 offline change point detection methods. *Signal Processing*,
575 167:107299, 2020.
- 576
- 577 Voss, H. U., Kolodner, P., Abel, M., and Kurths, J. Amplitude
578 equations from spatiotemporal binary-fluid convection
579 data. *Physical Review Letters*, 83(17):3422, 1999.
- 580
- 581 Wang, Y., Fang, H., Jin, J., Ma, G., He, X., Dai, X., Yue,
582 Z., Cheng, C., Zhang, H., Pu, D., Wu, D., Yuan, Y.,
583 Goncalves, J., Kurths, J., and Ding, H. Data-driven dis-
584 covery of stochastic differential equations. *Engineering*,
585 17:244–252, 2022.
- 586
- 587 Wei, S. and Xie, Y. Online kernel cusum for change-point
588 detection. *arXiv preprint arXiv:2211.15070*, 2022.
- 589
- 590 Yu, X., Baron, M., and Choudhary, P. K. Change-point de-
591 tection in binomial thinning processes, with applications
592 in epidemiology. *Sequential Analysis*, 32(3):350–367,
593 2013.
- 594
- 595 Yu, Y., Padilla, O. H. M., Wang, D., and Rinaldo, A. A note
596 on online change point detection. *Sequential Analysis*, 42
597 (4):438–471, 2023.
- 598
- 599 Yuan, Y., Ma, G., Cheng, C., Zhou, B., Zhao, H., Zhang,
600 H. T., and Ding, H. A general end-to-end diagnosis
601 framework for manufacturing systems. *National Science*
602 *Review*, 7(2):418–429, 2020.
- 603
- 604

605 A. The binning method for estimating the drift and diffusion values

606 As mentioned in the main text, a conventional technique for estimating KM coefficients involves discretizing the observed
 607 data into spatial intervals, a method often referred to as binning (Boninsegna et al., 2018). In this study, we concentrate
 608 on the univariate scenario ($d = 1$). Initially, we partition the MT_{tr} observed states into Q equally spaced bins. The states
 609 within each bin are then averaged as
 610

$$611 \quad 612 \quad \mathbf{X}(t_l) = \{\mathbf{x}^i(t_{l-j+1})\}_{\substack{i=1,2,\dots,M \\ j=1,2,\dots,T_{\text{tr}}}} \mapsto \{\bar{x}_q, \omega_q\}_{q=1,2,\dots,Q} = \bar{\mathbf{X}}(t_l),$$

613 where \bar{x}_q indicates the averaged value of the data in the q th bin and ω_q indicates the fraction of the data in the q th bin.
 614 Then we can estimate the expectations in Eq. (1) at positions $\{\bar{x}_q\}_{q=1,2,\dots,Q}$ by calculating the average linear and quadratic
 615 variation of the states in each bin, i.e.
 616

$$617 \quad \mathbf{D}^{(1)}(t_l) = \left\{ \frac{\mathbf{x}^i(t_{l-j+1}) - \mathbf{x}^i(t_{l-j})}{\Delta t} \right\}_{\substack{i=1,2,\dots,M \\ j=1,2,\dots,T_{\text{tr}}}} \mapsto \{\bar{f}_q\}_{q=1,2,\dots,Q} = \bar{\mathbf{D}}^{(1)}(t_l),$$

$$620 \quad \mathbf{D}^{(2)}(t_l) = \left\{ \frac{(\mathbf{x}^i(t_{l-j+1}) - \mathbf{x}^i(t_{l-j}))^2}{\Delta t} \right\}_{\substack{i=1,2,\dots,M \\ j=1,2,\dots,T_{\text{tr}}}} \mapsto \{\bar{G}_q\}_{q=1,2,\dots,Q} = \bar{\mathbf{D}}^{(2)}(t_l),$$

622 where \bar{f}_q and \bar{G}_q indicate the averaged value of the linear and quadratic variation of the states in the q th bin respectively.
 623

624 Next, we can utilize the bin-wise averaged data to construct a matrix based on the candidate library $\Theta[\bar{\mathbf{X}}(t_l)] \in \mathbb{R}^{Q \times P}$
 625 consisting of P pre-selected basis functions. For example, the P -order polynomial basis
 626

$$627 \quad 628 \quad 629 \quad 630 \quad \Theta[\bar{\mathbf{X}}(t_l)] = \begin{bmatrix} | & | & & \cdots & | \\ 1 & \bar{\mathbf{X}}(t_l) & \cdots & \bar{\mathbf{X}}^P(t_l) \\ | & | & \cdots & | \end{bmatrix}^\top,$$

631 where the higher polynomials are denoted as $\bar{\mathbf{X}}^2(t_l), \dots, \bar{\mathbf{X}}^P(t_l)$. If the basis functions in the candidate library encompass
 632 all the function terms present in the underlying SDE, we can determine the weight vectors corresponding to the drift term
 633 $\Xi^{(1)}$ and the diffusion term $\Xi^{(2)}$ by solving the following regression problems:
 634

$$635 \quad \mathbf{D}^{(1)}(t_l) = \Theta[\bar{\mathbf{X}}(t_l)]\Xi^{(1)}, \quad (10)$$

$$636 \quad \mathbf{D}^{(2)}(t_l) = \Theta[\bar{\mathbf{X}}(t_l)]\Xi^{(2)}.$$

638 Due to the variations in the data distribution across different space intervals, data in some bins are inevitably quite sparse,
 639 which will lead to substantial errors when estimating the drift and diffusion values in these bins using the binning method.
 640 Thus, some works combine the binning method with the filtering process, which involves eliminating bins with data
 641 quantities falling below a predetermined threshold, to improve the identification of the SDEs.
 642

643 In scenarios where there is a limited amount of training data, we can directly utilize all available data to construct the
 644 regression task as Eq. (3) without the need for bin averaging.
 645

646 It is important to emphasize that although bin-wise averaging is not required, grouping the data into distinct bins remains
 647 essential for executing the filtering procedure. For all the experiments described in this paper, we excluded all bins with data
 648 volume below the average to ensure the robustness of our algorithm. Through experimental validation, we confirm that the
 649 solution derived from Eq. (3) yields identification accuracy comparable to, or even superior to, that achieved by the average
 650 binning method.
 651

To employ the FTRL-proximal algorithm for solving Eq. (10), we define the regression loss of Eq. (10) at the time t_l as
 652

$$653 \quad L_{t_l}^{(1)}(\Xi^{(1)}) = \{\bar{\mathbf{D}}^{(1)}(t_l) - \Theta[\bar{\mathbf{X}}(t_l)]\Xi^{(1)}\}^\top \Omega \{\bar{\mathbf{D}}^{(1)}(t_l) - \Theta[\bar{\mathbf{X}}(t_l)]\Xi^{(1)}\},$$

$$654 \quad L_{t_l}^{(2)}(\Xi^{(2)}) = \{\bar{\mathbf{D}}^{(2)}(t_l) - \Theta[\bar{\mathbf{X}}(t_l)]\Xi^{(2)}\}^\top \Omega \{\bar{\mathbf{D}}^{(2)}(t_l) - \Theta[\bar{\mathbf{X}}(t_l)]\Xi^{(2)}\},$$

655 where the weight matrix Ω is defined as
 656

$$657 \quad \Omega = \text{diag}(\omega_1, \dots, \omega_Q).$$

660 B. The detailed derivation of the closed-form solution of the FTRL-proximal algorithm

661 In the following, our focus remains on the univariate scenario ($d = 1$). Besides, we assume that the stride step $\Delta T = 1$
 662 during the training process. The derivation process presented below is informed by the prior work of McMahan (2011; 2013)
 663 and Li (2023b).

664 Considering Eq. (4) as the regression loss, we can compute its corresponding gradient as
 665

$$\begin{aligned} \nabla L_{t_l}^{(1)} &= \frac{1}{MT_w} \Theta[\mathbf{X}(t_l)]^\top \{\Theta[\mathbf{X}(t_l)]\Xi^{(1)} - \mathbf{D}^{(1)}(t_l)\}, \\ \nabla L_{t_l}^{(2)} &= \frac{1}{MT_w} \Theta[\mathbf{X}(t_l)]^\top \{\Theta[\mathbf{X}(t_l)]\Xi^{(2)} - \mathbf{D}^{(2)}(t_l)\}. \end{aligned} \quad (11)$$

672 The FTRL-proximal algorithm constructs the following unconstrained optimization problems to address the online sparse
 673 regression tasks presented in Eq. (3), relying on the gradient described in Eq. (11):
 674

$$\Xi(t_{l+1}) = \arg \min_{\Xi} \sum_{i=1}^l \nabla L_{t_i} [\Xi(t_i)]^\top \Xi + \lambda_1 \|\Xi\|_1 + \frac{1}{2} \lambda_2 \|\Xi\|_2^2 + \frac{1}{2} \sum_{i=1}^l \sigma_i \|\Xi - \Xi(t_i)\|_2^2, \quad (12)$$

675 where $\Xi(t_l)$ denotes the weight vector of drift or diffusion at the time t_l of the training phase, $\lambda_1 > 0$ and $\lambda_2 > 0$ are the
 676 coefficients of the *origin-centered* L_1 and L_2 regularization respectively, whereas $\|\Xi - \Xi(t_i)\|_2^2$ denotes the *proximal* term
 677 that center additional regularization at the current weight rather than the origin one, which can prevent dramatic change of
 678 the weight vector during the training process, and $\sigma_i > 0$ serves as the learning rate that decays over time.
 679

680 We expand the *proximal* term and delete the constant term, then Eq. (12) can be rewritten as
 681

$$\Xi(t_{l+1}) = \arg \min_{\Xi} \zeta(t_l)^\top \Xi + \lambda_1 \|\Xi\|_1 + \frac{1}{2} (\lambda_2 + \sum_{i=1}^l \sigma_i) \|\Xi\|_2^2, \quad (13)$$

682 where
 683

$$\zeta(t_l) = \sum_{i=1}^l \{\nabla L_{t_i} [\Xi(t_i)] - \sigma_i \Xi(t_i)\}.$$

684 Note that we can decompose Eq. (13) into the following unconstrained optimization problems for each scalar component of
 685 the weight vector:
 686

$$\xi_j(t_{l+1}) = \arg \min_{\xi_j} \zeta_j(t_l) \xi_j + \lambda_1 |\xi_j| + \frac{1}{2} (\lambda_2 + \sum_{i=1}^l \sigma_i) \xi_j^2, \quad (14)$$

687 where $\xi_j(t_l)$ denotes the j th component of $\Xi(t_l)$, $\zeta_j(t_l)$ represents the j th component of $\zeta(t_l)$. Since $|\xi_j|$ is not differentiable
 688 at $\xi_j = 0$, we define its subgradient as
 689

$$\nabla |\xi_j| = \begin{cases} \tau = -1, & \text{if } \xi_j < 0, \\ -1 < \tau < 1, & \text{if } \xi_j = 0, \\ \tau = 1, & \text{if } \xi_j > 0. \end{cases}$$

700 Then we can differentiate the right-hand side of Eq. (14), substitute the subgradient, and set the derivative equal to zero,
 701 leading to the following expression:
 702

$$\zeta_j(t_l) + \lambda_1 \tau + (\lambda_2 + \sum_{i=1}^l \sigma_i) \xi_j = 0. \quad (15)$$

711 Now, we can categorize and discuss the different possible values of $\zeta_j(t_l)$ as follows.
 712

- 713 1. $|\zeta_j(t_l)| \leq \lambda_1$,

715 (a) $\xi_j > 0, \tau = 1, \zeta_j(t_l) + \lambda_1 + (\lambda_2 + \sum_{i=1}^l \sigma_i)\xi_j > 0$. Eq. (15) does not hold.

716 (b) $\xi_j = 0, \tau = -\frac{\zeta_j(t_l)}{\lambda_1} \in (-1, 1)$. Eq. (15) holds.

717 (c) $\xi_j < 0, \tau = -1, \zeta_j(t_l) - \lambda_1 + (\lambda_2 + \sum_{i=1}^l \sigma_i)\xi_j < 0$. Eq. (15) does not hold.

718 2. $\zeta_j(t_l) < -\lambda_1$,

719 (a) $\xi_j > 0, \tau = 1, \zeta_j(t_l) + \lambda_1 + (\lambda_2 + \sum_{i=1}^l \sigma_i)\xi_j = 0$ has a solution $\xi_j = -[\zeta_j(t_l) + \lambda_1](\lambda_2 + \sum_{i=1}^l \sigma_i)^{-1}$.

720 (b) $\xi_j = 0, \tau \in (-1, 1), \zeta_j(t_l) + \lambda_1\tau < 0$. Eq. (15) does not hold.

721 (c) $\xi_j < 0, \tau = -1, \zeta_j(t_l) + \lambda_1\tau + (\lambda_2 + \sum_{i=1}^l \sigma_i)\xi_j < 0$. Eq. (15) does not hold.

722 3. $\zeta_j(t_l) > \lambda_1$,

723 (a) $\xi_j > 0, \tau = 1, \zeta_j(t_l) + \lambda_1 + (\lambda_2 + \sum_{i=1}^l \sigma_i)\xi_j > 0$. Eq. (15) does not hold.

724 (b) $\xi_j = 0, \tau \in (-1, 1), \zeta_j(t_l) + \lambda_1\tau > 0$. Eq. (15) does not hold.

725 (c) $\xi_j < 0, \tau = -1, \zeta_j(t_l) - \lambda_1 + (\lambda_2 + \sum_{i=1}^l \sigma_i)\xi_j = 0$ has a solution $\xi_j = -[\zeta_j(t_l) - \lambda_1](\lambda_2 + \sum_{i=1}^l \sigma_i)^{-1}$.

726 To sum up, we can obtain the following closed-form solution of Eq. (15):

$$\xi_j(t_{l+1}) = \begin{cases} 0, & \text{if } |\zeta_j(t_l)| < \lambda_1, \\ -\frac{\zeta_j(t_l) - \text{sgn}[\zeta_j(t_l)]\lambda_1}{\lambda_2 + \sum_{i=1}^l \sigma_i}, & \text{if } |\zeta_j(t_l)| \geq \lambda_1, \end{cases} \quad (16)$$

727 where $\text{sgn}(\cdot)$ is the sign function.

728 Inspired by the learning rate setting in the prior work (McMahan et al., 2013; Li et al., 2023b), we set the value of the
729 learning rate σ_i as

$$\sigma_i = \frac{1}{\eta_i} - \frac{1}{\eta_{i-1}},$$

730 and

$$\eta_i = \frac{\alpha}{\beta + \sqrt{\sum_{k=1}^i \nabla L_{t_k}[\xi_j(t_k)]^2}}, \quad (17)$$

731 where $\alpha > 0$ and $\beta > 0$ are adjustable hyperparameters, $\nabla L_{t_k}[\xi_j(t_k)]$ denotes the partial derivative of the regression loss at
732 time t_k with respect to $\xi_j(t_k)$. Substituting Eq. (17) into Eq. (16), each scalar component of the weight vector adheres to the
733 following iterative update formula at each training step:

$$\xi_j(t_{l+1}) = \begin{cases} 0, & \text{if } |\zeta_j(t_l)| < \lambda_1, \\ -\frac{\zeta_j(t_l) - \text{sgn}[\zeta_j(t_l)]\lambda_1}{\lambda_2 + \frac{\beta + \sqrt{\sum_{k=1}^i \nabla L_{t_k}[\xi_j(t_k)]^2}}{\alpha}}, & \text{if } |\zeta_j(t_l)| \geq \lambda_1. \end{cases} \quad (18)$$

734 The selection of hyperparameters for the FTRL-proximal algorithm in each experiment is presented in Table 2.

735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
18010
18011
18012
18013
18014
18015
18016
18017
18018
18019
18020
18021
18022
18023
18024
18025
18026
18027
18028
18029
18030
18031
18032
18033
18034
18035
18036
18037
18038
18039
18040
18041
18042
18043
18044
18045
18046
18047
18048
18049
18050
18051
18052
18053
18054
18055
18056
18057
18058
18059
18060
18061
18062
18063
18064
18065
18066
18067
18068
18069
18070
18071
18072
18073
18074
18075
18076
18077
18078
18079
18080
18081
18082
18083
18084
18085
18086
18087
18088
18089
18090
18091
18092
18093
18094
18095
18096
18097
18098
18099
180100
180101
180102
180103
180104
180105
180106
180107
180108
180109
180110
180111
180112
180113
180114
180115
180116
180117
180118
180119
180120
180121
180122
180123
180124
180125
180126
180127
180128
180129
180130
180131
180132
180133
180134
180135
180136
180137
180138
180139
180140
180141
180142
180143
180144
180145
180146
180147
180148
180149
180150
180151
180152
180153
180154
180155
180156
180157
180158
180159
180160
180161
180162
180163
180164
180165
180166
180167
180168
180169
180170
180171
180172
180173
180174
180175
180176
180177
180178
180179
180180
180181
180182
180183
180184
180185
180186
180187
180188
180189
180190
180191
180192
180193
180194
180195
180196
180197
180198
180199
180200
180201
180202
180203
180204
180205
180206
180207
180208
180209
180210
180211
180212
180213
180214
180215
180216
180217
180218
180219
180220
180221
180222
180223
180224
180225
180226
180227
180228
180229
180230
180231
180232
180233
180234
180235
180236
180237
180238
180239
180240
180241
180242
180243
180244
180245
180246
180247
180248
180249
180250
180251
180252
180253
180254
180255
180256
180257
180258
180259
180260
180261
180262
180263
180264
180265
180266
180267
180268
180269
180270
180271
180272
180273
180274
180275
180276
180277
180278
180279
180280
180281
180282
180283
180284
180285
180286
180287
180288
180289
180290
180291
180292
180293
180294
180295
180296
180297
180298
180299
180300
180301
180302
180303
180304
180305
180306
180307
180308
180309
180310
180311
180312
180313
180314
180315
180316
180317
180318
180319
180320
180321
180322
180323
180324
180325
180326
180327
180328
180329
180330
180331
180332
180333
180334
180335
180336
180337
180338
180339
180340
180341
180342
180343
180344
180345
180346
180347
180348
180349
180350
180351
180352
180353
180354
180355
180356
180357
180358
180359
180360
180361
180362
180363
180364
180365
180366
180367
180368
180369
180370
180371
180372
180373
180374
180375
180376
180377
180378
180379
180380
180381
180382
180383
180384
180385
180386
180387
180388
180389
180390
180391
180392
180393
180394
180395
180396
180397
180398
180399
180400
180401
180402
180403
180404
180405
180406
18

770 C. Pseudocode for Online-SINSDy

771 In this section, we present the pseudocode of Online-SINSDy for the identification and change point detection tasks. To
 772 facilitate a clearer exposition, we divide the pseudocode into two parts. Algorithm 1 provides the pseudocode for utilizing
 773 Online-SINSDy to obtain the dynamics of the drift and the diffusion weights and perform system identification, while
 774 Algorithm 2 presents the pseudocode for change point detection utilizing the dynamics of the weights obtained in Algorithm
 775 1.

776 It is noteworthy that although the pseudocode is presented in two parts, system identification and change point detection
 777 are carried out simultaneously. For the sake of brevity, we only consider the case where the state dimension $d = 1$ in this
 778 context. However, the algorithms presented in both pseudocodes can be readily extended to the multivariate case. Moreover,
 779 we assume that there is only one change point in the time series data in Algorithm 2. If there are multiple change points in
 780 the streaming data, restarting Algorithm 2 after detecting one change point allows for the continued detection of subsequent
 781 change points.

785 Algorithm 1 Online-SINSDy algorithm for system identification

786 **Input:** $\mathbf{X}(t)$: the streaming data; $\alpha, \beta, \lambda_1, \lambda_2$: the hyperparameters for FTRL-proximal; T_{II} : the beginning time for
 787 Phase II; $\text{th}^{(1)}, \text{th}^{(2)}$: the truncation threshold for drift and diffusion in Phase II; T_{tr} : the sliding window length for
 788 training; ΔT : the stride step for training.
 789 **Output:** dynamics of the identified drift weight $\Xi^{(1)}(t)$; dynamics of the identified diffusion weight $\Xi^{(2)}(t)$.
 790 Assume the dimension of the data $d = 1$. Denote N as the total length of the streaming data, M as the number of
 791 trajectories, P as the number of basis functions in the library, $\xi_j^{(1)}(t)$ and $\xi_j^{(2)}(t)$ as the j th component of the identified
 792 drift and diffusion weight vectors at time t respectively.
 793 Initialize $\zeta(t_{1-\Delta T}) = 0, \xi_j^{(1)}(t_{1-\Delta T}) = \xi_j^{(2)}(t_{1-\Delta T}) = 0, \kappa_j^{(1)} = \kappa_j^{(2)} = 0$ ($j = 1, 2, \dots, P$).
 794 **for** time index l in $T_{\text{tr}} : \Delta T : N$ **do**
 795 Receive the collected data $\mathbf{X}(t_l) = [x^1(t_l), x^2(t_l), \dots, x^M(t_l), x^1(t_{l-1}), \dots, x^M(t_{l-T_{\text{tr}}+1})]$.
 796 Calculate the linear and quadratic variation $\mathbf{D}^{(1)}(t_l), \mathbf{D}^{(2)}(t_l)$.
 797 Calculate the library matrix $\Theta[\mathbf{X}(t_l)]$.
 798 **if** $t_l < T_{II}$ **then**
 799 $I^{(1)} = I^{(2)} = \{1, 2, \dots, P\}$.
 800 **else**
 801 $I^{(1)} = \{i \mid \xi_i^{(1)}(t_{l-\Delta T}) > \text{th}^{(1)}\}, I^{(2)} = \{i \mid \xi_i^{(2)}(t_{l-\Delta T}) > \text{th}^{(2)}\}$.
 802 $\xi_j^{(1)}(t_l) \leftarrow 0, \forall j \notin I^{(1)}; \xi_j^{(2)}(t_l) \leftarrow 0, \forall j \notin I^{(2)}$.
 803 **end if**
 804 **for** each component k in $I^{(1)}$ **do**
 805 Calculate $\xi_k^{(1)}(t_l)$ using Eq. (18).
 806 Calculate the gradient of the regression loss $\nabla L_{t_l}^{(1)}[\xi_k^{(1)}(t_l)]$ using Eq. (11).
 807 Calculate $\sigma_k^{(1)}(t_l) = \left\{ \sqrt{\kappa_k^{(1)} + \nabla L_{t_l}^{(1)}[\xi_k^{(1)}(t_l)]^2} - \sqrt{\kappa_k^{(1)}} \right\} / \alpha$.
 808 Calculate $\zeta_k^{(1)}(t_l) = \zeta_k^{(1)}(t_{l-\Delta T}) + \nabla L_{t_l}^{(1)}[\xi_k^{(1)}(t_l)] - \sigma_k^{(1)}(t_l)\xi_k^{(1)}(t_l)$.
 809 Update $\kappa_k^{(1)} \leftarrow \kappa_k^{(1)} + \nabla L_{t_l}^{(1)}[\xi_k^{(1)}(t_l)]^2$.
 810 **end for**
 811 **for** each component k in $I^{(2)}$ **do**
 812 Calculate $\xi_k^{(2)}(t_l)$ using Eq. (18).
 813 Calculate the gradient of the regression loss $\nabla L_{t_l}^{(2)}[\xi_k^{(2)}(t_l)]$ using Eq. (11).
 814 Calculate $\sigma_k^{(2)}(t_l) = \left\{ \sqrt{\kappa_k^{(2)} + \nabla L_{t_l}^{(2)}[\xi_k^{(2)}(t_l)]^2} - \sqrt{\kappa_k^{(2)}} \right\} / \alpha$.
 815 Calculate $\zeta_k^{(2)}(t_l) = \zeta_k^{(2)}(t_{l-\Delta T}) + \nabla L_{t_l}^{(2)}[\xi_k^{(2)}(t_l)] - \sigma_k^{(2)}(t_l)\xi_k^{(2)}(t_l)$.
 816 Update $\kappa_k^{(2)} \leftarrow \kappa_k^{(2)} + \nabla L_{t_l}^{(2)}[\xi_k^{(2)}(t_l)]^2$.
 817 **end for**
 818 **end for**

Algorithm 2 Online-SINSDy algorithm for change point detection

825
 826 **Input:** $\Xi(t)$: the streaming weight vectors of drift or diffusion obtained by Algorithm 1; T_{cpd} : the sliding window length
 827 for change point detection; h, γ : the hyperparameters for CUSUM; k : number of indicators in Phase II used to estimate
 828 the initial average m_R and standard deviation σ_R .
 829 **Output:** detected change point T_d .
 830 Initialize $U(t_{k-1}) = 0, t_{\text{CU}} = +\infty, R_{\text{CU}} = +\infty, T_d = 0$.
 831 Assume the dimension of the data $d = 1$ and only one change point exists in the time series data. For conciseness,
 832 denote the time series of the drift or diffusion weight vectors obtained by Algorithm 1 from the beginning of Phase II as
 833 $\Xi(t_1), \Xi(t_2), \dots, \Xi(t_{N_d})$ in sequence.
 834 **for** time index l **in** $2T_{\text{cpd}} + 1 : 1 : N_d$ **do**
 835 Calculate $\bar{\Xi}(t_{l-T_{\text{cpd}}}, t_l), \bar{\Xi}(t_{l-2T_{\text{cpd}}}, t_{l-T_{\text{cpd}}})$.
 836 Calculate $R(t_{l-T_{\text{cpd}}})$.
 837 **if** $l = k + T_{\text{cpd}} - 1$ **then**
 838 Calculate $m_R = \frac{1}{k} \sum_{i=T_{\text{cpd}}}^l R(t_i)$.
 839 Calculate $\sigma_R = \sqrt{\frac{1}{k-1} \sum_{i=T_{\text{cpd}}}^l [R(t_i) - m_R]^2}$.
 840 **end if**
 841 **if** $l > k + T_{\text{cpd}} - 1$ **then**
 842 Calculate $U(t_{l-T_{\text{cpd}}}) = \max[0, U(t_{l-T_{\text{cpd}}-1}) + R(t_{l-T_{\text{cpd}}}) - m_R - \frac{1}{2}h\sigma_R]$.
 843 **if** $U(t_{l-T_{\text{cpd}}}) > \gamma\sigma_R$ and $U(t_i) \leq \gamma\sigma_R (\forall i < l - T_{\text{cpd}})$ **then**
 844 Update $t_{\text{CU}} \leftarrow t_{l-T_{\text{cpd}}}$.
 845 Update $R_{\text{CU}} \leftarrow R(t_{l-T_{\text{cpd}}})$.
 846 **end if**
 847 **end if**
 848 **if** $t_{l-T_{\text{cpd}}} \geq t_{\text{CU}}$ and $R_{l-T_{\text{cpd}}} \geq R_{\text{CU}}$ **then**
 849 Update $T_d \leftarrow \operatorname{argmax}_{t \in \{T_d, t_{l-T_{\text{cpd}}}\}} R(t)$.
 850 **end if**
 851 **end for**

852 c

D. Results of online identification of SDEs

853
 854
 855 In this section, we validate the effectiveness of our Online-SINSDy algorithm in the identification task, which relies on
 856 the KM formulae and the modified FTRL-proximal algorithm, using both synthetic data and real-world systems. We also
 857 experimentally demonstrate the necessity and effectiveness of our two-phase training method.

D.1. The double-well system

858 We first illustrate the online data-driven discovery of SDEs by the example of a 1-d double-well system. The drift term and
 859 the diffusion term are given as

$$f(x) = -3x^3 + 2x, \quad G(x) = \frac{1}{4}. \quad (19)$$

860 We generate $M = 100$ independent trajectories by integrating Eq. (1) with the Euler-Maruyama method, each trajectory
 861 contains $N = 10^5$ recorded states with time intervals $\Delta t = 0.01$. The basis functions of the library consist of polynomials
 862 of the state up to the order of 9, i.e. $\Theta(x) = [1, x, \dots, x^9]$. We take the training time window as $T_{\text{tr}} = 10$ with the stride
 863 step $\Delta T = 10$. In the following experiments, we repeatedly generate 50 sets of data and construct the expressions of the
 864 drift function and the diffusion function of each dataset to verify the robustness of our algorithm. We use the mean value
 865 of the weight vectors as the identified coefficients of the function expressions and use the standard deviation to reflect the
 866 stability of the identification.

867 As is shown in Fig. 6(b), Online-SINSDy can accurately recover the analytical expressions for the drift and diffusion

functions with quite small standard deviation. Fig. 6(c) demonstrates that Online-SINSDy with truncation in Phase II can accurately discover the drift function of the double-well system as long as the threshold $\text{th}^{(1)} \geq 0.4$, whereas training without truncation fails to screen the terms appearing in the drift function correctly.

We show the dynamics of all the drift coefficients during the training process in Fig. 6(d) to further explore the reason for the failure of the training without truncation in Phase II. As is shown in the left figure, if the truncation operation of Phase II is not performed during the training process, the weight vector will converge and stabilize at a false surrogate stable state, resulting in the failure of the identification. However, if we conduct the truncation process in Phase II after the stability is achieved in Phase I, the coefficients will escape from this false equilibrium and converge toward their true values. As depicted in the right figure, after 5,000 training steps, we apply the truncation operation, where the non-zero coefficients of the interfering terms (x^8, x^9) below the threshold $\text{th}^{(1)} = 0.5$ are immediately set to zero. The truncation operation disturbs the previous stable state, allowing the actual terms (x, x^3) in the drift function to further approach their true values, while other interfering terms (x^5) whose coefficients are greater than the threshold gradually decrease until they ultimately become smaller than the threshold and are set to zero. At this point, we have filtered out all genuine terms from the library, and subsequent training will rapidly converge the coefficients of these terms towards their true values.

Actually, if we have more accurate prior knowledge about the range of the coefficients, we can use it to accelerate the convergence of the algorithm. As is shown in Fig. 6(e), we set a relative error threshold $\kappa_{\text{error}} = 0.1$ and record the number of training steps in Phase II required for the coefficients to converge within κ_{error} under different truncation thresholds. As the truncation threshold increases, the convergence rate increases as well.

To evaluate the robustness of our Online-SINSDy algorithm against noisy perturbations, we also introduce the Gaussian white noise with different standard deviation σ into the observational data. Fig. 6(f) illustrates that when $\sigma < \Delta t$, the relative error of the drift coefficients remains at a low level; when $\Delta t \leq \sigma < 2\Delta t$, the relative error experiences a steady increase, but Online-SINSDy still accurately identifies the terms appearing in the drift function; when $\sigma \geq 2\Delta t$, the inference accuracy begins to decrease. The inference accuracy remains 100% for the diffusion function with the noise intensity $0 \leq \sigma \leq 3\Delta t$, in spite of a gradual growth of the relative error. Thus, the algorithm we proposed is robust against the noise to a certain extent.

Note that all the experiments above are conducted both with and without bin-wise averaging of the data, and we find little difference between them.

D.2. The fish schooling dataset

Online-SINSDy can be applied to discover the governing SDEs of many real-world systems with the collected streaming data, which is quite challenging because of the strong observation noise and the existence of missing data points in real scenarios. Here we focus on the dynamics of collective alignment in groups of the cichlid fish *Etroplus suratensis* with different group sizes. We have mentioned in the main text that the highly polarized and coherent motion of the cichlid fish is a finite-size noise-induced effect, and the mesoscopic dynamics can be fitted as Eq. (8), where \mathbf{m} defined by Eq. (9) is a 2-d order parameter vector called group polarization, which represents the collective fish motion direction and the degree of alignment.

For Eq. (8), the deterministic dynamics of the fish ($N_f \rightarrow \infty$) promote the system to produce an isotropic motion pattern ($|\mathbf{m}| = 0$), while the existence of the intrinsic noise, which is caused by the probabilistic interactions between individuals, lead to the noise-introduced schooling behavior ($|\mathbf{m}| \lesssim 1$). Moreover, the intensity of the intrinsic noise that introduces the aligned motion is inversely proportional to the group size, meaning that the smaller the group size, the more likelihood of such highly polarized and coherent motion, as is shown in Fig. 7(a-c).

We here focus on three scenarios from small to medium group sizes with $N_f = 15, N_f = 30$, and $N_f = 60$ respectively. The dataset consists of four trajectories for $N_f = 15$ and $N_f = 30$ and three trajectories for $N_f = 60$. Each trajectory contains around 2.5×10^4 recorded data of the group polarization \mathbf{m} with uniform time interval $\Delta t = 0.12\text{s}$ and includes many missing data points. To estimate the KM coefficients, we take $40\Delta t$ as the time interval to calculate the drift value and Δt to calculate the diffusion value, which is the same as prior work (Jhawar et al., 2020). The library is set as $\Theta(x) = [1, x, \dots, x^9]$. We take the training window length $T_{\text{tr}} = 12\text{s}$ and the stride step length $\Delta T = 12\text{s}$. The first half of the trajectories is used for Phase I training, and the second half is used for Phase II training with truncation. The truncation threshold for both drift and diffusion are set as $\text{th}^{(1)} = \text{th}^{(2)} = 0.01$.

The experimental results are shown in Fig. 7(d-g), we calculate a series of accurate estimates of the drift values and the

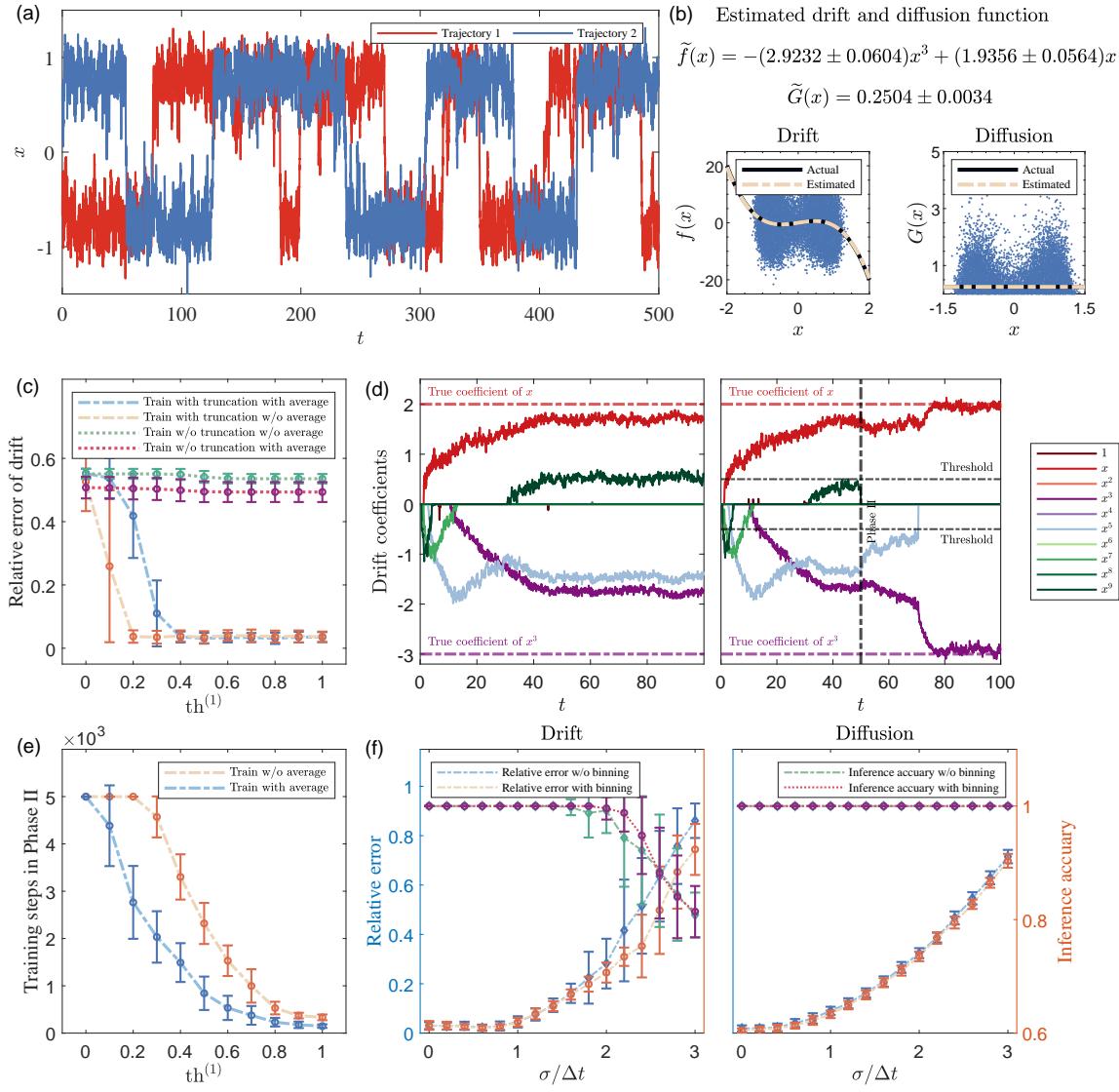


Figure 6. Identification results of the double-well system described by Eq. (19). (a) Two trajectories of the double-well system. (b) The estimated drift and diffusion functions using Online-SINSDy. The blue dots represent the linear and quadratic variation at each point calculated during the training process. (c) The relative error between the identified drift function and the true drift function. (d) The dynamics of the drift coefficients during the training process with (right) or without (left) the truncation operation in Phase II. (e) The number of training steps in Phase II required for the coefficients to converge within a small neighborhood around the true coefficients under different truncation thresholds. (f) The relative error and inference accuracy of the drift (left) and diffusion (right) terms under different noise levels.

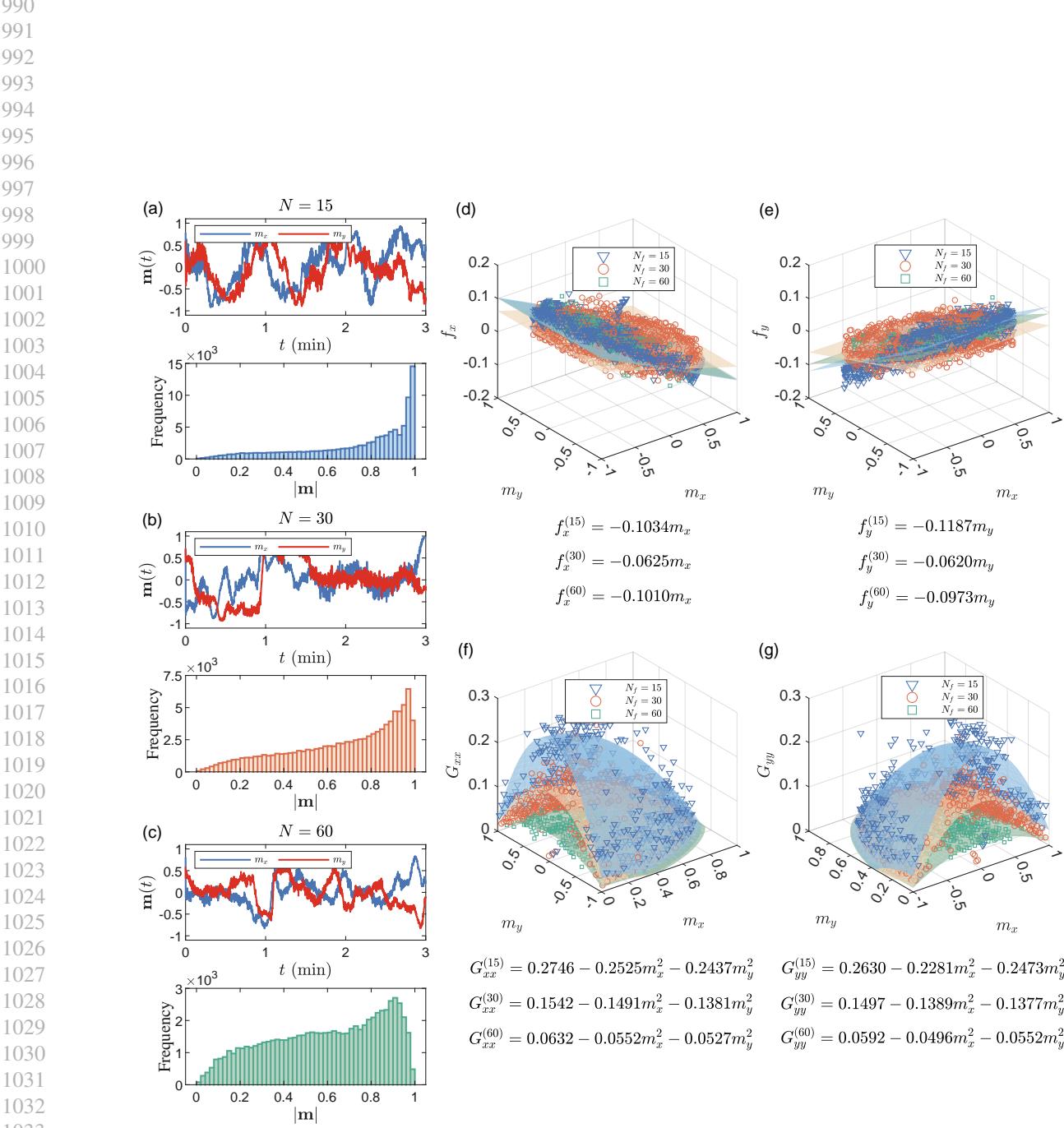


Figure 7. Identification results of the fish schooling dataset. (a-c) Collective dynamical trajectories and statistical histograms of group polarization for the cichlid fish with different group sizes. (d-e) The identified drift functions f_x and f_y for different group sizes. (f-g) The identified diffusion functions G_{xx} and G_{yy} for different group sizes.

1045 diffusion values by bin-wise averaging the whole dataset and find that the surface corresponding to the identified drift and
 1046 diffusion functions can fit these values well. In fact, the Online-SINSDy algorithm successfully identifies the correct sparse
 1047 functional terms from the redundant basis functions and ultimately constructs accurate and concise function expressions
 1048 of the underlying SDEs for different values of N_f . From the discovered SDEs, we can not only identify the deterministic
 1049 behavior of the fish but also gain insight into the inverse relationship between the noise intensity and the group size.
 1050 Moreover, we can also recover $\hat{\gamma}_2 = 3.8$ from the coefficients of the identified SDEs, which is consistent with the true value
 1051 $\gamma_2 \approx 4.0$.

1052
 1053 **D.3. The rolling bearing vibration dataset**
 1054 We also use Online-SINSDy to discover the governing equations of the rolling bearing vibration dynamics from the Case
 1055 Western Reserve University (CWRU) bearing dataset.

1056 We here focus on the trajectories of a normal fan-end bearing and a faulted fan-end bearing with a ball fault diameter
 1057 of 14 mils, which contain around 4.8×10^5 and 3.8×10^5 recorded data points respectively with the same time interval
 1058 $\Delta t = 1/48,000$. In each case, only one recorded trajectory is available for training. The library is set as $\Theta(x) =$
 1059 $[1, x, x^2, x^3, e^x, e^{2x}, e^{3x}]$. We take the training window length $T_{\text{tr}} = 1/48\text{s}$ with stride step length $\Delta T = 1/480\text{s}$. After
 1060 Phase I training with the first 60% data, Phase II training begins with truncation threshold $\text{th}^{(1)} = 0.3$ and $\text{th}^{(2)} = 0.01$.

1061 Although only one trajectory is available for training in each case, which is challenging for the online identification task,
 1062 our algorithm still provides reasonable and concise functional expressions for the bearing vibration dynamics, as is shown
 1063 in Table 3. Due to the availability of ample data in the trajectories, we employ the bin-wise averaging method to obtain
 1064 precise estimates of the drift and diffusion terms. These estimates are subsequently used as reference points to evaluate
 1065 the accuracy of the identification. The mean square error between the estimates derived from bin-wise averaging and the
 1066 corresponding identified drift and diffusion terms after training is then calculated. As depicted in Fig. 8, the SDEs identified
 1067 by Online-SINSDy effectively captured the authentic dynamics of the bearing from the recorded data points.

Table 3. Identification results of the rolling bearing vibration dataset.

Data set	Drift function	Drift MSE	Diffusion function	Diffusion MSE
Normal	$0.9134 - 5.2558x - 4.6694x^2 - 0.6426e^{3x}$	0.0019	$3.7311x^2 + 0.6454e^x$	0.0003
Faulted	$1.9914 + 11.4064x^2 + 0.6755e^x - 0.7715e^{2x} - 1.7427e^{3x}$	0.0013	$-4.4265x + 14.3748x^2 + 0.9046e^{3x}$	0.0109

E. The equivalence of the FTRL-proximal algorithm and the online gradient descent method

The parameter update formula in online gradient descent (OGD) algorithms is given by:

$$\Xi(t_{l+1}) = \Xi(t_l) - \eta_l \nabla L_{t_l}[\Xi(t_l)],$$

where η_l is the learning rate at t_l and L represents the loss function.

If we ignore the L_1 and L_2 regularization terms that are used to induce sparsity of the weight vectors, the Eq.(3) takes the following form:

$$\Xi(t_{l+1}) = \arg \min_{\Xi} \sum_{i=1}^l \nabla L_{t_i}[\Xi(t_i)]^\top \Xi + \frac{1}{2} \sum_{i=1}^l \sigma_i \|\Xi - \Xi(t_i)\|_2^2,$$

Taking the derivative of the right-hand side and setting it to zero, we obtain:

$$\Xi(t_{l+1}) = \frac{\sum_{i=1}^l \sigma_i \Xi(t_i)}{\sum_{i=1}^l \sigma_i} - \eta_l \sum_{i=1}^l \nabla L_{t_i}[\Xi(t_i)],$$

where $\eta_l = 1 / \sum_{i=1}^l \sigma_i$ can be regarded as the learning rate. It is evident that the parameter update form induced by our objective function closely resembles that of OGD. In the parameter update formula of OGD, we replace the weight vectors at t_l with the weighted average of all previous weight vectors and utilize the cumulative gradient in place of the instantaneous

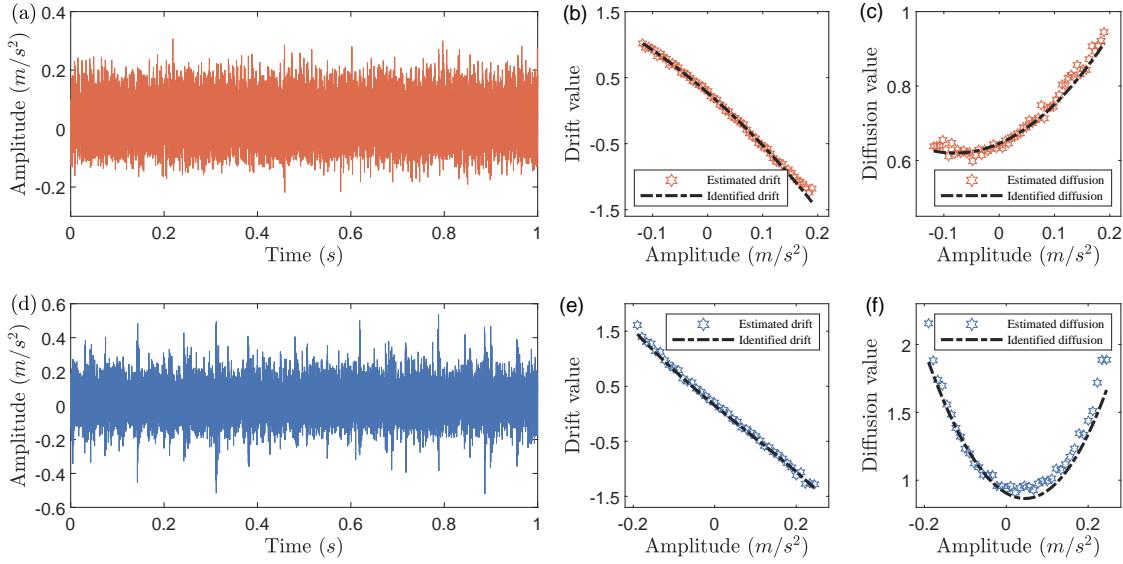


Figure 8. Identification results of the rolling bearing vibration dataset. (a) The trajectory of a normal fan-end bearing. (b-c) The estimated drift and diffusion values and the identified drift and diffusion functions of the normal fan-end bearing. The estimated values are calculated by the bin-wise averaging method using the whole dataset. (d) The trajectory of a faulted fan-end bearing with a ball fault diameter of 14 mils. (e-f) The estimated drift and diffusion values and the identified drift and diffusion functions of the faulted fan-end bearing.

gradient. Indeed, employing the weighted average weight vectors and cumulative gradient can enhance the stability of the OGD algorithm, reducing the impact of rare events in the dataset on training.

In summary, the existence of the first term and the proximal term in Eq.(3) enables the solution to be interpreted as an OGD of the weight vectors.

F. The detection delay of Online-SINSDy on synthetic and real-world datasets.

We note that many articles directly employ detection delay (the positive part of the stopping time minus the true change-point time) to measure the effectiveness of the detection. Here, we include statistical figures corresponding to the detection delay under different noise levels for both the synthetic dataset (Fig. 9) and the fish schooling dataset (Fig. 10) in Section 3.

(d)

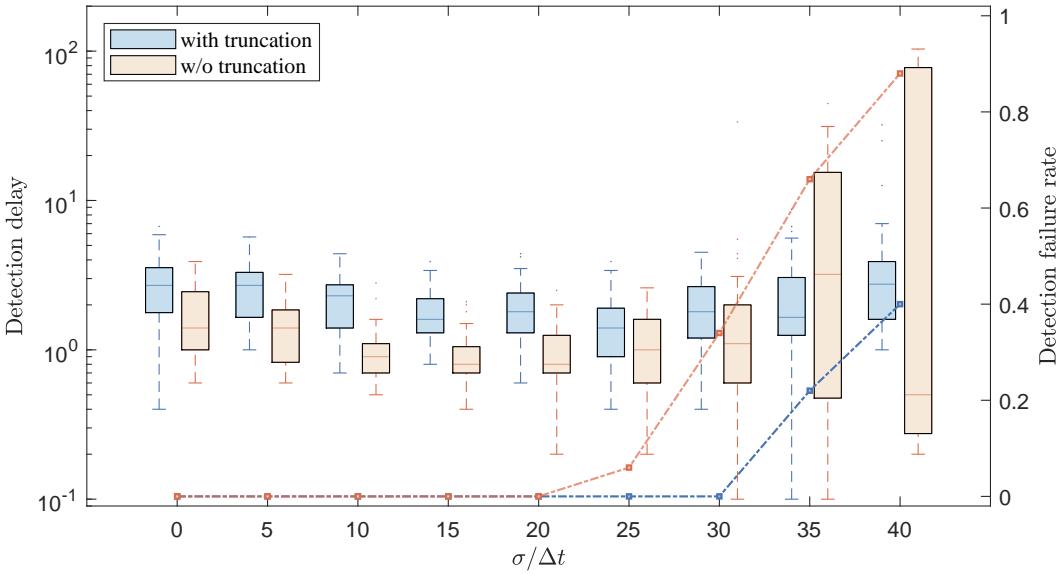


Figure 9. Box plot of the detection delay of Online-SINSDy on the synthetic dataset under different levels of noise.

(d)

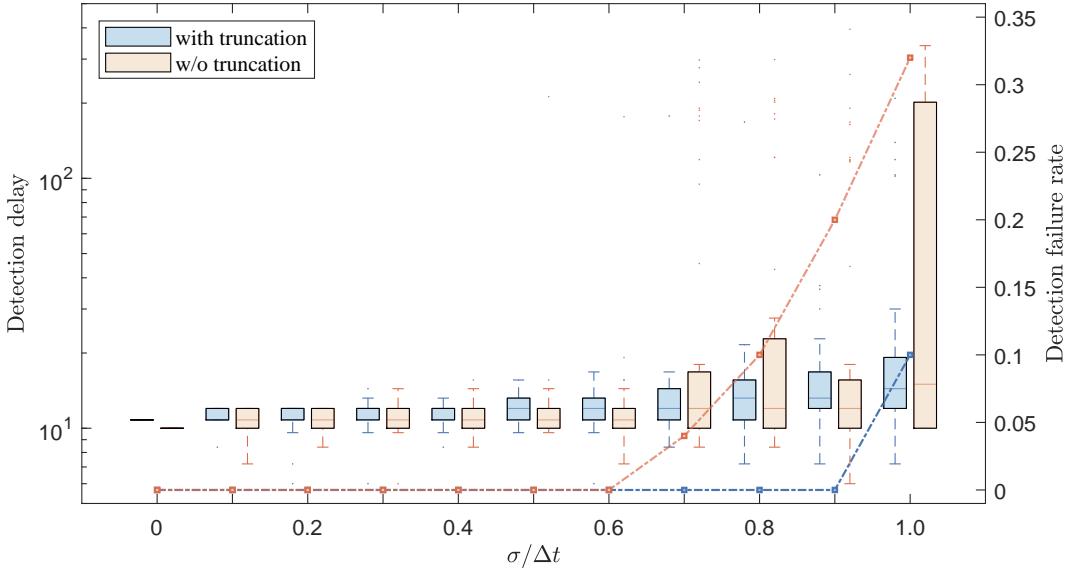


Figure 10. Box plot of the detection delay of Online-SINSDy on the fish schooling dataset using G_{xx} indicator under different levels of noise.