

Bayesian equation discovery of non-linear stochastic dynamical systems through sparse linear regression using spike and slab prior

Tapas Tripura¹ and Souvik Chakraborty^{1,2,*}

¹Department of Applied Mechanics, Indian Institute of Technology Delhi, 110016, India

²School of Artificial Intelligence (ScAI), Indian Institute of Technology Delhi, 110016, India

*souvik@am.iitd.ac.in

ABSTRACT

We propose a novel data-driven framework for discovering the governing physics of a non-linear dynamical system from the noisy measurement of system states in this work. A complete understanding about the governing physics of a dynamical system significantly enhances our ability to model and predict the systems behaviour in an unseen environment. The governing physics of many natural processes are generally modelled by putting together the first principles with scientific laws. These mathematical models are many a case approximate in nature. Often in complex environments the minimal knowledge of the physical process are not obtainable. In such cases, the modern advancement in the machine learning techniques has enabled the researchers and scientists to identify the governing physics of a complicated system from noisy data without invoking scientific assumptions. However recently developed physics discovery techniques requires the information about both the input force and output state measurement data. The driving input force in almost any natural process is not measurable directly due to limitation in the current state of art of the force estimation algorithms. To add, the forces in a naturally occurring process continuously evolves over time, making them intractable to measure. In such cases these techniques are bound to fail. Also measurement of all the systems states such as acceleration, displacement and velocities of a dynamical system can impose additional cost. Thus a novel framework should make use of as minimum states as possible and discover the governing physics without explicit measurement of the input force data. The absence of external force information essentially transfers the physics discovery problem from deterministic to stochastic domain. Hence we leverage the stochastic calculus theory and propose a novel framework to discover the governing physics of a non-linear stochastic dynamical system, in terms of the stochastic differential equations (SDE). In the proposed framework, a sparse Bayesian linear regression is formulated for identification of an SDE in terms of its drift and diffusion components. The target vector in the sparse linear regression is formulated as linear and quadratic variations of the noisy displacement time history. The spike and slab distributions are used as prior over the weight vector for promotion of sparsity in the solution. The statistical properties of the posterior distribution are utilized to obtain the final governing physics and associated confidence bounds. The fidelity of the proposed work is tested on variety of examples involving both fully and partially observable state measurement data. The results indicate the successful identification of actual governing physics.

Introduction

Physical systems are governed by the physical law, often represented either in form of ordinary differential equations (ODE) or partial differential equations (PDE). Once the governing ODE/PDE is defined, there exist a plethora of methods including finite element and finite volume methods for solving the same. The literature on solution of ODE/PDE is quite matured and given sufficient computational resources, it is possible to solve the governing ODE/PDE with sufficient accuracy. Modern techniques such a physics-informed deep learning and neural network based operator learning can also be used. However, due to various assumptions and approximations, the governing equations often fails to represent the exact physics of the system and results in modelling and prediction errors. With advancement in the sensor technology and the internet of things (IoT), we have access to data in abundance whereas the underlying governing physics often remains elusive specifically in domains such as climate science¹, biology², physics³, chemistry⁴ and finance⁵ to name a few. One possible alternative is obviously to train a purely data-driven machine learning based model to obtain an input-output mapping; however, such models often do not generalize to unseen environment and out-of-distribution inputs. To address this issue, methods for data-driven equation discovery has recently been proposed. A seminal work in this area was carried out by⁶ wherein symbolic regression was utilized for discovering underlying structure of the data. Another seminal work in this area includes work by⁷. However, most of the work in this area is limited by the fact that information on both input and state vector are needed. However, in practice, the input is often not measurable and only few states are accurately observable; this greatly limits the applicability of the available physics discovery techniques. To address this issue, we propose a novel approach that leverages sparse learning, Bayesian

statistics, and stochastic calculus and enables discovery of governing physics from only the state measurements.

The development in the equation discovery techniques has made significant progress from early 1970s when the equation discovery was mostly dependent on the expertise in the interested area, to the modern machine learning techniques that can extract the physical law of the underlying process with only little human supervision. In the early 1970s the models were selected by maintaining a balance between the model complexity and its fitness. The fitness of the models were measured using Akaike (AIC) and Bayesian information criterion (BIC)^{8,9}. Later advances in the machine learning tools in combination with new sophisticated data measurement instruments gave rise to the data-driven techniques for discovery of governing physics. Towards this, first attempt was made by combining symbolic regression with genetic programming to select the best combination of the functions from candidate functional forms that accurately represent the data^{6,10}. There are also equation-free methods that bypasses the need of formulating constitutive equations to track the time evolution of the dynamical systems in closed form. These methods provide a practical recipe for multiscale simulations through local microscopic simulations in time and space¹¹.

Following the work of symbolic regression^{6,10}, Sparse Identification of Nonlinear Dynamics (SINDy) was proposed later for discovery of the governing physics of a non-linear dynamical system⁷. In terms of accuracy and interpretability of the discovered physics, the SINDy framework overcame the shortcomings of the previous discovery techniques. The idea behind SINDy was to use sparse linear regression in the purview of least-square fitting for selecting the most dominating candidate functions that best represent the data. Due to sparsity promoting approach, SINDy proved to be computationally efficient and scalable with the increase in the input dimension. In the later years, SINDy has shown tremendous applications and area-specific developments of sparse linear regression in various fields. Few example are, sparse identification of biological networks in biology¹² and sparse identification of chemical reaction in chemistry^{13,14}, sparse modelling for state estimation in fluid mechanics^{15,16}, system identification of structures with hysteresis and sparse learning of aerodynamics of bridges in structural systems^{17,18}, sparse model selection using an integral formulation¹⁹, sparse model selection of dynamical system using information criteria²⁰, sparse identification for predictive control²¹, identification of structured dynamical systems with limited data²², model identification using recovery of differential equations from short impulse response time-series data²³, identifying stochastic dynamic equation²⁴, and discovery of partial differential equations^{25,26}, among others.

Apart from SINDy, identification of non-linear dynamical systems using deep neural networks were also proposed by the researchers²⁷⁻²⁹. The deep learning approaches for non-linear dynamical systems exists as a black-box which takes the data as a input and provides the prediction as output. Later, a different approach within the Bayesian framework for discovery of governing physics from noisy and/or limited data was proposed^{30,31}. In this approach the least-square based sparse linear regression was replaced by more accurate sparse Bayesian linear regression technique. Due to the ability of the Bayesian inference to perform simultaneous model selection and parameter estimation, this approach eliminated the possibility of overfitting of the governing model. The sparsity in the solution was promoted by assigning the appropriate sparse promoting priors over the weight vector of the sparse regression problem. This approach provided a natural elimination of the basis functions that do not present the data accurately. As a output, this framework provides the posterior distributions of the weight vector; this means one could accurately able to identify the uncertainty in the parameter estimation and construct a confidence interval for future predictions.

Despite the progress in equation discovery from data in recent times, almost all the approaches require measurements for both state and input variables; the only exceptions is perhaps the Stochastic SINDy²⁴. Unfortunately, we often only have access to the state variables only; this significantly limits the applicability of the available approaches. We hereby propose a equation discovery framework rooted in stochastic calculus, sparse learning, and Bayesian statistics. that require only the state estimates. The basic premise here is to treat the unknown input as a stochastic process, model it as a Gaussian white noise, and identify the stochastic differential equation (SDE); this essentially leads to identifying the drift component and diffusion component of a stochastic differential equation. We utilize stochastic calculus to decouple the drift and diffusion part; this allows parallel identification of the two components. Sparse learning and Bayesian statistics, on the other hand, ensures that the governing equation identified is sparse and interpretable; this is achieved by using the spike and slab prior^{32,33}. The proposed approach has several key features that can be encapsulated into the following points:

- First and foremost, unlike other approaches, the proposed approach require only the state measurements; no input measurements are needed.
- The proposed approach being Bayesian in nature estimates the posterior distribution. This allows computation of the epistemic uncertainty (aka predictive uncertainty) due to limited data. This feature is particularly important in design or decision making.
- Unlike non-Bayesian approaches, the proposed approach is autonomous in the sense that it doesn't require any human calibration and cross-validation.

- The proposed framework falls under the broad umbrella of interpretable machine learning and explainable AI. This allows generalization of the model to unseen environment and out-of-distribution data.

Results

Discovery of governing physics without input measurement

Natural dynamical processes are generally expressed in terms of higher order differential equations. However, in most of the data-driven machine learning techniques the actual system is not observed through its original space but identified in a projected space. In the current work, the actual dynamics of the underlying process is expressed in terms of the first order Itô . Let the projection can be realized by a map $\mathcal{T} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ such that $m < d$. If there exists a mapping \mathcal{T} , then any higher order system can be reduced to a set of SDEs of the form:

$$\dot{\mathbf{X}} = f(\mathbf{X}_t, t) + g(\mathbf{X}_t, t)\xi(t) \quad (1)$$

where, $f(\mathbf{X}_t, t)$ is the deterministic dynamics of the underlying process, $g(\mathbf{X}_t, t)$ is the volatility associated with the dynamics arising due to the stochastic external input, and $\xi(t)$ is the white noise. Let $\{B(t); t \geq 0\}$ be the Brownian motion with the properties $\langle B(t) \rangle = 0$, and $\langle B(s), B(t) \rangle = \min(s, t)$. Then the white noise can be defined as, $\xi(t) = \dot{B}(t)$. With the above definitions, an m -dimensional SDE is expressed as:

$$d\mathbf{X}_t = f(\mathbf{X}_t, t)dt + g(\mathbf{X}_t, t)d\mathbf{B}(t); \quad \mathbf{X}(t=t_0) = \mathbf{X}_0; \quad t \in [0, T] \quad (2)$$

where, $\mathbf{X}_t \in \mathbb{R}^m$ denotes the states of the process, $f(\mathbf{X}_t, t) : \mathbb{R}^m \mapsto \mathbb{R}^m$ is the drift vector, $g(\mathbf{X}_t, t) : \mathbb{R}^m \mapsto \mathbb{R}^{m \times n}$ is the diffusion matrix and $\mathbf{B}_j(t) \in \mathbb{R}^n$ is the n -dimensional Brownian motion. The Brownian motion $B(t)$ is generally not differential with respect to the process $X(t)$. However, their continuity is assumed to exists in mean square sense. Formally, if the interval $s \in [0, t]$ is partitioned into n -parts as, $\mathcal{P}_n(0, t) := [0 = s_0 < s_1 < \dots < s_n = t]$, then for a random process w the quadratic variation $Q_n(w, t) = \sum_i^n |w(s_i) - w(s_{i-1})|^2$ converges to $Q_n(w, t) \rightarrow Q(w, t)$. From the property of Brownian motions it can be found that $Q(w, t) = t$. This states that $B(t)$ has zero finite variation but non-vanishing quadratic variation. On the other hand, the deterministic functions have finite variation but zero quadratic variations. Leveraging this facts the drift and diffusion components of an SDE given in Eq. (2) can be expressed in terms of the sample time history as follows (more details are present in the *Methods* section and Appendix A):

$$\begin{aligned} f_i(\mathbf{X}_t, t) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[X_i(t + \Delta t) - \xi_i] \Big|_{x_k(t) = \xi_k} \quad \forall k = 1, 2, \dots, N \\ \Gamma_{ij}(\mathbf{X}_t, t) &= \frac{1}{2} \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[|x_i(t + \Delta t) - \xi_i| |x_j(t + \Delta t) - \xi_j|] \Big|_{x_k(t) = \xi_k} \quad \forall k = 1, 2, \dots, N \end{aligned} \quad (3)$$

where, $f_i(\mathbf{X}_t, t)$ is the i^{th} drift component and Γ_{ij} is the $(ij)^{th}$ component of the diffusion covaraince matrix $\boldsymbol{\Gamma} \in \mathbb{R}^{n \times n} := g(t, \mathbf{X}_t)g(t, \mathbf{X}_t)^T$. Let us assume that the analytical form of the drift and diffusion components can be obtained from measurement sample paths in terms of some basis functions. Let $\mathbf{L}^f \in \mathbb{R}^{N \times K} := [\nu_1^f(\mathbf{X}_t), \dots, \nu_K^f(\mathbf{X}_t)]$ and $\mathbf{L}^g \in \mathbb{R}^{N \times K} := [\nu_1^g(\mathbf{X}_t), \dots, \nu_K^g(\mathbf{X}_t)]$ are the library of candidate functions $\{\nu_k(\mathbf{X}_t), k = 1, \dots, K\}$. Here N is the length of sample path and the candidate functions can contain several functional forms such as polynomial, trigonometric, etc. Then, the drift and diffusion components can be expressed as linear combination of the library functions as,

$$\begin{aligned} f_i(\mathbf{X}_t, t) &= \theta_{i,1}^f \nu_1^f(\mathbf{X}_t) + \dots + \theta_{i,k}^f \nu_k^f(\mathbf{X}_t) + \dots + \theta_{i,K}^f \nu_K^f(\mathbf{X}_t) \\ \Gamma_{ij}(\mathbf{X}_t, t) &= \theta_{ij,1}^g \nu_1^g(\mathbf{X}_t) + \dots + \theta_{ij,k}^g \nu_k^g(\mathbf{X}_t) + \dots + \theta_{ij,K}^g \nu_K^g(\mathbf{X}_t) \end{aligned} \quad (4)$$

where, $\theta_{i,k}^f$ and $\theta_{ij,k}^g$ are the weights associated with the corresponding basis functions of drift and diffusion covariance components, respectively. In general, the libraries \mathbf{L}^f and \mathbf{L}^g can be different but one can use the same library for discovery of both drift and diffusion terms. The governing physics in the absence of input measurement is discovered in the proposed framework in terms of the SDEs. This requires identification of both the drift and diffusion components. Thus, one can construct two but separate sparse linear regression frameworks for the drift and diffusion terms and solve them independently. This yields,

$$\begin{aligned} \mathbf{Y}_i &= \mathbf{L}^f \boldsymbol{\theta}_i^f + \boldsymbol{\varepsilon}_i \\ \mathbf{Y}_{ij} &= \mathbf{L}^g \boldsymbol{\theta}_{ij}^g + \boldsymbol{\eta}_{ij} \end{aligned} \quad (5)$$

where, $\mathbf{Y}_i = f_i(\mathbf{X}_t, t)$ and $\mathbf{Y}_{ij} = \Gamma_{ij}(\mathbf{X}_t, t)$ are the target vectors of the sparse regression problem associated with i^{th} -drift component and $(ij)^{th}$ -diffusion covariance term, respectively, and, $\boldsymbol{\epsilon}_i$ and $\boldsymbol{\eta}_{ij}$ are the corresponding residual error vectors. The above equations are solved using sparse Bayesian linear regression. To understand, let us represent Eq. (5) as follows,

$$\mathbf{Y} = \mathbf{L}\boldsymbol{\theta} + \boldsymbol{\epsilon} \quad (6)$$

where, $\mathbf{Y} \in \mathbb{R}^N$ denotes the N-dimensional target vector, \mathbf{L} denotes the library of candidate functions, $\boldsymbol{\theta}$ is the weight vector, and $\boldsymbol{\epsilon} \in \mathbb{R}^N$ is the residual error vector representing the model mismatch error. On application of the Bayes formula yields,

$$P(\boldsymbol{\theta}|\mathbf{Y}) = \frac{P(\mathbf{Y}|\boldsymbol{\theta})}{P(\mathbf{Y})} \quad (7)$$

Modeling the mismatch error $\boldsymbol{\epsilon}$ as i.i.d Gaussian random variable with zero mean and variance σ^2 , the likelihood function is written as,

$$\mathbf{Y}|\boldsymbol{\theta}, \sigma^2 \sim \mathcal{N}(\mathbf{L}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_{N \times N}) \quad (8)$$

where, $\mathbf{I}_{N \times N}$ denotes the $N \times N$ identity matrix. The aim of the sparse regression is to find a governing model which contains only few terms in the final model. This is obtained by imposing sparsity promoting priors on the weight vector. In this work, the spike and slab (SS) distributions are used for promotion of sparsity in the solution. This suggest to take into account the effect of parameters of SS-prior on the probability of the weights, this is taken into account by constructing a bigger Bayesian hierarchical model. The structure of the Bayesian hierarchical model for the SS-prior is shown in Fig. 1. For the classification

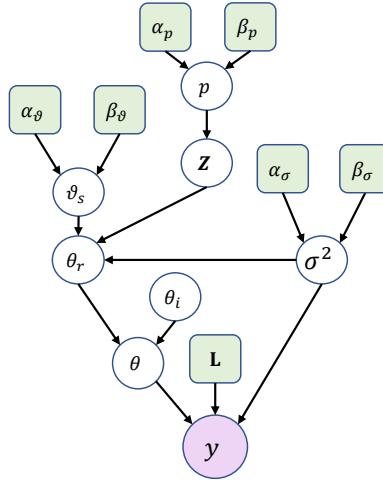


Figure 1. Hierarchical Bayesian network of the discontinuous spike and slab model for sparse linear regression. The variables in the green square boxes indicate the deterministic parameters and the variables in the white circles represents random variables. The term \mathbf{L} represents the library of candidate functions, and the parameters $\alpha_\theta, \beta_\theta, \alpha_p, \beta_p, \alpha_\sigma$, and β_σ indicates the hyperparameters of the priors of the hierarchical DSS model. Here p, ϑ_s and σ^2 scalar valued, and, the the variables Z and $\boldsymbol{\theta}$ are vector valued variables.

of the weights into spike and slab consider a latent variable Z such that it takes a value 1 if the weight fall into slab, otherwise, takes a value 0. Denoting $\boldsymbol{\theta}_r$ as the weight vector containing only those variables from $\boldsymbol{\theta}$ for which $Z_k = 1$, the DSS-prior is defined as^{31,32},

$$p(\boldsymbol{\theta}|Z) = p_{slab}(\boldsymbol{\theta}_r) \prod_{k, Z_k=0} p_{spike}(\theta_k) \quad (9)$$

where, the spike and slab distributions are defined as, $p_{spike}(\theta_k) = \delta_0$ and $p_{slab}(\boldsymbol{\theta}_r) = \mathcal{N}(\mathbf{0}, \sigma^2 \vartheta_s \mathbf{R}_{0,r})$ with $\mathbf{R}_{0,r} = \mathbf{I}_{r \times r}$. The hyperparameters $\alpha_\theta, \beta_\theta, \alpha_p, \beta_p, \alpha_\sigma$, and β_σ in the Fig. 1 are provided as a deterministic constants in the hierarchical model.

The random variables p_0 , ϑ_s , σ^2 and Z_k are simulated as follows,

$$p(\vartheta_s) = IG(\alpha_\vartheta, \beta_\vartheta) \quad (10)$$

$$p(Z_k|p_0) = Bern(p_0); k = 1 \dots K \quad (11)$$

$$p(p_0) = Beta(\alpha_p, \beta_p) \quad (12)$$

$$p(\sigma^2) = IG(\alpha_\sigma, \beta_\sigma) \quad (13)$$

The joint distribution of the random variables $\boldsymbol{\theta}$, \mathbf{Z} , ϑ_s , σ^2 , $p_0|Y$ are obtained from the Fig. 1 as,

$$\begin{aligned} p(\boldsymbol{\theta}, \mathbf{Z}, \vartheta_s, \sigma^2, p_0|Y) &= \frac{p(Y|\boldsymbol{\theta}, \sigma^2)p(\boldsymbol{\theta}|\mathbf{Z}, \vartheta_s, \sigma^2)p(\mathbf{Z}|p_0)p(\vartheta_s)p(\sigma^2)p(p_0)}{p(Y)} \\ &\propto p(Y|\boldsymbol{\theta}, \sigma^2)p(\boldsymbol{\theta}|\mathbf{Z}, \vartheta_s, \sigma^2)p(\mathbf{Z}|p_0)p(\vartheta_s)p(\sigma^2)p(p_0) \end{aligned} \quad (14)$$

where, $p(\boldsymbol{\theta}, \mathbf{Z}, \vartheta_s, \sigma^2, p_0|Y)$ denotes the joint distribution of the random variables, $p(Y|\boldsymbol{\theta}, \sigma^2)$ denotes the likelihood function, $p(\boldsymbol{\theta}|\mathbf{Z}, \vartheta_s, \sigma^2)$ is the prior distribution for the weight vector $\boldsymbol{\theta}$, $p(\mathbf{Z}|p_0)$ is the prior distribution for the latent vector \mathbf{Z} , $p(\vartheta_s)$ is the prior distribution for the slab variance ϑ_s , $p(\sigma^2)$ is the prior distribution for the noise variance, $p(p_0)$ is the prior distribution for the success probability p_0 and $p(Y)$ is the marginal likelihood. The Gibbs sampling technique is used to draw the random samples from the above joint distribution⁴⁴. For the Gibbs sampling the conditional distributions of the random variables are derived in Appendix B. For final selection of the model basis function marginal posterior probabilities are used, the details of which are briefly provided in *Methods* section.

The Gibbs sampler is initialised with the following initial values of the hyperparameters: $p_0^{(0)}=0.1$, $\vartheta^{(0)}=10$, and $\sigma^{2(0)}$ is set equal to the residual variance from ordinary least-squares regression. The initial vector of binary latent variables $\mathbf{Z}^{(0)}$ is computed by setting $\mathbf{Z}^{(0)} = [Z_1, \dots, Z_K]$ to zero and then activating the components of \mathbf{Z} that reduce the mean-squared error between the training data and the obtained model from ordinary least-squares. For this purpose a forward followed by a backward search algorithm is devised, where the backward search iterates through the activated components of initial latent vector in similar fashion to forward search. Given all the other parameters, the initial value of $\theta^{(0)}$ is obtained from the Eq. (14). For the commencement of the algorithm the deterministic prior parameters are set to the following values: $a_p=0.1$ and $b_p=1$ for the Beta prior on p_0 , $a_v=0.5$ and $b_v=0.5$ for inverse-Gamma prior on slab variance, and, $a_\sigma=10^4$ and $b_\sigma=10^4$ for inverse-Gamma prior on measurement noise. A Markov chains with 3000 samples is used for Gibbs sampling. The first 1000 samples are discarded as burn-in, and the remaining 2000 samples are used for posterior computation.

For all the demonstrations in this work, the data are simulated using the Euler-Maruyama scheme at a frequency of 1000Hz using the parameters listed in the Table 1. The noise in the measurements is modeled as N -dimensional sequence of zero-mean Gaussian white noise with a standard deviation equal to 5% of the standard deviation of the simulated quantities. In this work, the dictionary $\mathbf{L} \in \mathbb{R}^{N \times K}$ is constructed from 5 types of mathematical functions, each function representing a mapping of the m -dimensional state vector $\mathbf{X} = \{X_1, X_2, \dots, X_m\}$:

$$\mathbf{L}(\mathbf{X}) = \left[\begin{array}{ccccccc} \mathbf{1} & P^1(\mathbf{X}) & P^2(\mathbf{X}) & \dots & P^P(\mathbf{X}) & sgn(\mathbf{X}) & |\mathbf{X}| & \mathbf{X} \otimes |\mathbf{X}| \end{array} \right] \quad (15)$$

Here, $\mathbf{1} \in \mathbb{R}^N$ denotes the N -dimensional vector of 1, $P^P(\mathbf{X}) \in \mathbb{R}^{N \times m}$ denotes the set of terms present in the multinomial expansion $(X_1 + X_2 + \dots + X_m)^P$, $sgn(\mathbf{X}) \in \mathbb{R}^{N \times m}$ represents the signum functions of the form $sgn(X_i) \forall i = 1 \dots m$, $|\mathbf{X}| \in \mathbb{R}^{N \times m}$ denotes the absolute mapping of the states: $|X_i| \forall i = 1 \dots m$. The tensor product term $\mathbf{X} \otimes |\mathbf{X}| \in \mathbb{R}^{N \times 2m}$ represents the set of functions: $X_i \otimes |X_j| \forall i, j = 1 \dots m$. For the study, the length of the polynomial P is chosen as 6. The cardinality of the library is found as $K = (1 + |(X_1 + X_2 + \dots + X_m)^n| + 4m)$, where, the number of terms in the multinomial expansion can be found as, $|(X_1 + X_2 + \dots + X_m)^n| = {}^{n+m-1}C_{m-1}$. The complete architecture of the propose framework is illustrated in the Fig. 2.

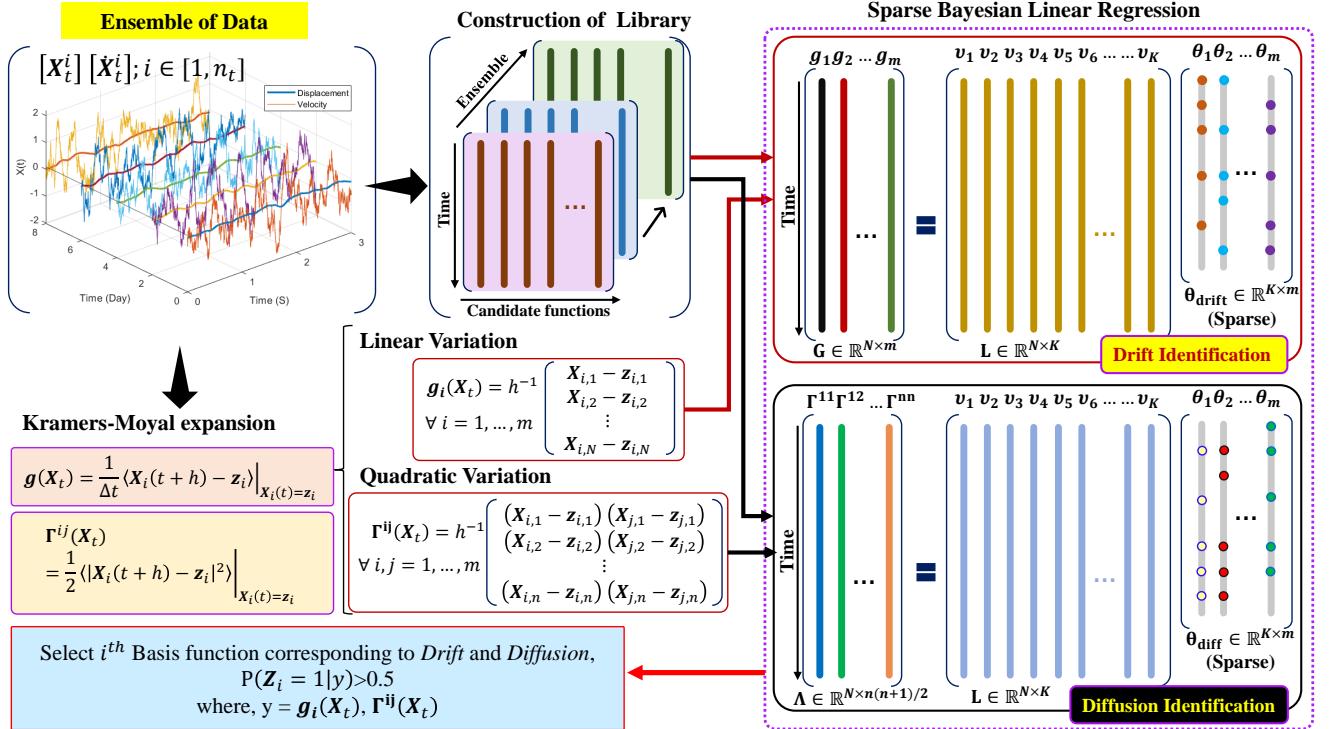


Figure 2. Schematic illustration of the proposed Bayesian-SDE framework for discovery of stochastic dynamical systems. The proposed framework integrates two key steps for identification of the stochastic systems in terms of the stochastic differential equations. The first key step, is the construction of the target vector for the sparse linear regression. For this purpose the proposed framework utilises the Kramers-Moyal expansion to estimate the drift $g(X_t, t)$ and diffusion terms $f(X_t, t)$ in terms of the linear and quadratic variations of the measured sample paths. The linear variation: $\lim_{\Delta t \rightarrow 0} \langle X(t + \Delta t) - z \rangle; X(t) = z$ indicates sum of the increments of the sample paths and the quadratic variation: $\lim_{\Delta t \rightarrow 0} \langle |X(t + \Delta t) - z|^2 \rangle; X(t) = z$ indicates sum of the square of the increments of the sample paths in limiting sense. The second step, includes formulation of the sparse Bayesian linear regression framework for the above target vector. Given an ensemble of time history of data the framework obtains a sequence of libraries by evaluating the candidate basis functions $\{v_k(\mathbf{X}_t); k = 1, 2, \dots, K\}$ over the ensembles. Ensemble mean is taken over all the libraries and the averaged library $\mathbf{L} \in \mathbb{R}^{N \times K}$ is utilized for performing the sparse linear regression. The candidate function in the library are parameterized by a weight vector θ . To identify each of the drift terms in an m -dimensional diffusion process this regression are performed m -times. The diffusion terms are not directly recoverable but the discovery of diffusion terms are possible in terms of its covariances. This requires $m(m+1)/2$ number of regression for discovery of all the terms in a covariance matrix. The elements of the weight vector θ is assigned a latent variable $Z_k; k = 1 \dots K$ in order to classify the weights into spikes and slabs. Post the Bayesian regression the marginal posterior inclusion probabilities (PIP) $p(Z_k = 1|Y)$ is estimated. In the final model only the basis functions whose corresponding marginal PIP values are higher than a threshold value is included.

Discovery of governing physics of example problems

In this section, the efficacy and robustness of the proposed Bayesian physics discovery framework is observed on a variety of representative stochastic non-linear dynamical systems. The examples taken include: (a) Black–Scholes SDE, (b) Duffing Van-der pol oscillator, and, (c) Two degree-of-freedom (DOF) base isolated shear building. For the equation discovery, it is assumed that the input forces $\dot{B}(t)$ are not measurable and only the noisy measurements of the displacements are available for physics discovery. The velocity component is obtained from the displacement vector through numerical differentiation such as fourth order central finite difference formula. Additionally a case study using Bouc-Wen oscillator is undertaken, where it is assumed that the hysteresis measurement is also not available. Therefore the proposed framework treats it as a partially observed system and tries to estimate the effect of the unobserved state without explicitly considering it in to the library of candidate functions. The results of the case studies undertaken in this work are presented in Fig. 3, 4 and Table 2. These results demonstrate sufficient accuracy in the discovered physics and robustness of the proposed framework to learn governing law from limited and noisy data.

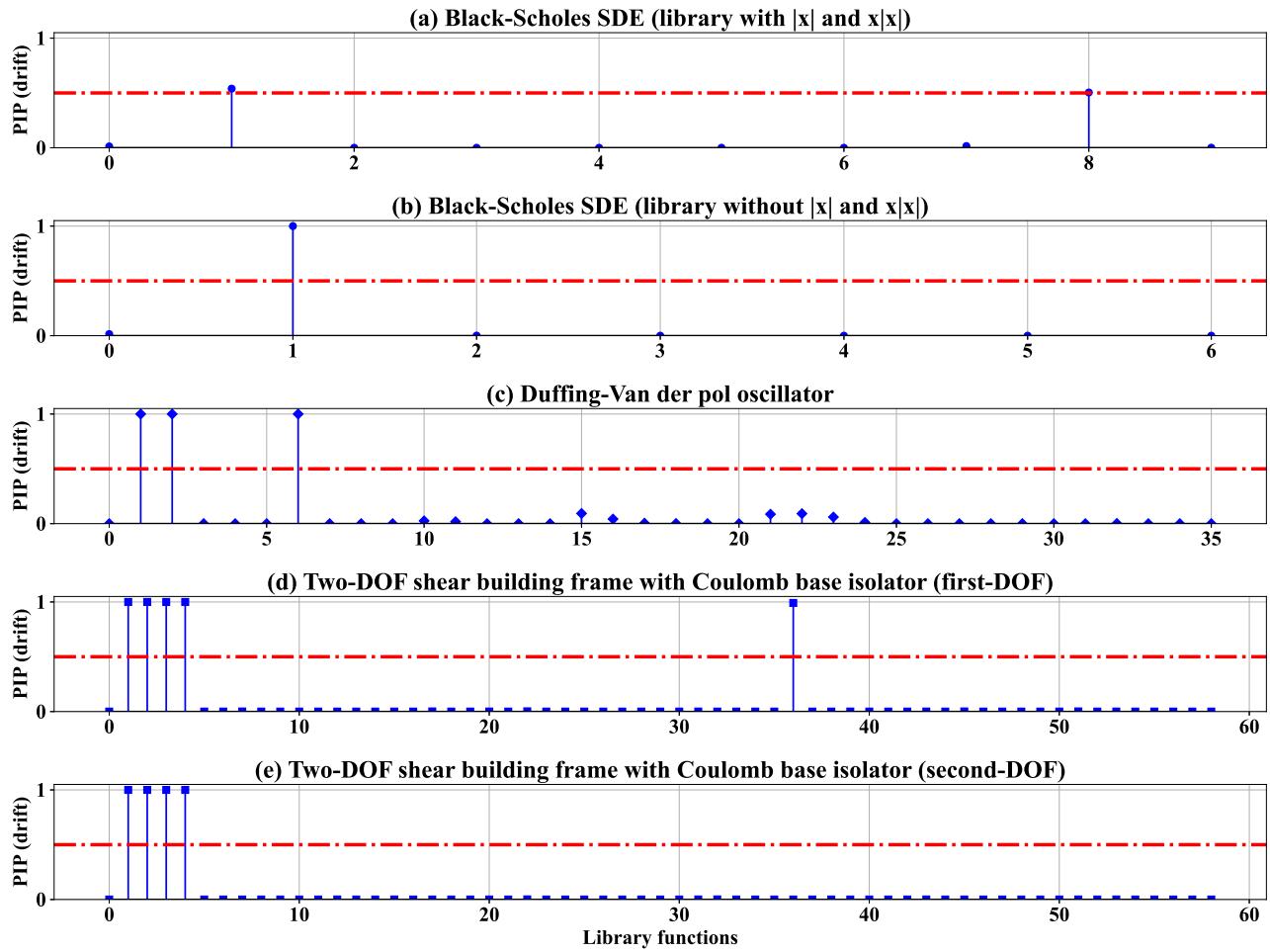


Figure 3. Discovery of the drift components of the governing physics from data corrupted with 5% white Gaussian noise based on marginal posterior inclusion probability (PIP), $P(Z_k = 1|Y)$. (a) Black-Scholes SDE: the library $\mathbf{L} \in \mathbb{R}^{N \times 9}$ consists 9 candidate basis functions out of which the basis functions x and $|x|$ are identified based on the criteria, $PIP > 0.5$. (b) Black-Scholes SDE: the library $\mathbf{L} \in \mathbb{R}^{N \times 6}$ is constructed as a collection of 6 basis functions. The correct drift component $X(t)$ is accurately identified with the almost sure probability $P(Z_1 = 1) = 1$. (c) Duffing-Van der pol oscillator: the library $\mathbf{L} \in \mathbb{R}^{N \times 36}$ has 36 candidate basis functions. The drift component is correctly identified as, $v(1) = X_1$, $v(2) = X_2$ and $v(6) = X_1^3$. (d) Two-DOF shear building (drift component of first DOF): the library $\mathbf{L} \in \mathbb{R}^{N \times 58}$ consists a total of 58 candidate basis functions, out of which 5 basis are selected for discovery of first drift component as, $v(1) = Y_1$, $v(2) = Y_2$, $v(3) = Y_3$, $v(4) = Y_4$ and $v(36) = \text{sgn}(Y_2)$. (e) Two-DOF shear building (drift component of first DOF): from the library $\mathbf{L} \in \mathbb{R}^{N \times 58}$ the basis functions are discovered as, $v(1) = Y_1$, $v(2) = Y_2$, $v(3) = Y_3$, $v(4) = Y_4$. The basis functions for both the drift components are selected with almost full probability $P(Z_k = 1|Y)$.

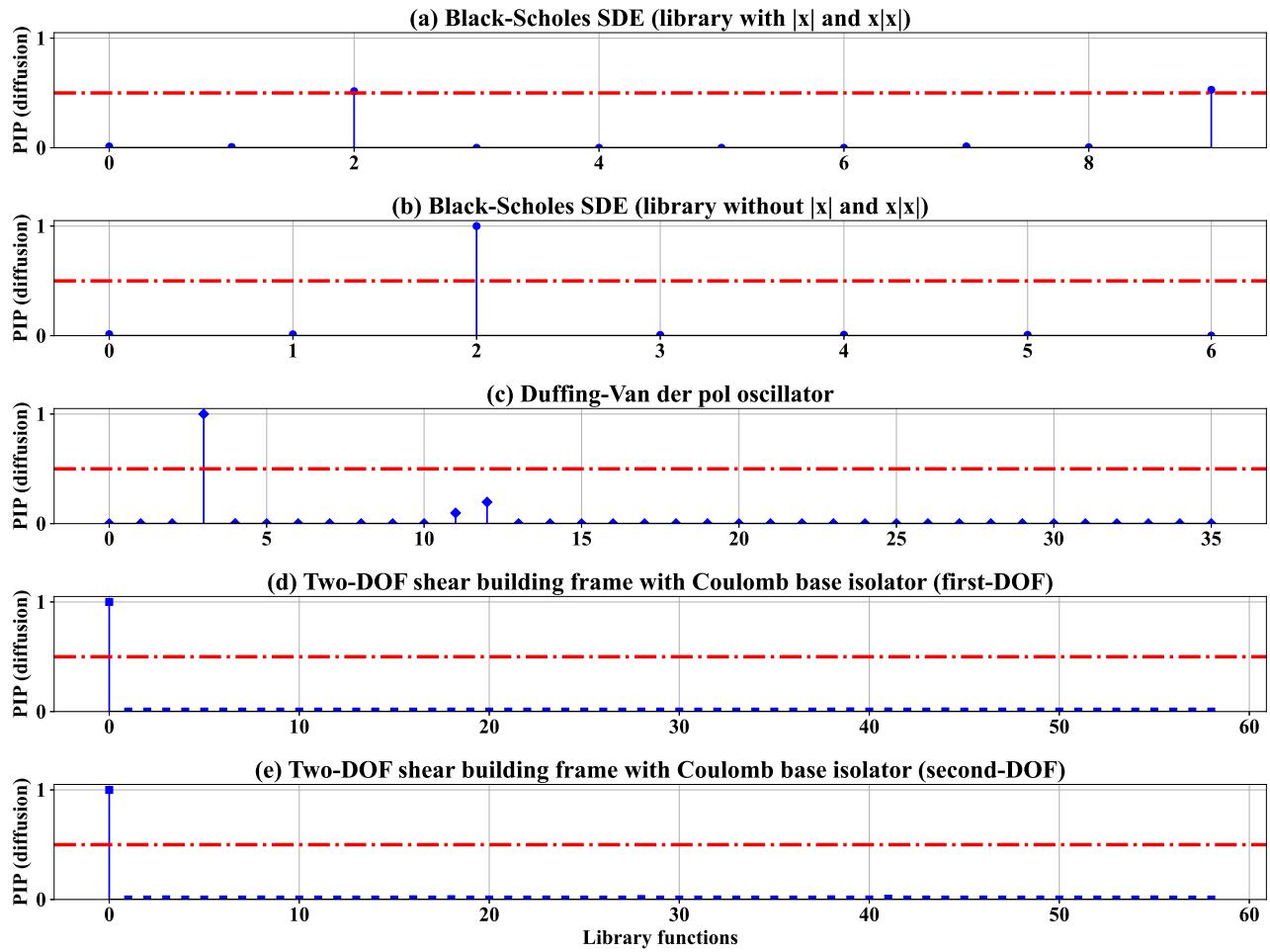


Figure 4. Discovery of the diffusion components of the governing physics from data corrupted with 5% white Gaussian noise based on marginal posterior inclusion probability (PIP), $P(Z_k = 1|Y)$. (a) Black-Scholes SDE: the weight vector $\theta \in \mathbb{R}^9$ for the library $\mathbf{L} \in \mathbb{R}^{N \times 9}$ is shown. The covariance of the diffusion component is identified as a combination of the terms $X^2(t)$ and $X|X|$. (b) Black-Scholes SDE: the weight vector $\theta \in \mathbb{R}^6$ has 6 elements out of which the correct covariance term $X^2(t)$ is discovered with almost sure probability $P(Z_2 = 1) = 1$. (c) Duffing-Van der pol oscillator: there are 36 elements in the weight vector $\theta \in \mathbb{R}^{36}$. The diffusion term is discovered correctly as $X_1^2(t)$ with almost sure probability $P(Z_2 = 1) = 1$. (d) Two-DOF shear building (drift component of first DOF): the weight vector $\theta \in \mathbb{R}^{58}$ has 58 elements corresponding to the 58 basis functions. The first diffusion component is accurately identified as $X_1^2(t)$ with almost sure probability $P(Z_2 = 1) = 1$. (e) Two-DOF shear building (drift component of second DOF): the second diffusion component $X_2^2(t)$ is accurately identified with almost sure probability $P(Z_2 = 1) = 1$.

Simulated systems	Drift parameters		Diffusion parameters
	Linear	Non-linear	
Black-Scholes	$\lambda = 2$	-	$\mu = 1$
Duffing-Van der pol	$m = 1, k = 2 \times 10^3, c = 2$	$\alpha = 10^5$	$\sigma = 10$
2-DOF non-linear	$m_1 = m_2 = 1, k_1 = 4 \times 10^3, k_2 = 2 \times 10^3$ $c_1 = c_2 = 2$	$\mu = 1, g = 9.81$	$\sigma_1 = 10, \sigma_2 = 10$
Bouc-Wen	$m = 1, k = 1 \times 10^4, c = 20, \lambda = 0.5$	$A_1 = A_2 = 0.5, A_3 = 1, \bar{n} = 3$	$\sigma = 2$

Table 1. Simulation parameters of the systems.

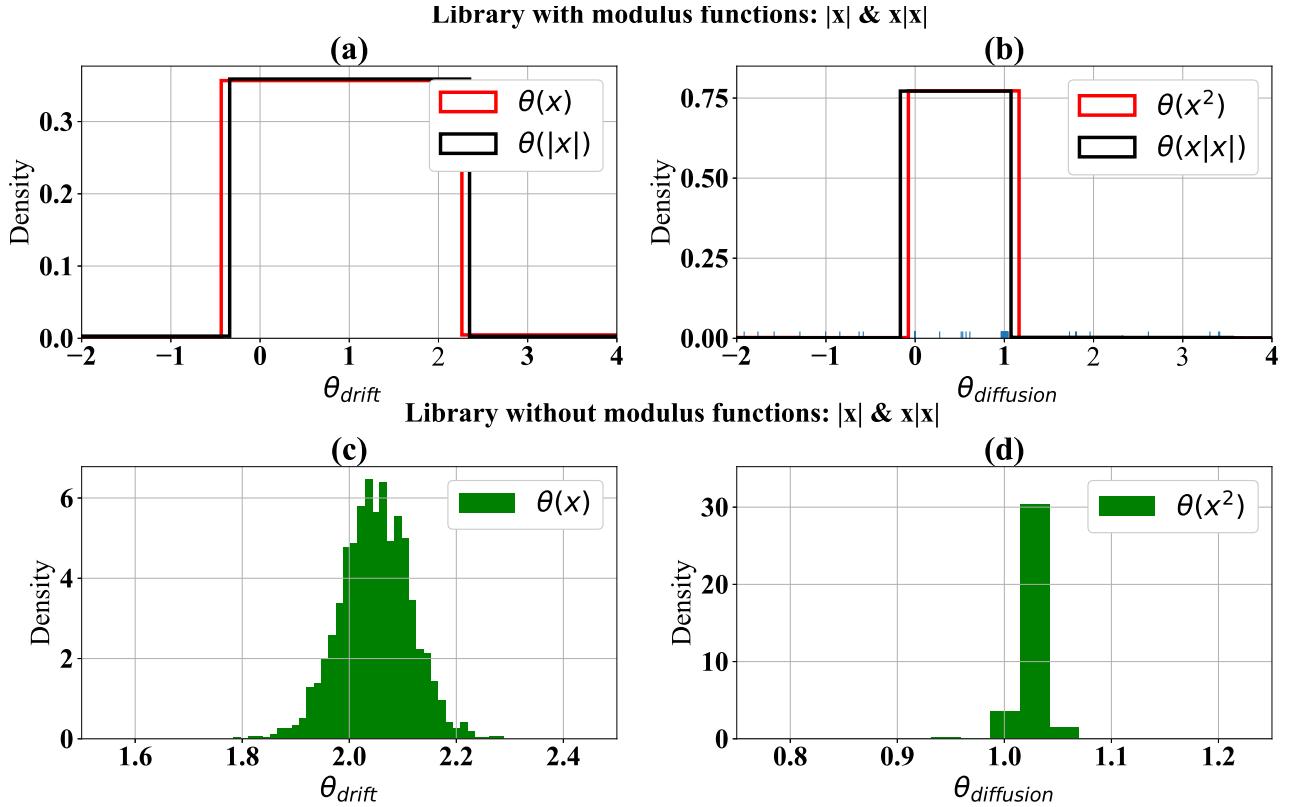


Figure 5. Black-Scholes SDE: posterior distributions of the weights of discovered basis functions. (a) The posterior distributions of θ_X and $\theta_{|X|}$. Posterior of both emulates the uniform distribution, whose means lie around the value 1. As a combination the basis θ_X and $\theta_{|X|}$ identifies the drift. (b) The posterior distributions of θ_{X^2} and $\theta_{X|X|}$. The posteriors emulate uniform distribution. The mean gives a value around 0.5. Thus, as a combination they discovered the diffusion. (c) The posterior of weight θ_X . The mean is obtained as 2.04, which exactly identifies the drift with almost sure probability $P(Z_1 = 1) = 1$. (d) The posterior of weight θ_{X^2} whose mean provides the value of the parameters μ as 1.02.

Correct SDEs:

Black-Schole: $dX(t) = 2X(t)dt + X(t)dB(t)$

Duff.-Van der pol: $dX_2(t) = -(1000X_1(t) + 2X_2(t) + 100000X_1^3(t))dt + 10X(t)dB(t)$

2-DOF non-linear: $\begin{cases} dY_2(t) = (-6000Y_1(t) - 4Y_2(t) - 9.81\text{sgn}(Y_2(t)) + 2000Y_3(t) + 2Y_4(t))dt + 10dB_1(t) \\ dY_4(t) = (2000Y_1(t) + 2Y_2(t) - 2000Y_3(t) - 2Y_4(t))dt + 10dB_2(t) \end{cases}$

Identified SDEs:

Black-Schole: $\begin{cases} dX(t) = 2.05X(t)dt + 1.02X(t)dB(t) \\ dX(t) = (0.47X(t) + 0.53|X(t)|)dt + (1.09X^2(t) + 0.95X(t)|X(t)|)dB(t) \end{cases}$

Duff.-Van der pol: $dX_2(t) = -(1000.02X_1(t) + 1.99X_2(t) + 99885.10X_1^3(t))dt + 10.31X(t)dB(t)$

2-DOF non-linear: $\begin{cases} dY_2(t) = (-6000.30Y_1(t) - 3.97Y_2(t) - 9.89\text{sgn}(Y_2(t)) + 1999.73Y_3(t) + 2.09Y_4(t))dt \\ + 10.03dB_1(t) \\ dY_4(t) = (2000.65Y_1(t) + 1.99Y_2(t) - 1999.06Y_3(t) - 1.99Y_4(t))dt + 10.11dB_2(t) \end{cases}$

Table 2. Summary of the results of Bayesian equation discovery of stochastic dynamical systems. The parameters of the identified systems denote the expected value of the weights after discarding the burn-in samples. Here, the first 1000 samples are taken as burn-in samples and therefore discarded for obtaining the final stationary distribution. In this work, the burn-in samples are taken as first 1000 MCMC samples.

Black–Scholes SDE

A Black–Scholes SDE (formulated as geometric Brownian motion) is a continuous-time stochastic process where the logarithm of the randomly varying quantity follows a Brownian motion. The Black–Scholes SDE is frequently used in stock market for modelling of the evolution of stock price of an underlying asset³⁴. The Black–Scholes SDE has the drift $f(t, X_t) = \lambda X$ and diffusion $g(t, X) = \mu X$ with $\lambda > 0$ and $\mu > 0$ being the real constants. Towards this the Black–Scholes SDE is defined as a geometric Brownian motion as follows:

$$\frac{dX(t)}{X(t)} = \lambda dt + \mu dB(t) \quad X(t=t_0) = X_0; \quad t \in [0, T] \quad (16)$$

where $B = \{B(t); t \geq 0\}$ is the Brownian motion. The solution to the above SDE is attempted using Euler Maruyama (EM) with a uniform time step size Δt . The one-step approximation of the above SDE is then expressed by EM scheme as:

$$X(n+1) = X(n) + \lambda X(n)\Delta t + \mu X(n)\Delta B(n) \quad (17)$$

where, $\Delta t = t_{n+1} - t_n$ is the time increment and $\Delta B(n) = B(t_{n+1}) - B(t_n)$ is the Brownian increment. Identifying the diffusion term as, $g(\mathbf{X}_t, t) = \mu X(n)$, the variance of the diffusion is obtained: $\Gamma = \mu^2 X^2(n)$. The results of this case study are plotted in the Figures 3(a) and 4(a). Figure 3(a) depicts that there are two identified basis function X and $|X|$. However, from the Eq. (16) it can be understood that only one function X should have been identified as the driving function of the Black–Scholes SDE. To understand this discrepancy one needs to consider that the evolution of the random variable $X(t)$ in the Black–Scholes SDE follows Geometric Brownian motion (GBM) and GBMs always assume positive values (e.g. real stock price). Since the characteristics of both the functions X and $|X|$ in this case are same, thus the algorithm identifies both the functions with almost equal probability. The similar phenomenon is observed in case of the identification of the diffusion term too. Figures 4(b) depicts almost equal contribution from the functions X^2 and $X|X|$, however, only the term X^2 should have been actually identified. Upon considering the corresponding parameter values of the basis functions X and $|X|$ for the drift component, and, X^2 and $X|X|$ for the diffusion component, from the Fig. 5(a) and (b), the Black–Scholes SDE in Eq. (16) can be easily simulated for any unseen environmental conditions.

To verify that this is not a limitation of the proposed scheme, this numerical demonstration is repeated in the Fig. 3(b) and 4(b) without considering the functions $|X|$ and $X|X|$ in the library. The parameters of the basis functions for the drift and diffusion components are portrayed in Fig. 5(c) and 5(d), respectively. The results clearly shows that the proposed scheme is able to identify the basis functions along with their associated parameters λ and μ without any significant error. As a consequence of the above case study, one can ask a question about what basis functions are to be considered in the library. One of the possibilities could be to visualize the the time series data and then take scientific and engineering judgements. For example, if the time history data show no zero crossing then one can opt for not-considering the functions which shares similar properties (for e.g. $|X|$ and X , and $X|X|$ and X^2). Besides one can also choose to consider due to the fact that the similar terms will have equal probability of occurrence and as a combination will demonstrate the observed phenomenon. However, for future prediction the model will be applicable to process depicting zero crossing only.

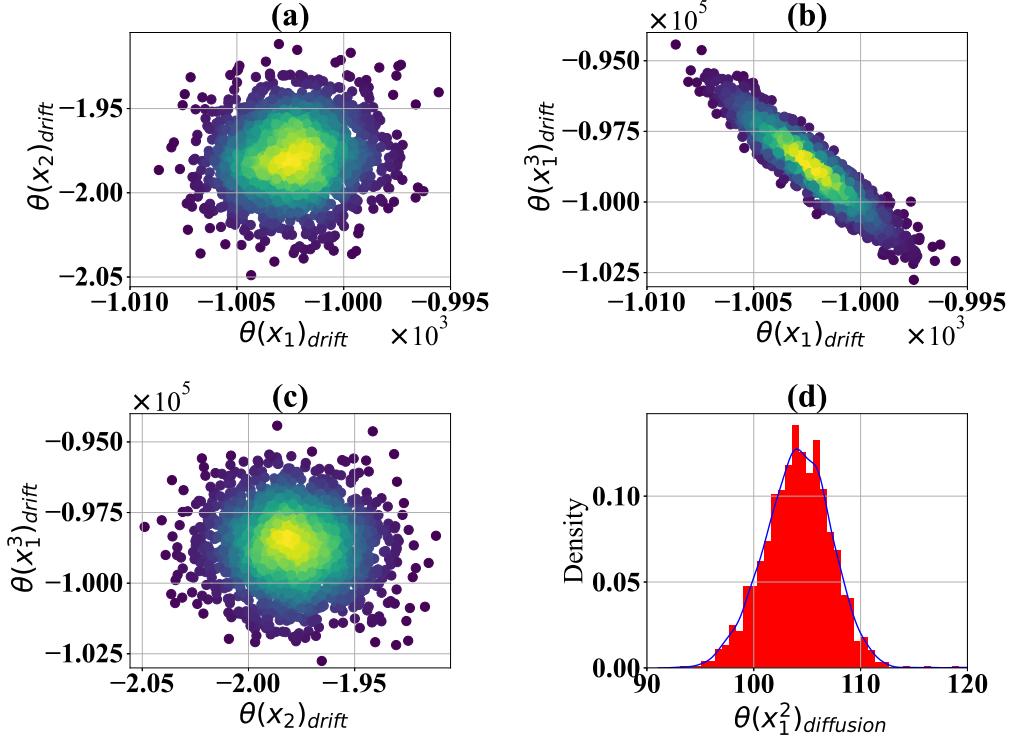


Figure 6. Duffing-Van der pol oscillator: Posterior distributions of the weights of the selected basis functions. (a) Joint posterior distribution of $\theta(X_1)$ and $\theta(X_2)$. The yellow region indicates the mean of the joint distribution. This weights $\theta(X_1)$ and $\theta(X_2)$ represents the parameters k and c . The expected values of the identified weights are, $E(\theta(X_1)) = 1000.2$ and $E(\theta(X_2)) = 1.99$. (b) Joint posterior distribution of $\theta(X_1)$ and $\theta(X_1^3)$. The expected values of the weights are $E(\theta(X_1)) = 1000.2$ and $E(\theta(X_1^3)) = 99885.1$. (c) Joint posterior distribution of $\theta(X_2)$ and $\theta(X_1^3)$. he expected values of the weights are $E(\theta(X_2)) = 1.99$ and $E(\theta(X_1^3)) = 99885.1$. (d) The histogram of the diffusion parameter. The mean value of the histogram is 106.47 which represents the term $g^2(X_t, t)$. Upon conversion the parameters σ/m can be obtained as 10.29.

Duffing-Van der pol oscillator

The second order non-linear hardening Duffing-Van der pol oscillator is considered, in this section. The Duffing-Van der pol oscillator draws its physical relevance from the models in flow-induced structural vibration problems. The Duffing-Van der pol oscillator has a cubic dissipating force and it is driven by multiplicative white noise. Towards this, the equation of motion of the undertaken system is expressed as:

$$m\ddot{X}(t) + c\dot{X}(t) + kX(t) + \alpha X^3(t) = \sigma X(t)\dot{B}(t) \quad X(t=t_0) = X_0; \quad t \in [0, T] \quad (18)$$

where, m , c and k are the mass, damping and stiffness parameters of the oscillator. Here, α is a real-valued parameter associated with the cubic non-linearity, $\sigma \geq 0$ is the strength of the multiplicative white noise, and $B = \{B(t); t \geq 0\}$ is the Brownian motion. The time derivative of the Brownian $\dot{B}(t)$ here represents the white noise. With a state-space transformation $X = X_1$, and $\dot{X} = X_2$, where X_1 and X_2 represents the displacement and velocity, the corresponding two first order Itô-stochastic SDEs for the dynamical system are derived as:

$$\begin{bmatrix} dX_1(t) \\ dX_2(t) \end{bmatrix} = \begin{bmatrix} X_2(t) \\ -\frac{1}{m}(kX_1(t) + cX_2(t) + \alpha X_1^3(t)) \end{bmatrix} dt + \begin{bmatrix} 0 \\ \frac{\sigma}{m}X(t) \end{bmatrix} dB(t) \quad (19)$$

In the above equation it be noticed that the first equation provides only information of the velocity component, thus not considered for discovery of equation. For estimating the drift and diffusion terms only the evolution of second variable $X_2(t)$ is used to construct the target linear and quadratic variation vectors from the sample paths using the Kramers-Moyal formulae. Here, the variance of the diffusion term is identified as, $\Gamma = (\sigma^2 X^2(t))/m^2$. The results for the basis function selection are displayed in Fig. 3(c) and 4(c), respectively, while their associated parameters are plotted in Fig. 6. From the Fig. 3(c) and 4(c),

it is evident that the proposed framework is able to accurately discover the basis functions for the drift and diffusion terms as, $\{X_1(t), X_2(t), X_3^3(t)\}$ and $\{X_1(t), X_2(t)\}$, respectively. Afterwards, using the reverse statespace transformation the governing physics in terms of the second order differential equation can be easily obtained.

The pairwise joint posterior distributions of weights associated with the discovered basis functions are presented in Fig. 6. In the subplots (a), (b) and (c), the weights of discovered drift functions are shown. It is evident in this figures that the actual parameters associated with the drift component of DVP oscillator as listed in Table 1 is correctly identified with very small relative error. The subplot (d) depicts the posterior distribution of the identified diffusion basis function. The mean value of the posterior distribution represents the term $g^2(X_t, t)$, which in this case can be noticed as $g(X_t, t) = X(t)$. The mean value here is obtained as approximately 106.47. By performing the square root operation over the mean, one can approximately get the diffusion value $\sigma = 10.31$. This results, verifies that the proposed scheme can be successfully applied to discover the governing physics of non-linear oscillators when subjected to parametric excitation effectively.

Two-DOF base isolated Shear Building

The responses of civil engineering structures often exhibit the characteristics of stochastic non-linear dynamical systems due to the random external excitation. For the purpose of dynamical analysis the building structures are dominantly represented through a spring-mass-dashpot model. Here, a simple 2-DOF base isolated structure idealised as spring-mass-dashpot is considered. The system is a linear base isolated structure where the structure is connected to the foundation through a Coulomb friction-type base isolator. Under this considerations, the governing equation of motion of the system is written as,

$$\begin{aligned} m_1 \ddot{X}_1(t) + c_1 \dot{X}_1(t) + \mu m_1 g \text{sgn}(\dot{X}_1(t)) + k_1 X_1(t) + c_2 (\dot{X}_1(t) - \dot{X}_2(t)) + k_2 (X_1(t) - X_2(t)) &= \sigma_1 \dot{B}_1(t) \\ m_2 \ddot{X}_2(t) + c_2 (\dot{X}_2(t) - \dot{X}_1(t)) + k_2 (X_2(t) - X_1(t)) &= \sigma_2 \dot{B}_2(t) \\ X(t = t_0) = X_0; \quad t \in [0, T] \end{aligned} \quad (20)$$

where, $\mu m_1 g \text{sgn}(\dot{X}_1(t))$ is the Coulomb friction force arising due to the sliding of bearings in Coulomb oscillator, and, $\{m_i, i = 1, 2\}$, $\{c_i, i = 1, 2\}$ and $\{k_i, i = 1, 2\}$ are the mass, damping and stiffness of i^{th} -DOF. Here, $\sigma_1 \geq 0$ and $\sigma_2 \geq 0$ are the strength of the white noise $\dot{B}_1(t)$ and $\dot{B}_2(t)$, where $B_1(t)$ and $B_2(t)$ are the independent Brownian motion. A statespace transformation is considered for the shear building as, $X_1 = Y_1$, $\dot{X}_1 = Y_2$, $X_2 = Y_3$, and $\dot{X}_2 = Y_4$, where Y_1 and Y_3 denotes the displacements, and, Y_2 and Y_4 denotes the velocities of the two DOF. By representing the complete state vector as $\mathbf{Y} = [Y_1, Y_2, Y_3, Y_4]$ the first order Itô-stochastic SDEs are obtained as follows:

$$\begin{bmatrix} dY_1(t) \\ dY_2(t) \\ dY_3(t) \\ dY_4(t) \end{bmatrix} = \begin{bmatrix} Y_2(t) \\ -\frac{k_1}{m_1} Y_1(t) - \frac{c_1}{m_1} Y_2(t) - \mu g \text{sgn}(Y_2(t)) - \frac{k_2}{m_1} (Y_1(t) - Y_3(t)) - \frac{c_2}{m_1} (Y_2(t) - Y_4(t)) \\ Y_4(t) \\ -\frac{k_2}{m_2} (Y_3(t) - Y_1(t)) - \frac{c_2}{m_2} (Y_4(t) - Y_2(t)) \end{bmatrix} dt + \begin{bmatrix} 0 & 0 \\ \frac{\sigma_1}{m_1} & 0 \\ 0 & 0 \\ 0 & \frac{\sigma_2}{m_2} \end{bmatrix} \begin{bmatrix} dB_1(t) \\ dB_2(t) \end{bmatrix} \quad (21)$$

For the interest of the proposed scheme, the covariance of the diffusion matrix is obtained in Eq. (22). In this case study, the equation discovery involves identification of two drift and two covariance terms (one for each of the DOFs). Following to the case study of Duffing-Van der pol oscillator, the target linear variation vector for discovery of $i^{th}; \forall i = 1, \dots, 2$ drift term needs to be constructed from the velocity component of the corresponding DOF. For discovery of the $(ij)^{th}; \forall i, j = 1, \dots, 2$ component of the covariance matrix, the quadratic variation matrix is constructed from the pointwise product of i^{th} and j^{th} response vectors.

$$\Gamma = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{\sigma_1^2}{m_1^2} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{\sigma_2^2}{m_2^2} \end{bmatrix} \quad (22)$$

The results for the identification of corresponding basis functions for both the drift terms are presented in Fig. 3(d) and (e). From the Fig. 3(d), it can be observed that the identified basis functions are $\{Y_1(t), Y_2(t), Y_3(t), Y_4(t), \text{sgn}(Y_2(t))\}$, which correctly matches with the terms in the equation of motion of first DOF. Similarly, from Fig. 3(e) it can be verified that the proposed scheme has correctly identified the basis functions of the second drift component as $\{Y_1(t), Y_2(t), Y_3(t), Y_4(t)\}$. Figure 7 and 8, shows the pairwise joint posterior distributions of the weights associated with identified first and second drift basis functions,

respectively. From these figures it is evident that the proposed scheme is able to identify the parameters of the basis functions as listed in Table 1 with sufficient accuracy. In order to find the explicit values of system stiffness and damping parameters, the parameters of second DOF can be identified first from the Figures 7(h), 7(e) and 7(f). Thereafter, the parameters of first DOF can be extracted easily by following the relations given in Eq. (21).

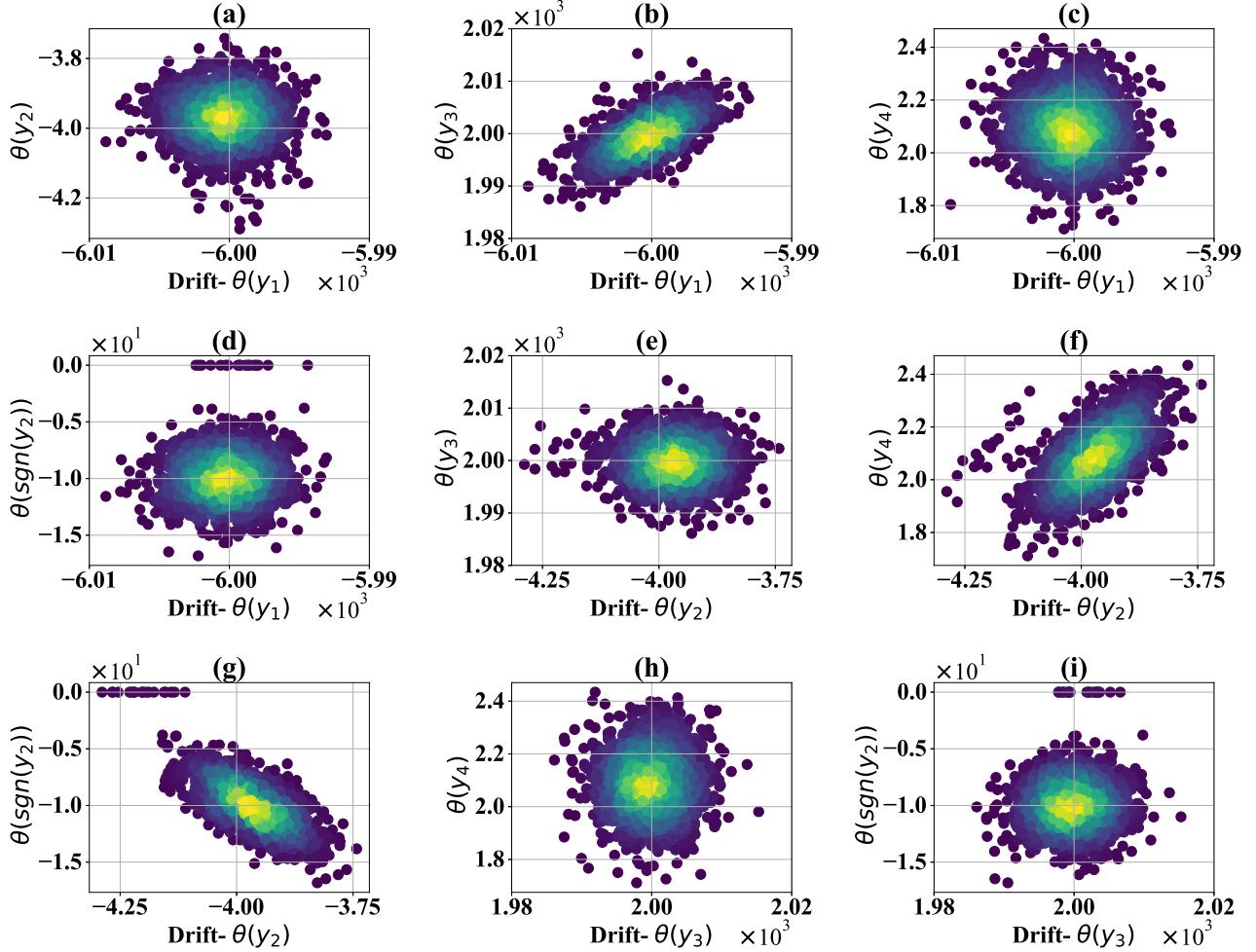


Figure 7. Non-linear two-DOF oscillator (first drift component): Pairwise joint posterior distributions of the parameters θ_k of the selected basis functions. (a) Joint posterior of the weights $\theta(Y_1)$ and $\theta(Y_2)$. The yellow color denotes the mean value of the joint distribution. The expected values of the parameters are identified as $E(\theta(Y_1)) = -6000.30$ and $E(\theta(Y_2)) = -3.97$. (b) Joint posterior of the weights θ_{Y_1} and θ_{Y_3} . The expected values of the parameters are identified as $E(\theta(Y_1)) = -6000.30$ and $E(\theta(Y_3)) = 1999.73$. (c) Joint posterior of the weights θ_{Y_1} and θ_{v_1} . The expected values of the parameters are identified as $E(\theta(Y_1)) = -6000.30$ and $E(\theta(Y_4)) = 2.09$. (d) Joint posterior of the weights θ_{Y_1} and $\theta_{sgn(Y_2)}$. The expected values of the parameters are identified as $E(\theta(Y_1)) = -6000.30$ and $E(\theta(sgn(Y_2))) = -9.89$. (e) Joint posterior of the weights θ_{Y_2} and θ_{Y_3} . The expected values of the parameters are identified as $E(\theta(Y_2)) = -3.97$ and $E(\theta(Y_3)) = 1999.73$. (f) Joint posterior of the weights θ_{Y_2} and θ_{Y_4} . The expected values of the parameters are identified as $E(\theta(Y_2)) = -3.97$ and $E(\theta(Y_4)) = 2.09$. (g) Joint posterior of the weights θ_{Y_2} and $\theta_{sgn(Y_2)}$. The expected values of the parameters are identified as $E(\theta(Y_2)) = 1999.73$ and $E(\theta(sgn(Y_2))) = -9.89$. (h) Joint posterior of the weights θ_{Y_3} and θ_{Y_4} . The expected values of the parameters are identified as $E(\theta(Y_3)) = 1999.73$ and $E(\theta(Y_4)) = 2.09$. (i) Joint posterior of the weights θ_{Y_3} and $\theta_{sgn(Y_2)}$. The expected values of the parameters are identified as $E(\theta(Y_3)) = 1999.73$ and $E(\theta(sgn(Y_2))) = -9.89$. The weights $\{\theta_{Y_1}, \theta_{Y_2}, \theta_{Y_3}, \theta_{Y_4}, \theta_{sgn(Y_2)}\}$ corresponds to the basis functions $\{v_1, v_2, v_3, v_4, v_5, v_6\}$. They represents the relative value of the parameters $\{k_1, k_2, k_3, k_4, \mu mg\}$. The negative values arises due to the coupling of the systems states.

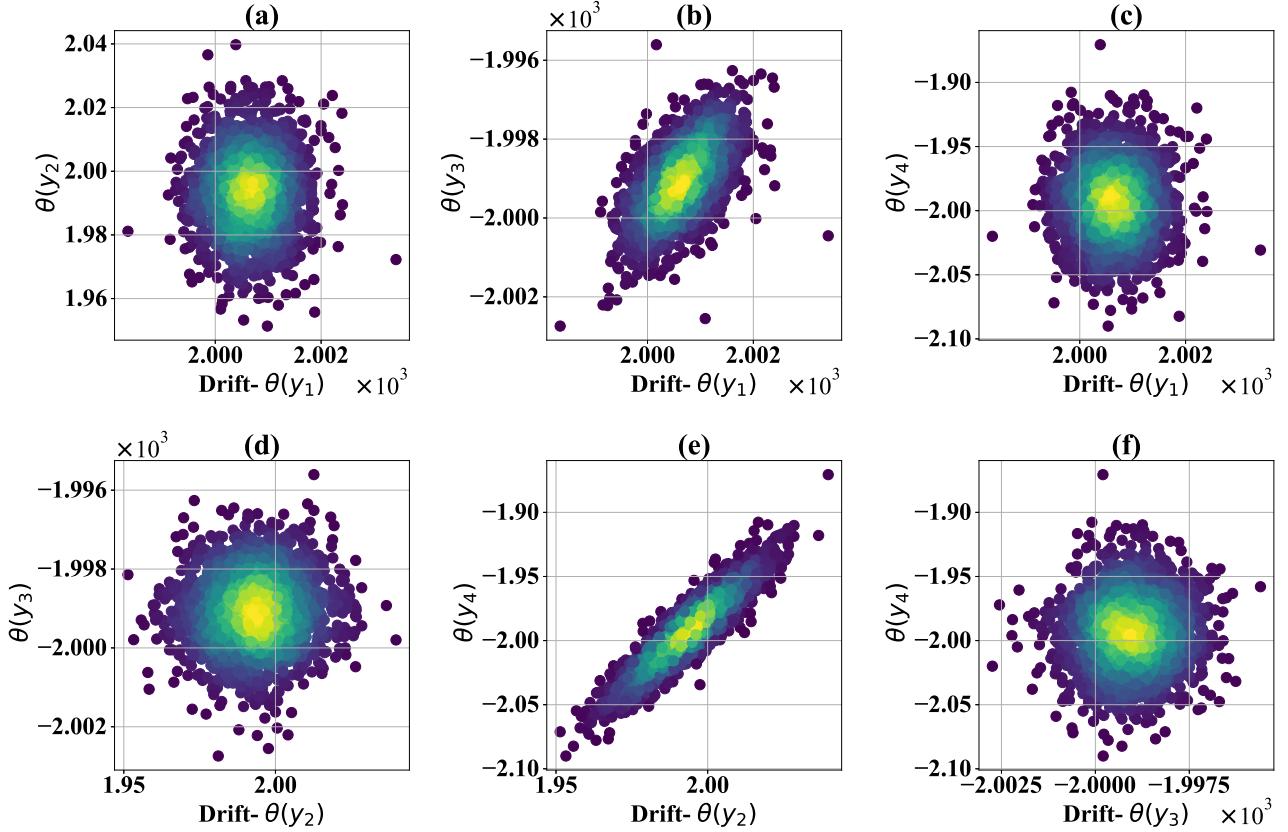


Figure 8. Non-linear two-DOF oscillator (second drift component): Pairwise joint posterior distributions of the parameters θ_k of the selected basis functions. (a) Joint posterior of the weights $\theta(Y_1)$ and $\theta(Y_2)$ corresponding to the basis functions v_1 and v_2 . The yellow color denotes the mean value of the joint distribution. The expected values of the parameters are identified as $E(\theta(Y_1))= 2000.65$ and $E(\theta(Y_2))= 1.99$. (b) Joint posterior of the weights θ_{Y_1} and θ_{Y_3} corresponding to the basis functions v_1 and v_3 . The expected values of the parameters are identified as $E(\theta(Y_1))= 2000.65$ and $E(\theta(Y_3))= -1999.06$. (c) Joint posterior of the weights θ_{Y_1} and θ_{Y_4} corresponding to the basis functions v_1 and v_4 . The expected values of the parameters are identified as $E(\theta(Y_1))= 2000.65$ and $E(\theta(Y_4))= -1.99$. (d) Joint posterior of the weights θ_{Y_2} and θ_{Y_3} corresponding to the basis functions v_2 and v_3 . The expected values of the parameters are identified as $E(\theta(Y_2))= 1.99$ and $E(\theta(Y_3))= -1999.06$. (e) Joint posterior of the weights θ_{Y_2} and θ_{Y_4} corresponding to the basis functions v_2 and v_4 . The expected values of the parameters are identified as $E(\theta(Y_2))= 1.99$ and $E(\theta(Y_4))= -1.99$. (f) Joint posterior of the weights θ_{Y_3} and θ_{Y_4} corresponding to the basis functions v_3 and v_4 . The expected values of the parameters are identified as $E(\theta(Y_3))= -1999.06$ and $E(\theta(Y_4))= -1.99$. The weights $\{\theta_{Y_1}, \theta_{Y_2}, \theta_{Y_3}, \theta_{Y_4}\}$ represents the relative value of the parameters $\{k_1, k_2, k_3, k_4\}$. The negative values arises due to the coupling of the systems states.

The identification results for the basis functions along with the posterior distribution of their parameters for both the diffusion terms are depicted in Fig. 9. Figures 9(a) and 9(b) presents the results of the first diffusion term. The identification results of the second diffusion term are shown in Figures 9(c) and 9(d). In both the cases, it is evident that the proposed scheme has correctly identified the basis function as 1. The posterior distributions the covariance terms, Γ^{11} and Γ^{22} is plotted in Fig. 9(b) and 9(d) whose mean values represents σ_1^2/m_1^2 and σ_2^2/m_2^2 . The summery of the results are presented in Table 2.

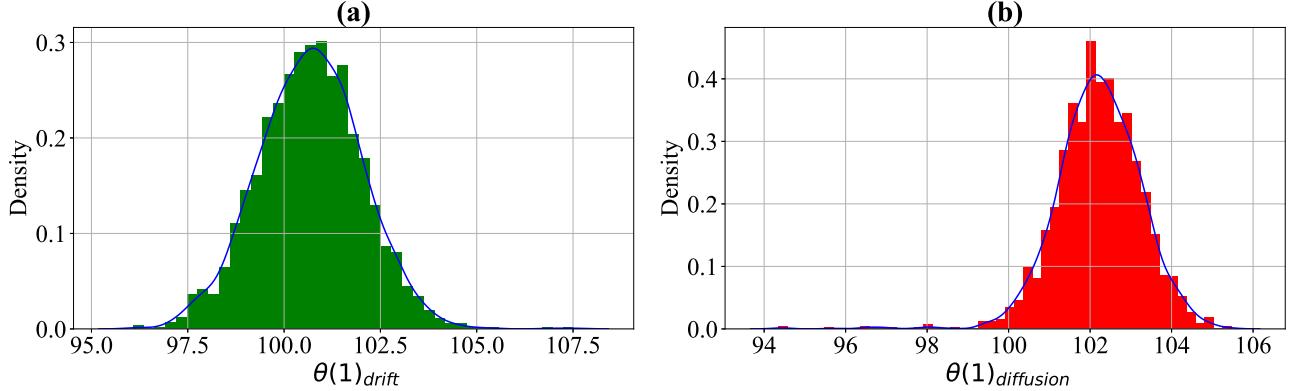


Figure 9. Two-DOF shear building: posterior distributions of the weights associated with the discovered basis functions of diffusion components. (a) First DOF: posterior distribution of the basis function θ_1 . The mean of the distribution denotes the expected value of σ_1^2/m_1^2 , which in this case is observed from the figure as approximately 101. (b) Second-DOF: posterior distribution of the basis function θ_1 . The mean is observed as 102, which is equivalent to the term σ_2^2/m_2^2 . By performing the square root operation over the mean values one can easily obtain the desired diffusion components.

Case study: Stochastic non-linear system with partially observed variables

In this section, a case study has been undertaken when a measurement of the states of a system is not possible. Let us consider a Bouc-Wen oscillator where the system is described using three states namely, displacement and velocity of the main mass, and the hysteresis displacement associated with the isolator^{35,36}. In this context, it can be understood that the measurement of the hysteresis displacement is practically intractable and only the displacement and velocity states of the main system is available. Further, it can be appreciated that measuring the velocity state of a mechanical system using sensors are practically challenging task which leaves the only available option is displacement. Thus, a noble methodology should be able to identify the physics of the system using only the directly measurable states only, since the derivative quantities may not represent the motion of the original system appropriately. Towards this, a first step is taken using the proposed framework in which a base isolated sdof spring-mass-dashpot system is undertaken and only the displacement state is observed. Consider, the governing equation of motion of a sdof Bouc-Wen oscillator as³⁵,

$$m\ddot{X}(t) + c\dot{X}(t) + k(1 - \lambda)Z(t) = \sigma\dot{B}(t) \quad X(t = t_0) = X_0; \quad t \in [0, T] \quad (23)$$

where, m , c , and k are the mass, damping and stiffness of the system, $X(t)$ is the system state, λ is a factor that defines the participation of the elastic force F_e and hysteresis force F_h , $F_e(t) = k\lambda X(t)$, $F_h(t) = k(1 - \lambda)Z(t)$, and $Z(t)$ is the hysteresis displacement. The evolution of the hysteresis parameter is defined using the non-linear differential equation:

$$Z(t) = -A_1 Z(t) |\dot{X}(t)| |Z(t)|^{\bar{n}-1} - A_2 \dot{X}(t) |Z(t)|^{\bar{n}} + A_3 \dot{X}(t) \quad Z(t = t_0) = Z_0; \quad t \in [0, T] \quad (24)$$

The positive exponential parameter n defines the smoothness of the transition from elastic to the post-elastic branch is smooth, and the parameters A_1 , A_2 and A_3 control the size and shape of the hysteresis loop. For the simulation a statespace transformation is effected as $X = X_1$, $\dot{X} = X_2$, $Z = X_3$ to obtain the first order Itô-stochastic differential equations as,

$$\begin{bmatrix} dX_1(t) \\ dX_2(t) \\ dX_3(t) \end{bmatrix} = \begin{bmatrix} X_2(t) \\ \frac{-1}{m}(k\lambda X_1(t) + cX_2(t) + k(1 - \lambda)X_3(t)) \\ -A_1 X_3(t) |X_2(t)| |X_3(t)|^{\bar{n}-1} - A_2 X_2(t) |X_3(t)|^{\bar{n}} + A_3 X_2(t) \end{bmatrix} dt + \begin{bmatrix} 0 \\ \frac{\sigma}{m} \\ 0 \end{bmatrix} dB(t) \quad (25)$$

One can easily simulate the system using the following Euler discretization,

$$\begin{bmatrix} X_1(n+1) \\ X_2(n+1) \\ X_3(n+1) \end{bmatrix} = \begin{bmatrix} X_1(n) \\ X_2(n) \\ X_3(n) \end{bmatrix} + \begin{bmatrix} X_2(t) \\ \frac{-1}{m}(k\lambda X_1(t) + cX_2(t) + k(1 - \lambda)X_3(t)) \\ -A_1 X_3(t) |X_2(t)| |X_3(t)|^{\bar{n}-1} - A_2 X_2(t) |X_3(t)|^{\bar{n}} + A_3 X_2(t) \end{bmatrix} \Delta t + \begin{bmatrix} 0 \\ \frac{\sigma}{m} \\ 0 \end{bmatrix} \Delta B(t) \quad (26)$$

where, the drift and diffusion terms can be identified as,

$$f(t, X_t) = \begin{bmatrix} X_2(t) \\ -\frac{1}{m}(k\lambda X_1(t) + cX_2(t) + k(1-\lambda)X_3(t)) \\ -A_1 X_3(t)|X_2(t)||X_3(t)|^{\bar{n}-1} - A_2 X_2(t)|X_3(t)|^{\bar{n}} + A_3 X_2(t) \end{bmatrix}; \quad g(t, X_t) = \begin{bmatrix} 0 \\ \frac{\sigma}{m} \\ 0 \end{bmatrix} \quad (27)$$

Then, the covariance of the diffusion matrix is obtained as,

$$\Gamma = g(t, \mathbf{X}_t)g(t, \mathbf{X}_t)^T = \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{\sigma^2}{m^2} & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (28)$$

It is straightforward to note that the key equation is, $dX_2(t) = \frac{-1}{m}(k\lambda X_1(t) + cX_2(t) + k(1-\lambda)X_3(t))dt + \frac{\sigma}{m}dB(t)$, because the first equation in Eq. (25) does not convey anything about system motion and the variable $Z(t)$ is not observable. Thus, for the identification of the system only the displacement value $X_2(t)$ is considered as input to the algorithm. The system is simulated for $T = 1$ s at a frequency of 1000Hz using the parameter values, $m = 1$, $c = 20$, $k = 100000$, $\lambda = 0.5$, $A_1 = 0.5$, $A_2 = 0.5$, $A_3 = 1$, $\bar{n} = 3$ and $\sigma_1 = 2$. Once, the displacement time history of the state $X_2(t)$ is obtained, the linear and quadratic variation of the displacement history is obtained using the formula in Eq. (34)³⁷. The identified basis functions for the drift and diffusion terms are presented in Fig. 10. In sub-figure (a) it is evident that there are more number of basis functions than the input equation. It is straightforward to note that this extra basis functions arises to take care of the non-observable hysteresis parameter $Z(t)$. One can then identify the basis functions $\Theta_f(\mathbf{X})$ and $\Theta_g(\mathbf{X})$ for the drift and diffusion, respectively, that best represent the data as:

$$\begin{aligned} \Theta_f(\mathbf{X}) &= [1 \quad X_1 \quad X_2 \quad X_1^2 \quad X_2^2 \quad X_1^3 \quad X_1^4 \quad \text{sgn}(X_1) \quad |X_1| \quad |X_2| \quad X_1|X_1|] \\ \Theta_g(\mathbf{X}) &= [1 \quad X_1 \quad |X_1|] \end{aligned} \quad (29)$$

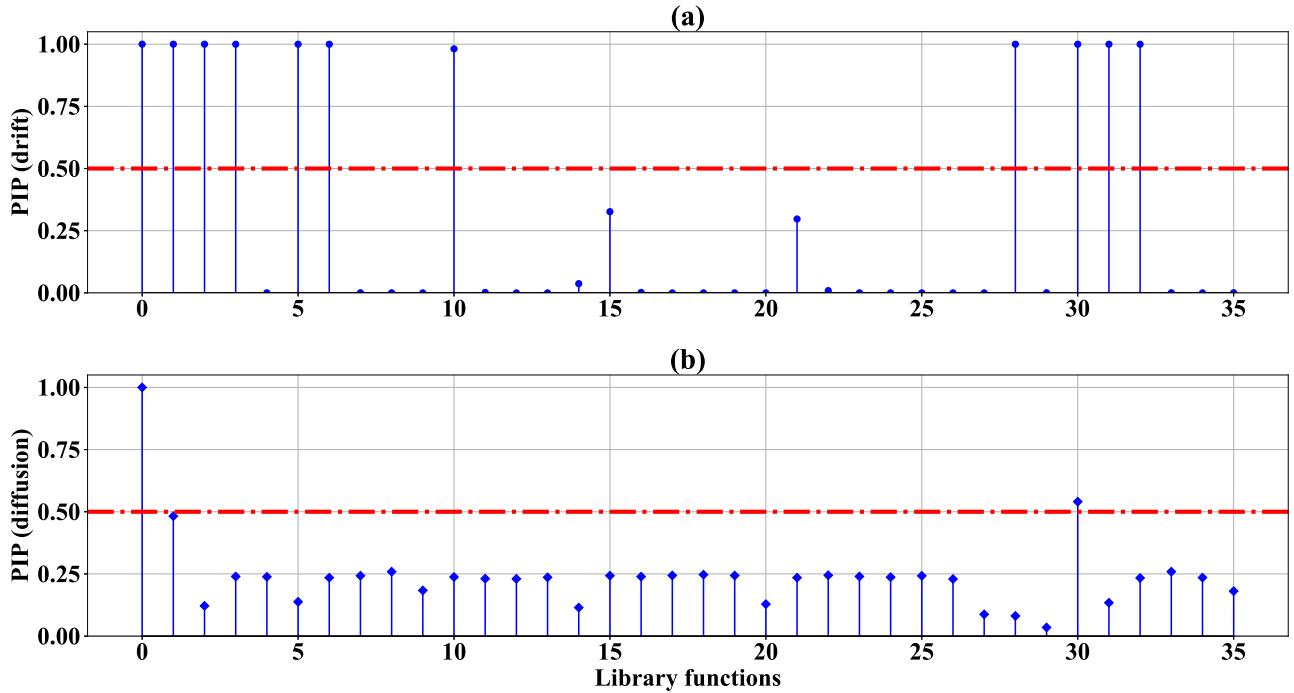


Figure 10. Bouc-Wen oscillator: Basis function selection for the Bouc-Wen base isolator system based on marginal posterior inclusion probability (PIP), $P(Z_k = 1|Y); k = 1, \dots, 36$. (a) The identified basis functions for the drift component. The black horizontal axes represent the collection of 36 basis functions, and the red line represents marginal PIP > 0.5 . The identified basis functions are $\{v(0), v(1), v(2), v(3), v(5), v(6), v(10), v(28), v(30), v(31), v(32)\}$. These bases corresponds to the functions $\{1, X_1, X_2, X_1^2, X_2^2, X_1^3, X_1^4, \text{sgn}(X_1), |X_1|, |X_2|, X_1|X_1|\}$. (b) The identified basis functions for the diffusion component. The identified function are $\{1, X_1, |X_1|\}$.

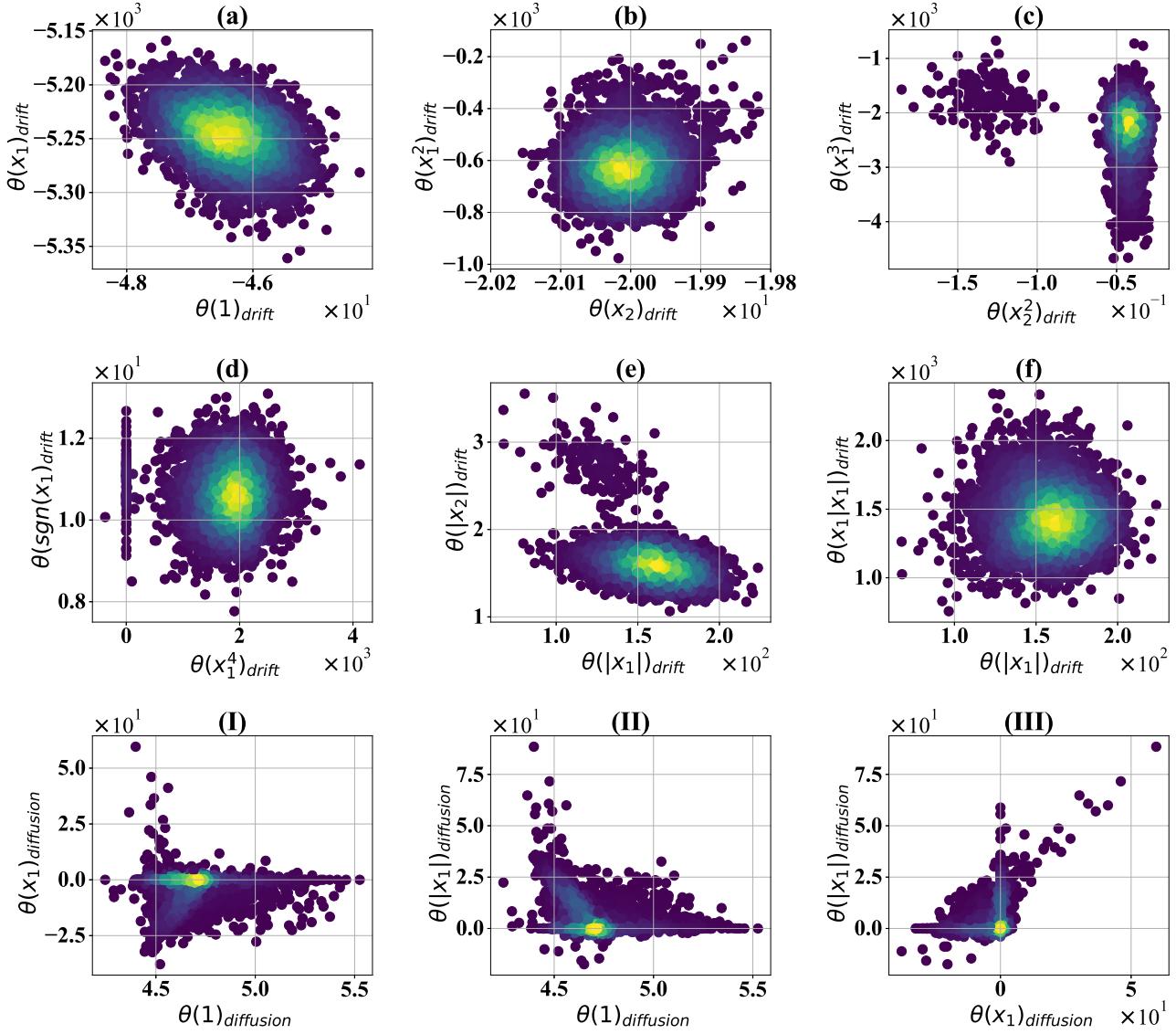


Figure 11. Bouc-Wen oscillator: Pairwise joint posterior distributions of the parameters θ_k of the selected basis functions. The subplots (a)→(f) presents the joint posterior distributions of drift parameters whereas the same is presented in the subplots (I)→(III). (a) Joint posterior distribution of $\theta(1)$ and $\theta(X_1)$. The parameters are identified as $E(\theta(1)) = -46.4$ and $E(\theta(X_1)) = -5246$. (b) Joint posterior distribution of $\theta(X_2)$ and $\theta(X_1^2)$. The expected value of the parameters are identified as $E(\theta(X_2)) = -20$ and $E(\theta(X_1^2)) = -614$. (c) Joint posterior distribution of $\theta(X_2^2)$ and $\theta(X_1^3)$. The parameters are identified as $E(\theta(X_2^2)) = -0.05$ and $E(\theta(X_1^3)) = -2328$. (d) Joint posterior distribution of $\theta(X_1^4)$ and $\theta(\text{sgn}(X_1))$. The expected value of the parameters are identified as $E(\theta(X_1^4)) = 1829$ and $E(\theta(\text{sgn}(X_1))) = 10.64$. (e) Joint posterior distribution of $\theta(|X_1|)$ and $\theta(|X_2|)$. The parameters are identified as $E(\theta(|X_1|)) = 159$ and $E(\theta(\text{sgn}(X_1))) = 10.64$. (f) Joint posterior distribution of $\theta(|X_1|)$ and $\theta(|X_1||X_1|)$. The parameters are identified as $E(\theta(|X_1|)) = 159$ and $E(\theta(|X_1||X_1|)) = 1439$. (I) The joint posterior of the weights θ_1 and θ_2 corresponding to the basis functions 1 and X_1 is shown. The yellow color denotes the mean values of the joint distribution. (II) The joint posterior of the weights θ_1 and θ_{30} corresponding to the basis functions 1 and $|X_1|$ is shown. (III) The joint posterior of the weights θ_2 and θ_{30} corresponding to the basis functions X_1 and $|X_1|$ is shown.

The joint posterior distributions of the retained basis functions in Fig. 10, (a) and (b) are plotted in the Fig. 11. Based on the basis selection in Fig. 10 and from the joint posteriors of the parameter θ_i plotted in Fig. 11, the final drift and diffusion

fields for the Bouc-Wen system can be identified as follows:

$$\begin{aligned} f(t, X) = & -46.4 - 5246X_1 - 20X_2 - 614X_1^2 - 0.05X_2^2 - 2328X_1^3 + 1829X_1^4 + 10.64\operatorname{sgn}(X_1) + 159|X_1| \\ & + 1.625|X_2| + 1439X_1|X_1| \\ g(t, X) = & 4.65 - 4.872X_1 + 6.786|X_1| \end{aligned} \quad (30)$$

Correct SDEs:	$\begin{cases} dX_2(t) = -(5000X_1(t) + 20X_2(t) + 5000X_3(t))dt + 2dB(t) \\ dX_3(t) = (-0.5X_3(t) X_2(t) X_3(t) ^2 - 0.5X_2(t) X_3(t) ^3 + X_2(t))dt \end{cases}$
Identified SDEs:	$\begin{cases} dX_2(t) = (-46.4 - 5246X_1 - 20X_2 - 614X_1^2 - 2328X_1^3 + 1829X_1^4 + 10.64\operatorname{sgn}(X_1) + 159 X_1 \\ \quad + 1.63 X_2 + 1439X_1 X_1)dt + \sqrt{4.65 - 3.79X_1 + 5.76 X_1 }dB(t) \end{cases}$

Table 3. Identification results of the Bayesian physics discovery for Bouc-Wen oscillator

For further treatment the term $0.05X_2^2$ is neglected due to its non-significant parameter value. With the above assumptions, finally the identified equation of motion of the Bouc-Wen system can be derived as in Eq. (31). In order to verify the fidelity of the identified system the learned Bouc-Wen equation in Eq. (31) is simulated using the Euler Maruyama (EM) scheme with a time step of $\Delta t = 0.001$. However, in ahead of time the external excitation will not remain same. In the present framework this is modelled as Brownian force with different intensity. A Brownian noise intensity of $\sigma = 6$ is used to simulate the new unseen environment. The results of the comparison are plotted in Fig. 12.

$$\begin{aligned} \ddot{X} + 20\dot{X} + 5246X + 614X^2 + 2328X^3 - 1829X^4 - 10.64\operatorname{sgn}(X) - 159|X| \\ - 1.625|\dot{X}| - 1439X|X| + 4.65 = 4.65 - 4.872X + 6.786|X| \end{aligned} \quad (31)$$

Figure 12 provides evidence that the proposed system correctly learns the SDEs corresponding to Bouc-Wen system with only 1 second of data (1000 samples). The learned system not only identifies the input responses of the system but also able to emulates the predicted response of the original system by learning the system from only 1 second of data. This presents the ability of the proposed scheme to not only learn a partially observed system but additionally, to be able predict the future in an unseen environment using the new basis functions. The summary of the cases studies undertaken in this work is presented in the Table 2.

Discussion

We propose a novel data-driven framework for discovery of governing physics of multi-dimensional non-linear stochastic dynamical systems from limited and noisy data. In many natural process the estimate of external input is intractable due to the limitations in the present measurement technologies. Further, measurement of all the systems are often not possible due to high operational cost. The proposed novel framework tries to emulate such situations by assuming that the input information is not required and only the measurements of noisy displacement data is available for physics discovery. The velocity state is obtained by performing the numerical differentiation on the displacement vector. The library of candidate functions are then constructed from the functional evaluated over obtained displacement and velocity states. To the knowledge of authors the discovery of governing physics of without the explicit knowledge of input excitation is not present in the literature, and, therefore a key novelty of the proposed novel framework. The absence of information about external force renders the physics discovery problem to a stochastic equation discovery problem. Towards this the proposed framework discovers the governing physics of the underlying stochastic process in terms of the stochastic differential equations (SDEs). The deterministic drift component of an SDE captures the dynamics of the underlying dynamical whereas the external stochastic force is identified in terms of the diffusion component.

The proposed framework employs the sparse Bayesian linear regression in conjunction with the Gibbs sampler to simultaneously obtain the basis functions of the model and their associated parameters. For this purpose the Kramers-Moyal expansion is utilised to express the drift and diffusion dynamics of an SDE in terms of the measured samples paths. The drift and diffusion vectors obtained from the sample paths are then used as a target vector in the sparse regression. Use of the Kramers-Moyal expansion within the purview of sparse Bayesian linear regression for discovery of an explainable governing physics is second novelty of the proposed framework. The unified framework possesses the ability of the Bayesian inference to update the probability of observing the correct basis function based on the new information. The ability to update the information is missing in recently published least-square based physics discovery schemes. This further introduces the natural elimination of the basis functions that do not present the observed data. In a noisy and unseen environment, one would be more interested in

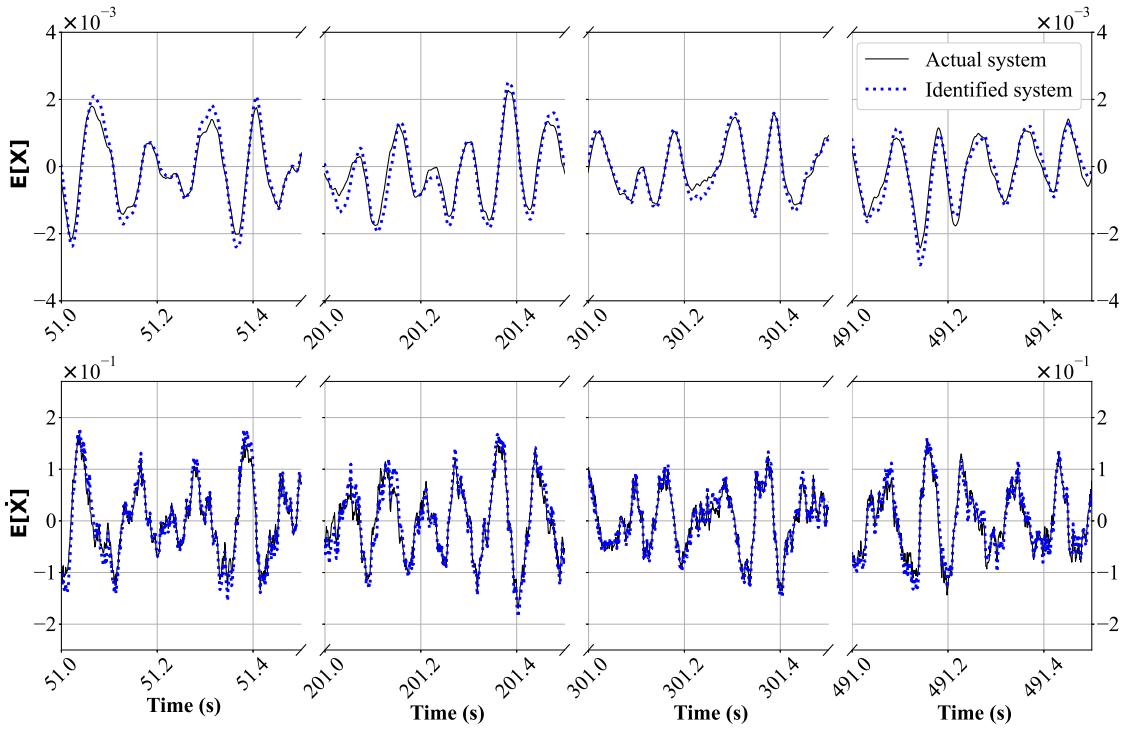


Figure 12. Bouc-Wen oscillator: Future prediction in an unseen scenario using the discovered equation in the absence of the hysteresis data. The governing equation of the Bouc-Wen system learned using 1 second of data sampled at 1000Hz. The prediction is performed for 500 seconds ahead of time using the discovered Bouc-Wen model. In the figure, the ensemble of the displacement and velocity prediction data over 500 seconds is presented. The simulation is performed using a different Brownian noise intensity of $\sigma = 6$, which is significantly higher than the training data. The dotted red line shows the prediction using the discovered equation and the solid blue line denotes the accurate simulated data. The predicted data almost accurately emulates the original system even though the measurements for all the states are not available.

dealing with the probability distribution of a random variable rather than a deterministic value. The proposed framework thus takes a leap in such situations and provides a probability distribution for each of the weights associated with corresponding basis functions. The obtained standard deviations of the distributions can then be used for defining the lower and upper bounds on the estimates of unknown parameters. Another aspect of the proposed framework is its advantages over the Neural network based grey models, where the explicit expression of the discovered governing physics is not known. In cases of partially observed processes, the proposed framework seems to provide an explainable equation for the governing physics. Although the basis functions in the obtained equation is significantly different from the one determined from first principal laws, the obtained model is able to predict the distant future without significant error. To the knowledge of authors, the discovery of physics in a partially observed process from the purview of probability theory is one of the key contributions of the proposed novel framework.

Despite of the advantages of the proposed framework over presently available physics discovery techniques there are certain issues that require special attention for accurate identification of the governing physics. The first issue is the judicious selection of the candidate basis function. The presence of basis functions with high correlation may sometimes yield a physical law different from the actual physics. This is evident from the identification example on the Black-Schole SDE. The second issue is the associated computational cost. Although the computational cost of the proposed scheme is significantly less than the time required to train its neural network alternatives, further improvements can be made. One such alternative is to use the marginal standard deviation criteria, instead of taking the mean of the latent vector. Another improvement in terms of computational efficiency can be made by devising an appropriate filter to process out the signal noise synchronously with the sparse Bayesian linear regression. The third issue is related to the quality of the data. The low-fidelity data are less expensive to obtain but the representation of the physical model in the data poor due to the presence of noise. On the other hand, the

the the high-fidelity data are highly accurate but expensive to obtain. The proposed framework in its present state is not well equipped to make use of both the low and high-fidelity data. Thus, to further enhance the fidelity of the proposed framework for field applications a more robust algorithm needs to be formulated. The use of low-fidelity data ensure low operational cost while the high-fidelity data will ensure the accuracy of the discovered physics. Lastly, the proposed framework leverages the Kramers-Moyal expansion to express the drift and diffusion components of an SDE in terms of the first and second order moments of the sample paths. The moments are approximated in a limiting sense by making use of the absolute and quadratic identities of Brownian motion. In this context, future improvements in the proposed framework can be made by exploiting the higher order moments to retain the higher order terms from Stochastic-Taylor expansions of the drift and diffusion terms. Overall, the proposed framework presented a novel methodology for discovering the governing physics of a stochastic process whenever there is not possible to measure the input force data. The learned equations have shown accurate identification of the physical process and demonstrated good prediction ability in distant future.

Methods

The structure of the proposed Bayesian physics discovery framework is presented in the Fig. 2. The proposed framework treats the state measurements from the sensors as stochastic process and applies Kramers-Moyal formula to obtain the target vector for sparse linear regression. In the sparse regression, the library is obtained by taking ensemble mean over all the libraries constructed from the collected data set. Then the proposed framework constructs two sparse regression problem using these libraries and target vectors for identification of the governing physics, which are independent and can be solved simultaneously. The discovery of governing physics here involves identification of the system dynamics and the volatility associated arising due to environmental uncertainty. The system dynamics and the volatility are represented in this framework as drift and diffusion components of an SDE. The methodology of the proposed framework is further illustrated though mathematical equations as follows. Let $\mathbf{L} \in \mathbb{R}^{N \times K}$ be the library of candidate functions, where, K is the total number of independent basis functions in the library and N is the sample length. Then the key idea of the sparse Bayesian linear regression is to select k basis functions from the total number of K candidate library functions such that $k \ll K$. If $\mathbf{Y} \in \mathbb{R}^N$ is the target vector then the linear combinations of the identified basis functions should represent the target vector \mathbf{Y} in best possible manner. The coefficients of the linear coefficient is denoted by a weight vector $\boldsymbol{\theta}$. For the Bayesian model discovery, the condition $k \ll K$ is maintained by assigning certain type of prior distributions on the weight vector $\boldsymbol{\theta}$, such that they favour the sparsity in the solution. In this context, the SS-prior have shown high shrinkage property due to its sharp spike at zero and a diffused density spanned over a large range of possible parameter value. The Dirac-delta spike concentrates most of the probability mass at zero thus allowing most of the samples to take value zero. On the other hand the diffused tail distributes a small amount of probability mass over a large range of possible values allowing only very few samples with very high probability to escape the shrinkage. The alternate flavours of the SS-prior constructed from various combination of candidate spike and slab functions are well documented in the literature^{32,33,38}. In this work, the authors have considered the discontinuous spike and slab prior (DSS-prior) where the spike at zero is modelled as Dirac-delta function and the tail distribution as independent Student's-t distribution. Next, a sparse regression framework is formulated for stochastic differential equations.

Sparse learning of Stochastic differential equations

The natural dynamical systems are often not observed directly but through some projected space. The projected space contains all the states of a system which are directly observable. For instance, the second order dynamical systems are often expressed in terms of its displacement and velocity components. In cases of non-linear stochastic dynamical systems the stochastic differential equations provide an splendid approach for expressing the behavior of the underlying dynamical system. Stochastic differential equations (SDEs) arises naturally in non-linear dynamical systems subjected to stochastic excitation such as earthquake, wind force, wave force etc.^{39,40}. An SDE is defined in term of its deterministic drift and the additive stochastic diffusion components where both the components are allowed to depend on time and systems states. Let (Ω, \mathcal{F}, P) be the probability space and $\{\mathcal{F}_t, 0 \leq t \leq T\}$ be the natural filtration constructed from sub σ -algebras of \mathcal{F} . Consider the m -dimensional n -factor SDE driven by n -dimensional Brownian motion $\{\mathbf{B}_j(t), j = 1, \dots, n\}$:

$$dX_t = f(X_t, t) dt + \sum_{j=1}^n g_j(X_t, t) dB_j(t); \quad X(t=t_0) = X_0; \quad t \in [0, T] \quad (32)$$

where, $X_t \in \mathbb{R}^m$ denotes the \mathcal{F}_t -measurable state vector, $f(X_t, t) : \mathbb{R}^m \mapsto \mathbb{R}^m$ is the drift vector, $g(X_t, t) : \mathbb{R}^m \mapsto \mathbb{R}^{m \times n}$ is the diffusion matrix and $\mathbf{B}_j(t) \in \mathbb{R}^n$ is the Brownian motion. In a compact matrix notation the SDE can be expressed as:

$$dX_t = f(X_t, t) dt + g(X_t, t) dB(t); \quad X(t=t_0) = X_0; \quad t \in [0, T] \quad (33)$$

The solution to Eq. (33) can be attempted using various stochastic integration schemes^{35,41–43}. However, in the interest of the present work one can try to find the probability distribution function (PDF) of the random variable $X(t)$. In general the

stochastic Brownian motion $B(t)$ is not a well defined mathematical function since it is not differential everywhere with respect to the process $X(t)$. This kind of non-differentiable functions have zero finite variation but non-vanishing quadratic variation. The Kramers-Moyal formula suggests that the drift and diffusion components of an SDE can be expressed in terms of the linear and quadratic variations of the sample time history as (the simplified detailed derivation is presented in A):

$$\begin{aligned} f_i(X_t, t) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[X_i(t + \Delta t) - \xi_i] \Big|_{x_k(t) = \xi_k} \quad \forall k = 1, 2, \dots, N \\ \Gamma_{ij}(X_t, t) &= \frac{1}{2} \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[|X_i(t + \Delta t) - \xi_i| |X_j(t + \Delta t) - \xi_j|] \Big|_{x_k(t) = \xi_k} \quad \forall k = 1, 2, \dots, N \end{aligned} \quad (34)$$

where, $f_i(X_t, t)$ is the i^{th} drift component and Γ_{ij} is the $(ij)^{th}$ component of the diffusion covaraince matrix $\Gamma \in \mathbb{R}^{n \times n} := g(t, X_t)g(t, X_t)^T$. The discovery of an SDE requires identification of the drift and diffusion components in terms of the candidate basis functions from the library. Since the discovery problem demands the identified physics to have a interpretable form, sparse regression using Bayesian framework can be utilized in this context. Thus for the discovery of the drift and diffusion terms two independent linear regression problems are constructed next.

For the discovery of the i^{th} drift terms, the target vector is defined as, $\mathbf{Y}_i := \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[X_i(t + \Delta t) - \xi_i] \Big|_{x_k(t) = \xi_k}$. Further, in a m -dimensional diffusion process each of the drift term can be expressed in terms of the basis functions as, $f_i(X_t, t) = \sum_{k=1}^K v_k(X_t) \theta_{i,k}; i = 1 \dots m$ where, $\{v_k(X_t); k = 1, 2, \dots, K\}$ denotes the candidate basis functions and $\{\theta_{i,k}; k = 1 \dots K\}$ denotes the weight vector corresponding to the i^{th} drift component. A large library is constructed from the candidate basis functions as, $\mathbf{L} = \{v_1(X_t), \dots, v_K(X_t)\}$, where, $v_k(X_t)$ represents the various linear and non-linear mathematical functions evaluated on the system states. In a general setting, the identification of the drift terms culminates in the following regression problem.

$$[\mathbf{Y}_1 \quad \mathbf{Y}_2 \quad \dots \quad \mathbf{Y}_m] = \underbrace{\begin{bmatrix} v_1(X_{1,1} \dots X_{m,1}) & v_2(X_{1,1} \dots X_{m,1}) & \dots & v_K(X_{1,1} \dots X_{m,1}) \\ v_1(X_{1,2} \dots X_{m,2}) & v_2(X_{1,2} \dots X_{m,2}) & \dots & v_K(X_{1,2} \dots X_{m,2}) \\ \vdots & \vdots & \ddots & \vdots \\ v_1(X_{1,N} \dots X_{m,N}) & v_2(X_{1,N} \dots X_{m,N}) & \dots & v_K(X_{1,N} \dots X_{m,N}) \end{bmatrix}}_{\mathbf{L} \in \mathbb{R}^{N \times K}} \underbrace{\begin{bmatrix} \theta_{1,1} & \theta_{2,1} & \dots & \theta_{m,1} \\ \theta_{1,2} & \theta_{2,2} & \dots & \theta_{m,2} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{1,k} & \theta_{2,k} & \dots & \theta_{m,k} \end{bmatrix}}_{\theta \in \mathbb{R}^{K \times m}} \quad (35)$$

As a subset of the above regression one can solve the following equation and identify the drift components independently:

$$\mathbf{Y}_i = \mathbf{L}\boldsymbol{\theta}_i^f + \boldsymbol{\varepsilon}_i \quad (36)$$

where, the i^{th} target vector \mathbf{Y}_i is constructed from the linear variations of the i^{th} sample path as: $\mathbf{Y}_i = [(X_{i,1} - \xi_{i,1}) \dots (X_{i,N} - \xi_{i,N})]^T$ and the drift weight vector is $\boldsymbol{\theta}_i^f = [\theta_{1,1} \quad \theta_{1,2} \quad \dots \quad \theta_{1,K}]^T$. The straightforward application of the Gibbs sampling mentioned in the Eqs. (47) → (52) can be performed to discover the drift components of a system of SDEs.

The diffusion components does not have finite variations but are bounded by their quadratic variations. Thus the diffusion components unlike the drifts of an SDE are recoverable only through its covariation terms. In the formulation of the sparse regression problem for diffusion components the target vector is defined as, $\mathbf{Y}_{ij} := \frac{1}{2} \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[|X_i(t + \Delta t) - \xi_i| |X_j(t + \Delta t) - \xi_j|] \Big|_{x_k(t) = \xi_k}$. Also, the $(ij)^{th}$ component of the covariance matrix Γ can be expressed in terms of the basis functions as: $\Gamma_{ij} = \sum_{k=1}^K v_k(X_t) \theta_{ijk}; i, j = 1, 2, \dots, m$. Noting that the covariance matrix of the $m \times n$ -diffusion matrix has a dimension $m \times m$ the following regression problem can be constructed,

$$\mathbf{Y}_{ij} = \mathbf{L}\boldsymbol{\theta}_{ij}^g + \boldsymbol{\eta}_{ij} \quad (37)$$

where, $\mathbf{Y}_{ij} = \Gamma_{ij} = [(X_{i,1} - \xi_{i,1})(X_{j,1} - \xi_{j,1}) \dots (X_{i,N} - \xi_{i,N})(X_{j,N} - \xi_{j,N})]^T$ is the quadratic covariation of the i^{th} and j^{th} sample path, and $\boldsymbol{\theta}_{ij}^g = [\theta_1^{ij} \quad \theta_2^{ij} \quad \dots \quad \theta_K^{ij}]^T$ is the weight vector for the diffusion terms. With the above regression statement the Gibbs sampling can be performed to obtain the sparse solution. Due to the symmetry of the covariance matrix the above sparse regression is required to be performed for $m(m+1)/2$ times to completely identify the diffusion space.

Discovery of SDE by sparse Bayesian regression

The equation discovery using Bayesian linear regression essentially involves classification of the weights of the basis functions in the library into either spike or slab components of the SS-priors. This is done by introducing a latent indicator variable $\mathbf{Z} = [Z_1, \dots, Z_K]$ for each of the component $\theta_k; k = 1, \dots, K$ of the weight vector $\boldsymbol{\theta}$. The latent indicator variable Z_k is assigned

a value 1 if the weight corresponds to the slab component, otherwise, it takes a value 0 when the weight belongs to spike component. For the understanding on the sparse Bayesian linear regression, consider the following one-dimensional regression,

$$\mathbf{Y} = \mathbf{L}\boldsymbol{\theta} + \epsilon \quad (38)$$

where, $\mathbf{Y} \in \mathbb{R}^N$ denotes the N -dimensional target vector, \mathbf{L} denotes the library of candidate functions, $\boldsymbol{\theta}$ is the weight vector, and $\epsilon \in \mathbb{R}^N$ is the residual error vector representing the model mismatch error. The model error ϵ is modelled as i.i.d Gaussian random variable with zero mean and variance σ^2 . As the name suggests only a few function from the candidate library will actively participate in the final representation of the target vector \mathbf{Y} . This will yield a sparse solution such that only the components of weight vector that consistently takes a value sufficiently different from zero will be considered. For estimating the weight vector $\boldsymbol{\theta}$, the Bayes formula can be applied as,

$$P(\boldsymbol{\theta}|\mathbf{Y}) = \frac{P(\mathbf{Y}|\boldsymbol{\theta})}{P(\mathbf{Y})} \quad (39)$$

However, the effect of the parameters of the DSS-prior on the probability of the weights also needs to accounted. For this one can construct a Bayesian hierarchical model as presented in Fig. 1. Then, the likelihood function is written as,

$$\mathbf{Y}|\boldsymbol{\theta}, \sigma^2 \sim \mathcal{N}(\mathbf{L}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_{N \times N}) \quad (40)$$

where, $\mathbf{I}_{N \times N}$ denotes the $N \times N$ identity matrix. For the estimate of the posterior distribution of the parameters of the hierarchical the components of the weight vector that belongs to spike does not contribute. For this purpose let us consider a vector $\boldsymbol{\theta}_r \in \mathbb{R}^r : \{r \ll K\}$, composed from the elements of the weight vector $\boldsymbol{\theta}$ for which $Z_k = 1$. Thus, if $\boldsymbol{\theta}_r$ denotes the weight vector corresponding to slab component of DSS-prior then the SS-prior is defined as^{31,32},

$$p(\boldsymbol{\theta}|\mathbf{Z}) = p_{slab}(\boldsymbol{\theta}_r) \prod_{k, Z_k=0} p_{spike}(\theta_k) \quad (41)$$

where, p_{spike} and p_{slab} denotes the spike and slab distributions. These are defined as, $p_{spike}(\theta_k) = \delta_0$ and $p_{slab}(\boldsymbol{\theta}_r) = \mathcal{N}(\mathbf{0}, \sigma^2 \vartheta_s \mathbf{R}_{0,r})$ with $\mathbf{R}_{0,r} = \mathbf{I}_{r \times r}$. It is to be noted that the spike and slab distributions are assumed as independent. Given the noise variance σ^2 , the the slab variance ϑ_s is assigned the Inverse-gamma prior with the hyperparameters α_ϑ and β_ϑ . The latent variables Z_k takes a value from the set $\{0, 1\}$, thus each of the element in the latent vector is assigned the Bernoulli prior with common hyperparameter p_0 . However, the hyperparameter p_0 is allowed to adapt according to the Beta prior with the hyperparameters α_p , and β_p . Similarly the variance of the measurement noise is simulated from the Inverse-gamma distribution with the hyperparameters α_σ , and β_σ . The hyperparameters α_ϑ , β_ϑ , α_p , β_p , α_σ , and β_σ are provided as a deterministic constants in the hierarchical model. The summery of the DSS-prior model is specified in the Fig. 1 whose corresponding equations are given in Eqs. (42), (43), (44) and (45).

$$p(\vartheta_s) = IG(\alpha_\vartheta, \beta_\vartheta) \quad (42)$$

$$p(Z_k|p_0) = Bern(p_0); k = 1 \dots K \quad (43)$$

$$p(p_0) = Beta(\alpha_p, \beta_p) \quad (44)$$

$$p(\sigma^2) = IG(\alpha_\sigma, \beta_\sigma) \quad (45)$$

From the DAG structure of the DSS-model in Fig. 1, the joint distribution of the random variables $\boldsymbol{\theta}, \mathbf{Z}, \vartheta_s, \sigma^2, p_0 | \mathbf{Y}$ can be expanded using the Bayes formula as,

$$\begin{aligned} p(\boldsymbol{\theta}, \mathbf{Z}, \vartheta_s, \sigma^2, p_0 | \mathbf{Y}) &= \frac{p(\mathbf{Y}|\boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta}|\mathbf{Z}, \vartheta_s, \sigma^2) p(\mathbf{Z}|p_0) p(\vartheta_s) p(\sigma^2) p(p_0)}{p(\mathbf{Y})} \\ &\propto p(\mathbf{Y}|\boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta}|\mathbf{Z}, \vartheta_s, \sigma^2) p(\mathbf{Z}|p_0) p(\vartheta_s) p(\sigma^2) p(p_0) \end{aligned} \quad (46)$$

where, $p(\boldsymbol{\theta}, \mathbf{Z}, \vartheta_s, \sigma^2, p_0 | \mathbf{Y})$ denotes the joint distribution of the random variables, $p(\mathbf{Y}|\boldsymbol{\theta}, \sigma^2)$ denotes the likelihood function, $p(\boldsymbol{\theta}|\mathbf{Z}, \vartheta_s, \sigma^2)$ is the prior distribution for the weight vector $\boldsymbol{\theta}$, $p(\mathbf{Z}|p_0)$ is the prior distribution for the latent vector \mathbf{Z} , $p(\vartheta_s)$ is the prior distribution for the slab variance ϑ_s , $p(\sigma^2)$ is the prior distribution for the noise variance, $p(p_0)$ is the prior distribution for the success probability p_0 and $p(\mathbf{Y})$ is the marginal likelihood. Direct sampling from the joint distribution function is intractable in this case due to the spike and slab distribution. Thus, Gibbs sampling technique is used to draw the random

samples from the joint distribution⁴⁴. For the Gibbs sampling the conditional distributions of the random variables are derived in Appendix B. The weights corresponding to the spike distribution does not contribute to the selection of the correct basis functions. Thus only the weights corresponding to $Z_k \neq 0$ i.e. the θ_r vector is samples using the Gibbs sampling. Referring to the Eqs. (42), (43), (44) and (45), the sequence of the random variables $\theta^{(0)}, \sigma^{2(0)}, \vartheta_s^{(0)}, p_0^{(0)}, \mathbf{Z}^{(0)}, \dots, \theta^{(1)}, \sigma^{2(1)}, \vartheta_s^{(1)}, p_0^{(1)}, \mathbf{Z}^{(1)}, \dots$, using the Gibbs sampling technique is obtained by following steps,

1. The weight vector $\theta_r^{(i)}$ is sampled from the Gaussian distribution with mean μ_θ and variance Σ_θ as,

$$\theta_r^{(i)} | \mathbf{Y}, \vartheta_s^{(i)}, \sigma^{2(i)} \sim N(\mu_\theta^{(i)}, \Sigma_\theta^{(i)}) \quad (47)$$

where, the mean and covariance is defined as, $\mu_\theta^{(i)} = \Sigma_\theta^{(i)} \mathbf{L}_r^{(i)T} \mathbf{Y}$ and $\Sigma_\theta^{(i)} = \sigma^{2(i)} (\mathbf{L}_r^{(i)T} \mathbf{L}_r^{(i)} + \vartheta_s^{(i)-1} \mathbf{R}_{0,r}^{(i)-1})^{-1}$, respectively.

2. The latent variable $Z_k^{(i+1)}$ is assigned the values from the set, {0, 1} by using the Bernoulli distribution as,

$$Z_k^{(i+1)} | \mathbf{Y}, \vartheta_s^{(i)}, p_0^{(i)} \sim Bern(u_k) \quad (48)$$

where, $u_k = \frac{p_0}{p_0 + \lambda(1-p_0)}$ and $\lambda = \frac{p(\mathbf{Y}|Z_k^{(i)}=0, \mathbf{Z}_{-k}^{(i)}, \vartheta_s^{(i)})}{p(\mathbf{Y}|Z_k^{(i)}=1, \mathbf{Z}_{-k}^{(i)}, \vartheta_s^{(i)})}$. Here, $\mathbf{Z}_{-k}^{(i)} \in \mathbb{R}^{K-1}$ denotes the latent variable vector \mathbf{Z} consisting of all the elements except the k^{th} component. The probability that the k^{th} latent variable $Z_k^{(i)}$ takes a value 0 or 1 is estimated as follows,

$$\begin{aligned} p(\mathbf{Y}|\mathbf{Z}^{(i)}, \vartheta_s^{(i)}) &= \frac{\Gamma\left(\alpha_\sigma + \frac{N}{2}\right) \beta_\sigma^{\alpha_\sigma}}{\Gamma(\alpha_\sigma)(2\pi)^{\frac{N}{2}} \left(\beta_\sigma + \frac{1}{2} \mathbf{Y}^T \mathbf{Y}\right)^{\frac{N}{2}}}; \quad \text{when all } \{Z_k^{(i)} : k = 1, \dots, K\} = 0 \\ &= \frac{\Gamma\left(\alpha_\sigma + \frac{N}{2}\right) \beta_\sigma^{\alpha_\sigma} \left(\left|\mathbf{R}_{0,r}^{(i)-1}\right| \left|\Sigma_\theta^{(i)}\right|\right)^{\frac{1}{2}}}{\Gamma(\alpha_\sigma)(2\pi)^{\frac{N}{2}} \vartheta_s^{\frac{N}{2}} \left(\beta_\sigma + \frac{1}{2} \mathbf{Y}^T (\mathbf{I}_{N \times N} - \mathbf{L}_r^{(i)T} \Sigma_\theta^{(i)} \mathbf{L}_r^{(i)}) \mathbf{Y}\right)^{\frac{N}{2}}}; \quad \text{otherwise} \end{aligned} \quad (49)$$

3. The noise variance $\sigma^{2(i+1)}$ is simulated from the Inverse-gamma distribution as,

$$\sigma^{2(i+1)} | \mathbf{Y}, \mathbf{Z}^{(i+1)}, \vartheta_s^{(i)} \sim IG\left(\alpha_\sigma + \frac{N}{2}, \beta_\sigma + \frac{1}{2} (\mathbf{Y}^T \mathbf{Y} - \mu_\theta^{(i)T} \Sigma_\theta^{(i)-1} \mu_\theta^{(i)})\right) \quad (50)$$

4. The slab variance $\vartheta_s^{(i+1)}$ is sampled from the Inverse-gamma distribution as,

$$\vartheta_s^{(i+1)} | \theta_r^{(i)}, \mathbf{Z}^{(i+1)}, \sigma^{2(i+1)} \sim IG\left(\alpha_\vartheta + \frac{h_z}{2}, \beta_\vartheta + \frac{1}{2\sigma^2} \theta_r^{(i)T} \mathbf{R}_{0,r}^{(i)-1} \theta_r^{(i)}\right) \quad (51)$$

5. The success rate $p_0^{(i+1)}$ is sample from the Beta distribution as,

$$p_0^{(i+1)} | \mathbf{Z}^{(i+1)} \sim Beta(\alpha_p + h_z, \beta_p + K - h_z) \quad (52)$$

where, $h_z = \sum_{k=1}^K Z_k^{(i+1)}$.

6. The weight vector $\theta_r^{(i+1)}$ is updated using the step 1.

This MCMC is performed for a total of N_M simulations out of which initial 1000 samples are discarded as the burn-in samples. Let N_s denote the number of MCMC required to achieve the stationary distribution after the burn-in samples are discarded. Then the marginal posterior inclusion probability (PIP):= $p(Z_k = 1 | \mathbf{Y})$ for each of the K basis functions can be estimated by taking mean over the Gibbs samples for each of the k^{th} latent vector Z_k ³¹ as,

$$p(Z_k = 1 | \mathbf{Y}) \approx \frac{1}{N_s} \sum_{j=1}^{N_s} Z_k^j; k = 1, \dots, K \quad (53)$$

The basis functions whose corresponding PIP values are more than 0.5 i.e. $p(Z_k = 1|Y) > 0.5$ are included in the final model of discovered equation. The PIP value greater than 0.5 indicates that the selected basis functions are observed more than half of the times in the MCMC simulations. Higher PIP value suggests that in case of unseen scenarios the corresponding basis functions are highly likely to occur in the data representation target vector. The mean and covariance of the weight vector gives the expected value and standard deviation of the system parameters. The standard deviation provides the confidence interval of the parameters which can be used to design the lower and upper bounds of a random variable in an unseen environment. A pseudo code for the proposed framework, is provided in the Algorithm 1.

Algorithm 1 Pseudo code of the proposed Bayesian framework for discovery of governing physics of stochastic non-linear systems

Input: Sample paths: $\mathbf{X}(t) \in \mathbb{R}^{N \times m}$, hyperparameters: $\alpha_p, \beta_p, \alpha_\sigma, \beta_\sigma, \alpha_\theta, \beta_\theta, p_0^{(0)}, \vartheta_s^{(0)}$

- 1: Estimate the linear and quadratic variation vectors. ▷ Eq. (34)
- 2: For drift obtain the target vectors \mathbf{Y}_i and the library \mathbf{L} using the candidate basis functions. ▷ Eq. (36)
- 3: For diffusion obtain the target vectors \mathbf{Y}_{ij} and the library \mathbf{L} . ▷ Eq. (37)
- 4: Estimate the initial variance of noise from residual variance: $\sigma^{2,(0)} = \text{Var}(\mathbf{L}\boldsymbol{\theta} - \mathbf{Y})$
- 5: Estimate the initial latent vector $\mathbf{Z}^{(0)} = [Z_1^{(0)}, Z_2^{(0)}, \dots, Z_K^{(0)}]$ such that $\text{MSE}(\mathbf{L}\boldsymbol{\theta} - \mathbf{Y})$ is minimum.
- 6: Find $\boldsymbol{\theta}^{(0)}, \mu_\theta$ and Σ_θ ▷ Eq. (47)
- 7: **for** $i = 1, \dots, N_M$ **do**
- 8: Update the latent variable vector $\mathbf{Z}^{(i)}$. ▷ Eq. (48)
- 9: Update the noise variance $\sigma^{2(i)}$. ▷ Eq. (50)
- 10: Update the slab variance $\vartheta_s^{(i)}$. ▷ Eq. (51)
- 11: Update the success rate $p_0^{(i)}$. ▷ Eq. (52)
- 12: Update the weight vector $\boldsymbol{\theta}^{(i)}$. ▷ Eq. (47)
- 13: Repeat steps 8→12
- 14: Discard the burn-in MCMC samples
- 15: Estimate the marginal PIP values $p(Z_k = 1|Y)$ ▷ Eq. (53)
- 16: Include the basis functions with higher PIP values
- 17: Estimate the expected value $E[\theta_k]$ and standard deviation $\sigma[\theta_k]$ of the system parameters

Output: $\{\theta_k; k = 1 \dots K\}, \{\vartheta_k(X_j); j = 1 \dots \bar{m}\}$, where, K and \bar{m} are the number of library functions and process states, respectively.

Data availability

On acceptance all the used datasets in this study will be made public on GitHub by the corresponding author.

Code availability

On acceptance all the source codes to reproduce the results in this study will be made available to public on GitHub by the corresponding author.

References

1. Zanna, L. & Bolton, T. Data-driven equation discovery of ocean mesoscale closures. *Geophys. Res. Lett.* **47**, e2020GL088376 (2020).
2. Zheng, Y. *et al.* The maximum likelihood climate change for global warming under the influence of greenhouse effect and lévy noise. *Chaos: An Interdiscip. J. Nonlinear Sci.* **30**, 013132 (2020).
3. Jia, C., Zhang, M. Q. & Qian, H. Emergent lévy behavior in single-cell stochastic gene expression. *Phys. Rev. E* **96**, 040402 (2017).
4. Noé, F. & Clementi, C. Collective variables for the study of long-time kinetics from molecular trajectories: theory and methods. *Curr. opinion structural biology* **43**, 141–147 (2017).
5. Zhang, X. F. Information uncertainty and stock returns. *The J. Finance* **61**, 105–137 (2006).
6. Bongard, J. & Lipson, H. Automated reverse engineering of nonlinear dynamical systems. *Proc. Natl. Acad. Sci.* **104**, 9943–9948 (2007).

7. Brunton, S. L., Proctor, J. L. & Kutz, J. N. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proc. national academy sciences* **113**, 3932–3937 (2016).
8. Akaike, H. A new look at the statistical model identification. *IEEE transactions on automatic control* **19**, 716–723 (1974).
9. Schwarz, G. Estimating the dimension of a model. *The annals statistics* 461–464 (1978).
10. Schmidt, M. & Lipson, H. Distilling free-form natural laws from experimental data. *science* **324**, 81–85 (2009).
11. Kevrekidis, I. G. & Samaey, G. Equation-free multiscale computation: Algorithms and applications. *Annu. review physical chemistry* **60**, 321–344 (2009).
12. Mangan, N. M., Brunton, S. L., Proctor, J. L. & Kutz, J. N. Inferring biological networks by sparse identification of nonlinear dynamics. *IEEE Transactions on Mol. Biol. Multi-Scale Commun.* **2**, 52–63 (2016).
13. Hoffmann, M., Fröhner, C. & Noé, F. Reactive sindy: Discovering governing reactions from concentration data. *The J. chemical physics* **150**, 025101 (2019).
14. Bhadriraju, B., Narasingam, A. & Kwon, J. S.-I. Machine learning-based adaptive model identification of systems: Application to a chemical process. *Chem. Eng. Res. Des.* **152**, 372–383 (2019).
15. Loiseau, J.-C. & Brunton, S. L. Constrained sparse galerkin regression. *J. Fluid Mech.* **838**, 42–67 (2018).
16. Loiseau, J.-C., Noack, B. R. & Brunton, S. L. Sparse reduced-order modelling: sensor-based dynamics to full-state estimation. *J. Fluid Mech.* **844**, 459–490 (2018).
17. Lai, Z. & Nagarajaiah, S. Sparse structural system identification method for nonlinear dynamic systems with hysteresis/inelastic behavior. *Mech. Syst. Signal Process.* **117**, 813–842 (2019).
18. Li, S. *et al.* Discovering time-varying aerodynamics of a prototype bridge by sparse identification of nonlinear dynamical systems. *Phys. Rev. E* **100**, 022220 (2019).
19. Schaeffer, H. & McCalla, S. G. Sparse model selection via integral terms. *Phys. Rev. E* **96**, 023302 (2017).
20. Mangan, N. M., Kutz, J. N., Brunton, S. L. & Proctor, J. L. Model selection for dynamical systems via sparse regression and information criteria. *Proc. Royal Soc. A: Math. Phys. Eng. Sci.* **473**, 20170009 (2017).
21. Kaiser, E., Kutz, J. N. & Brunton, S. L. Sparse identification of nonlinear dynamics for model predictive control in the low-data limit. *Proc. Royal Soc. A* **474**, 20180335 (2018).
22. Schaeffer, H., Tran, G., Ward, R. & Zhang, L. Extracting structured dynamical systems using sparse optimization with very few samples. *Multiscale Model. & Simul.* **18**, 1435–1461 (2020).
23. Stender, M., Oberst, S. & Hoffmann, N. Recovery of differential equations from impulse response time series data for model identification and feature extraction. *Vibration* **2**, 25–46 (2019).
24. Boninsegna, L., Nüske, F. & Clementi, C. Sparse learning of stochastic dynamical equations. *The J. chemical physics* **148**, 241723 (2018).
25. Rudy, S. H., Brunton, S. L., Proctor, J. L. & Kutz, J. N. Data-driven discovery of partial differential equations. *Sci. Adv.* **3**, e1602614 (2017).
26. Zhang, S. & Lin, G. Robust data-driven discovery of governing physical laws with error bars. *Proc. Royal Soc. A: Math. Phys. Eng. Sci.* **474**, 20180305 (2018).
27. Rudy, S. H., Kutz, J. N. & Brunton, S. L. Deep learning of dynamics and signal-noise decomposition with time-stepping constraints. *J. Comput. Phys.* **396**, 483–506 (2019).
28. Raissi, M. Deep hidden physics models: Deep learning of nonlinear partial differential equations. *The J. Mach. Learn. Res.* **19**, 932–955 (2018).
29. Raissi, M., Perdikaris, P. & Karniadakis, G. E. Multistep neural networks for data-driven discovery of nonlinear dynamical systems. *arXiv preprint arXiv:1801.01236* (2018).
30. Fuentes, R. *et al.* Equation discovery for nonlinear dynamical systems: A bayesian viewpoint. *Mech. Syst. Signal Process.* **154**, 107528 (2021).
31. Nayek, R., Fuentes, R., Worden, K. & Cross, E. J. On spike-and-slab priors for bayesian equation discovery of nonlinear dynamical systems via sparse linear regression. *Mech. Syst. Signal Process.* **161**, 107986 (2021).
32. Mitchell, T. J. & Beauchamp, J. J. Bayesian variable selection in linear regression. *J. american statistical association* **83**, 1023–1032 (1988).

33. George, E. I. & McCulloch, R. E. Approaches for bayesian variable selection. *Stat. sinica* 339–373 (1997).
34. Oksendal, B. *Stochastic differential equations: an introduction with applications* (Springer Science & Business Media, 2013).
35. Tripura, T., Gogoi, A. & Hazra, B. An ito-taylor weak 3.0 method for stochastic dynamics of nonlinear systems. *Appl. Math. Model.* (2020).
36. Tripura, T., Bhowmik, B., Pakrashi, V. & Hazra, B. Real-time damage detection of degrading systems. *Struct. Heal. Monit.* **19**, 810–837 (2020).
37. Risken, H. Fokker-planck equation. In *The Fokker-Planck Equation*, 63–95 (Springer, 1996).
38. O’Hara, R. B. & Sillanpää, M. J. A review of bayesian variable selection methods: what, how and which. *Bayesian analysis* **4**, 85–117 (2009).
39. Calin, O. *An informal introduction to stochastic calculus with applications* (World Scientific, 2015).
40. Klebaner, F. C. *Introduction to stochastic calculus with applications* (World Scientific Publishing Company, 2005).
41. Kloeden, P. E. & Platen, E. Higher-order implicit strong numerical schemes for stochastic differential equations. *J. statistical physics* **66**, 283–314 (1992).
42. Tripura, T., Imran, M., Hazra, B. & Chakraborty, S. A change of measure enhanced near exact euler maruyama scheme for the solution to nonlinear stochastic dynamical systems. *arXiv preprint arXiv:2108.10655* (2021).
43. Tripura, T., Hazra, B. & Chakraborty, S. Generalized weakly corrected milstein solutions to stochastic differential equations. *arXiv preprint arXiv:2108.10681* (2021).
44. Casella, G. & George, E. I. Explaining the gibbs sampler. *The Am. Stat.* **46**, 167–174 (1992).

Acknowledgements

SC acknowledges the financial support received from IIT Delhi in form of seed grant.

Author contributions statement

T. Tripura: Conceptualization, Methodology, Software, Formal analysis, Writing - original draft. S. Chakraborty: Conceptualization, Methodology, Software, review and editing, Supervision, Funding acquisition.

Competing interests

The authors declare no competing interests. The corresponding author is responsible for submitting a [competing interests statement](#) on behalf of all authors of the paper. This statement must be included in the submitted article file.

Additional information

Correspondence and requests for materials should be addressed to Souvik Chakraborty.

A Kramers-Moyal expansion for estimation of the drift and diffusion terms of an SDE from sample paths

Let us consider, $p(X, t) = P(X, t|X_0, t_0)$ be the transition probability density of the solution of the SDE in Eq. (33). Then the Kramers-Moyal expansion is written as³⁷:

$$\frac{\partial P(X, t)}{\partial t} = \sum_{n=1}^{\infty} \left(-\frac{\partial}{\partial X} \right)^n D^{(n)}(X, t) p(X, t) \quad (54)$$

where the coefficients in the expansion are given as:

$$D^{(n)}(X) = \frac{1}{n!} M_n(t) \Big|_{t=0} = \frac{1}{n!} \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \langle |X(t + \Delta t) - z|^n \rangle \Big|_{X(t)=z} \quad (55)$$

In order to know how many terms in the Eq. (54) will be active for the SDE in Eq. (33) an Fokker-Planck equation for the pdf of the solution of Eq. (33) needs to be constructed. Towards this, let us consider a well behaved generic function $F(X, t)$ which is at least twice differential. For this function the Itô's lemma is written as:

$$dF(X, t) = \frac{\partial F(X, t)}{\partial X} dX(t) + \frac{1}{2} \frac{\partial^2 F(X, t)}{\partial X^2} g^2(X, t) dt \quad (56)$$

Substituting the Eq. (33) in the above equation gives the following:

$$dF(X, t) = \left(\frac{\partial F(X, t)}{\partial X} f(X, t) + \frac{1}{2} \frac{\partial^2 F(X, t)}{\partial X^2} g^2(X, t) \right) dt + \frac{\partial F(X, t)}{\partial X} g(X, t) dB(t) \quad (57)$$

Upon taking the derivative with respect to t on both sides yields:

$$\frac{dE[F(X, t)]}{dt} = E \left(\frac{\partial F(X, t)}{\partial X} f(X, t) + \frac{1}{2} \frac{\partial^2 F(X, t)}{\partial X^2} g^2(X, t) \right) \quad (58)$$

Expanding the above equation using expectation operator and noting that the last term in the expression is a stochastic integral, one can invoke the result $E \left[\int_0^t \frac{\partial F(X, s)}{\partial X} g(X, s) dB(s) \right] = 0$, where $E(\cdot)$ is the expectation operator to obtain the following,

$$\int_{-\infty}^{\infty} p(X, t) \frac{dF(X)}{dt} dX = \int_{-\infty}^{\infty} p(X, t) \left(\frac{\partial F(X, t)}{\partial X} f(X, t) + \frac{1}{2} \frac{\partial^2 F(X, t)}{\partial X^2} g^2(X, t) \right) dX \quad (59)$$

By expanding the terms in the above equation using integration by parts and assuming that $\lim_{X \rightarrow 0} p(X, t) \rightarrow 0$, the following weak form is obtained:

$$\int_{-\infty}^{\infty} F(X) \left(\frac{\partial p(X, t)}{\partial t} + \frac{\partial(p(X, t)f(X, t))}{\partial X} - \frac{1}{2} \frac{\partial^2(p(X, t)g^2(X, t))}{\partial X^2} \right) dX = 0 \quad (60)$$

As the function $F(X, t)$ is arbitrarily selected, the FPK equation is obtained:

$$\frac{\partial p(X, t)}{\partial t} = - \frac{\partial(p(X, t)f(X, t))}{\partial X} + \frac{1}{2} \frac{\partial^2(p(X, t)g^2(X, t))}{\partial X^2}; \quad p(0, X) = p_0(X) \quad (61)$$

In a general setting the above equation for higher dimensional diffusion process is expressed as:

$$\frac{\partial p(X, t)}{\partial t} = L_t p(X, t) = - \sum_i^m \frac{\partial(p(X, t)f_i(X, t))}{\partial X_i} + \frac{1}{2} \sum_i^m \sum_j^m \sum_k^n \frac{\partial^2(p(X, t)g_{i,k}(X, t)g_{j,k}(X, t))}{\partial X_i \partial X_j} \quad (62)$$

where, the random variable X satisfies the SDE in Eq. (33). At this point it is clear that there will be two terms active in the Eq. (33). Thus for $n=2$, the Eq. (54) gives,

$$\frac{\partial P(X, t)}{\partial t} = - \frac{\partial}{\partial X} [D^{(1)}(X, t)p(X, t)] + \frac{\partial^2}{\partial X^2} [D^{(2)}(X, t)p(X, t)] \quad (63)$$

On comparison of Eq. (54) and 63, it is straightforward to note that to estimate the drift and diffusion terms it suffices to estimate the first and second order moments of the variations of the random variable $X(t)$ and then calculate the coefficients $D^{(1)}$ and $D^{(2)}$, formally,

$$f(X, t) = D^{(1)}, \quad g^2(X, t) = D^{(2)} \quad (64)$$

To derive these coefficients in the Kramers-Moyal expansion, it is imperative to understand the one-step Itô-Taylor expansion of the random variable $X(t)$. Referring to the Eq. (57), the integral form the Itô-lemma can be expressed as follows:

$$F(X_{t+h}, t+h) = F(X_t, t) + \int_t^{t+h} \left\{ f(X_s, s)F'(X_s, s) + \frac{1}{2} g^2(X_s, s)F''(X_s, s) \right\} ds + \int_t^{t+h} g(X_s, s)F'(X_s, s)dB(s) \quad (65)$$

where, $F'(X_s, s)$ and $F''(X_s, s)$ denotes the first and second order partial derivatives with respect to system states. For the simplicity two stochastic operators are defined in the following form:

$$\begin{aligned}\mathfrak{J}^0(.) &= \frac{\partial(.)}{\partial t} + \sum_i^m f_i(X_t, t) \frac{\partial(.)}{\partial X_i} + \frac{1}{2} \sum_i^m \sum_j^n g_{i,k}(X_t, t) g_{j,k}(X_t, t) \frac{\partial^2(.)}{\partial X_i \partial X_j} \\ \mathfrak{J}^1(.) &= \sum_i^m \sum_k^n g_{i,k}(X_t, t) \frac{\partial(.)}{\partial X_i}\end{aligned}\quad (66)$$

Under which the Eq. (65) can be rephrased as:

$$F(X_{t+h}, t+h) = F(X_t, t) + \int_t^{t+h} \mathfrak{J}^0(F(X_s, s)) ds + \int_t^{t+h} \mathfrak{J}^1(F(X_s, s)) dB(s) \quad (67)$$

Then substituting $F(X_t, t) = X(t)$, one verifies that $\mathfrak{J}^0 X(t) = f(X_t, t)$, and $\mathfrak{J}^1 X(t) = g(X_t, t)$. This yields the first iteration:

$$\begin{aligned}X(t+h) &= X(t) + \int_t^{t+h} \mathfrak{J}^0(X(s)) ds + \int_t^{t+h} \mathfrak{J}^1(X(s)) dB(s) \\ &= X(t) + \int_t^{t+h} f(X_s, s) ds + \int_t^{t+h} g(X_s, s) dB(s)\end{aligned}\quad (68)$$

In order to perform the second iteration it is required to find the stochastic expansion of the terms $F(X_t, t) = f(X_s, s)$, and $F(X_t, t) = g(X_s, s)$. Thus, expanding $f(X_t, t)$, and $g(X_t, t)$ and using the operators in Eq. (66) yields,

$$\begin{aligned}f(X_s, s) &= f(X, t) + \int_t^{s_1} \mathfrak{J}^0(f(X_{s_2}, s_2)) ds_2 + \int_t^{s_1} \mathfrak{J}^1(f(X_{s_2}, s_2)) dB(s_2) \\ g(X_s, s) &= g(X, t) + \int_t^{s_1} \mathfrak{J}^0(g(X_{s_2}, s_2)) ds_2 + \int_t^{s_1} \mathfrak{J}^1(g(X_{s_2}, s_2)) dB(s_2)\end{aligned}\quad (69)$$

On substituting the above result in Eq. (68) and further iterating gives,

$$X(t+h) - X(t) = f(X, t) \int_t^{t+h} ds_1 + g(X, t) \int_t^{t+h} dB(s_1) + \mathfrak{J}^1(g(X, t)) \int_t^{t+h} \int_t^{s_1} dB(s_2) dB(s_1) + R \quad (70)$$

where the remainder term R is,

$$R = \int_t^{t+h} \int_t^{s_1} \mathfrak{J}^0(f(X_{s_2}, s_2)) ds_2 ds_1 + \left\{ \int_t^{t+h} \int_t^{s_1} \mathfrak{J}^1(f(X_{s_2}, s_2)) dB(s_2) ds_1 + \int_t^{t+h} \int_t^{s_1} \mathfrak{J}^0(g(X_{s_2}, s_2)) ds_2 dB(s_1) \right\} + \int_0^t \int_{t_0}^{s_1} \int_t^{s_2} \mathfrak{J}^1(g(X_{s_3}, s_3)) ds_3 dB(s_2) dB(s_1) + \int_0^t \int_{t_0}^{s_1} \int_t^{s_2} \mathfrak{J}^1 \mathfrak{J}^1(g(X_{s_3}, s_3)) dB(s_3) dB(s_2) dB(s_1) \quad (71)$$

It can be noticed that the above equation is infinitely expandable using the Eq. (66). For further treatment, the first moment is taken on both side. Noting the results $\langle \int_t^{t+h} dB(s_1) \rangle = 0$, and $\langle \int_{t_0}^t \int_{t_0}^{s_1} dB(s_2) dB(s_1) \rangle = \langle \int_{t_0}^t B(s_1) dB(s_1) - B(t) \int_{t_0}^t dB(s_1) \rangle = 0$,

$$\begin{aligned}\langle X(t+h) - X(t) \rangle &= f(X, t) h + g(X, t) \left\langle \int_t^{t+h} dB(s_1) \right\rangle + \mathfrak{J}^1(g(X, t)) \left\langle \int_t^{t+h} \int_t^{s_1} dB(s_2) dB(s_1) \right\rangle + \langle R \rangle \\ &= f(X, t) h + \langle R \rangle\end{aligned}\quad (72)$$

If the number of iteration is k then the Brownian integrals $\int_t^{t+h} \int_t^{s_1} \dots \int_t^{s_k} dB(s_{k+1}) \dots dB(s_2) dB(s_1)$ of multiplicity k have a contribution proportional to $h^{k/2}$, the time integrals $\int_t^{t+h} \int_t^{s_1} \dots \int_t^{s_k} ds_{k+1} \dots ds_2 ds_1$ have a contribution proportional to h^k and the combination between them shares a contribution between $h^{k/2}$ and h^k . Thus it is easy to infer that as $h \rightarrow 0$, the higher order terms in the remainder R will vanish. After invoking this facts in the above expression the first coefficient in Kramers-Moyal expansion is obtained as,

$$D^{(1)} = f(X, t) = \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \left\langle X^{(1)}(t + \Delta t) - z \right\rangle \Big|_{X(t)=z} \quad (73)$$

To derive the second coefficient it is only required to find the quadratic variation of the increment process of $X(t)$. Upon taking second moment on both sides of Eq. (70) the following is obtained,

$$\begin{aligned}\langle |X(t+h) - X(t)|^2 \rangle &= f(X, t) h + g(X, t) \Delta B + \langle R \rangle \\ &= f^2(X, t) h^2 + 2f(X, t)g(X, t)h\Delta B + g^2(X, t)(\Delta B)^2 + \langle R \rangle \\ &= g^2(X, t)h + \langle R \rangle\end{aligned}\quad (74)$$

In the above prove the Itô identifies are utilized which states that under the mean square convergence theory the quadratic variation of the time and Brownian increments are given as $ds_2 ds_1 = 0$, $ds_1 dW(s_1) = 0$, $dW(s_2) ds_1 = 0$, $dW(s_1) dW(s_2) = 0$. The deterministic time integrals will vanish automatically since they have finite variation and zero quadratic covariation and the other higher order Brownian integrals will vanish as $h \rightarrow 0$. Thus the second coefficient in the Kramers-Moyal expansion is obtained as,

$$D^{(2)} = g^2(X, t) = \frac{1}{2} \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \left\langle |X^{(1)}(t + \Delta t) - z|^2 \right\rangle_{X(t)=z} \quad (75)$$

From the ergodic assumption of the evolution of process $X(t)$ the time average $\langle \cdot \rangle$ is often replaced by expectation operator $E(\cdot)$. For m -variables $\mathbf{X} = \{X_1, X_2, \dots, X_m\}$, the diffusion process has the form,

$$dX_i(t) = f_i(\mathbf{X}_t, t) + \sum_j^n g_{ij}(\mathbf{X}_t, t) dB_j(t); \quad i = 1, 2, \dots, m \quad (76)$$

The properties of the Brownian motion are: $\langle B_i(t) \rangle = 0$ and $\langle B_i(t) B_j(s) \rangle = \min(t, s)$. If the covariation matrix of the diffusion components is given by $\Gamma(\mathbf{X}_t, t) = g(\mathbf{X}_t, t)g(\mathbf{X}_t, t)^T$, then, the drift $f_i(\mathbf{X}, t)$ of i^{th} -diffusion process and the ij^{th} -element of the covariation matrix $\Gamma(\mathbf{X}, t)$ can be estimated as,

$$\begin{aligned} f_i(\mathbf{X}, t) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[X_i(t + \Delta t) - z_i] \Big|_{X_k(t)=z_k} \quad \forall k = 1, 2, \dots, m \\ \Gamma_{ij}(\mathbf{X}, t) &= \frac{1}{2} \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} E[(X_i(t + \Delta t) - z_i)(X_j(t + \Delta t) - z_j)] \Big|_{X_k(t)=z_k} \quad \forall k = 1, 2, \dots, m \end{aligned} \quad (77)$$

B Conditional probability distributions of the discontinuous Spike and Slab priors (DSS)

The DSS prior model is described in section . Drawing sample from the joint distribution is intractable but possible through Monte Carlo Markov Chain methods. In the present work, the Gibbs sampling technique is utilized to draw the sample for the random variables \mathbf{Y} , $\boldsymbol{\theta}$, \mathbf{Z} , ϑ_s , σ^2 and p_0 . For drawing samples using Gibbs sampling the random variables need to be conditioned on other variables. However, due to the DAG structure in Fig. 1, the dependencies of the conditional distributions on other random variables can be relaxed to some extent as follows:

$$\begin{aligned} p(p_0 | \mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, \vartheta_s, \sigma^2) &= p(p_0 | \mathbf{Z}) \\ p(\vartheta_s | \mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, p_0, \sigma^2) &= p(\vartheta_s | \boldsymbol{\theta}, \mathbf{Z}, \sigma^2) \\ p(\sigma^2 | \mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, p_0, \vartheta_s) &= p(\sigma^2 | \mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, \vartheta_s) \\ p(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{Z}, p_0, \vartheta_s, \sigma^2) &= p(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{Z}, \vartheta_s, \sigma^2) \\ p(\mathbf{Z} | \mathbf{Y}, \boldsymbol{\theta}, p_0, \vartheta_s, \sigma^2) &= p(\mathbf{Z} | \boldsymbol{\theta}, p_0, \vartheta_s, \sigma^2) \end{aligned} \quad (78)$$

The elements of the weight vector $\boldsymbol{\theta}$ for which the latent variable $Z_k \neq 1$, becomes an absorbing state in the Markov chain. However to achieve a stationary distribution one needs to construct a irreducible Markov chain. Thus the conditionals over weight vector is eliminated by marginalizing the conditional distributions with respect to the vector $\boldsymbol{\theta}$.

$$\begin{aligned} \int p(\sigma^2 | \mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, \vartheta_s) d\boldsymbol{\theta} &= \int \frac{p(\sigma^2, \mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, \vartheta_s)}{p(\mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, \vartheta_s)} d\boldsymbol{\theta} = p(\sigma^2 | \mathbf{Y}, \mathbf{Z}, \vartheta_s) \\ \int p(\mathbf{Z} | \mathbf{Y}, \boldsymbol{\theta}, p_0, \vartheta_s, \sigma^2) d\boldsymbol{\theta} &= \int \frac{p(\mathbf{Z}, \mathbf{Y}, \boldsymbol{\theta}, p_0, \vartheta_s, \sigma^2)}{p(\mathbf{Y}, \boldsymbol{\theta}, p_0, \vartheta_s, \sigma^2)} d\boldsymbol{\theta} = \frac{p(\mathbf{Z}, \mathbf{Y}, p_0, \vartheta_s, \sigma^2)}{p(\mathbf{Y}, p_0, \vartheta_s, \sigma^2)} = p(\mathbf{Z} | \mathbf{Y}, p_0, \vartheta_s, \sigma^2) \end{aligned} \quad (79)$$

The conditional distribution $p(p_0 | \mathbf{Z})$:

$$\begin{aligned} p(p_0 | \mathbf{Z}) &= p(\mathbf{Z} | p_0) p(p_0) \\ &\propto \left(\prod_{k=1}^K p_0^{Z_k} (1-p_0)^{1-Z_k} \right) \left(p_0^{\alpha_p-1} (1-p_0)^{\beta_p-1} \right) \\ &\propto p_0^{\alpha_p+h_Z} (1-p_0)^{\beta_p+K-h_Z} \end{aligned} \quad (80)$$

where, $h_Z = \sum_{k=1}^K Z_k$. Given the latent vector \mathbf{Z} , p_0 is sampled as, $p_0 | \mathbf{Z} \sim Beta(\alpha_p + h_Z, \beta_p + K - h_Z)$.

The conditional distribution $p(\vartheta_s | \boldsymbol{\theta}, \mathbf{Z}, \sigma^2)$:

$$\begin{aligned}
p(\vartheta_s | \boldsymbol{\theta}, \mathbf{Z}, \sigma^2) &\propto p(\boldsymbol{\theta} | \mathbf{Z}, \vartheta_s, \sigma^2) p(\vartheta_s) p(\mathbf{Z}) p(\sigma^2) \\
&\propto p(\boldsymbol{\theta} | \mathbf{Z}, \vartheta_s, \sigma^2) p(\vartheta_s) \\
&\propto \mathcal{N}(\boldsymbol{\theta}_r | \mathbf{0}, \sigma^2 \boldsymbol{\vartheta}_s \mathbf{R}_{0,r}) \mathcal{IG}(\alpha_\theta, \beta_\theta) \\
&\propto \frac{1}{(\vartheta_s \sigma^2)^{r/2} |\mathbf{R}_{0,r}|^{1/2}} \exp\left(-\frac{\boldsymbol{\theta}_r^T \mathbf{R}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2 \vartheta_s}\right) \vartheta_s^{-\alpha_\theta - 1} \exp\left(-\frac{\beta_\theta}{\vartheta_s}\right) \\
&\propto \vartheta_s^{-\left(\alpha_\theta + \frac{r}{2}\right) - 1} \exp\left(-\frac{\beta_\theta + (\boldsymbol{\theta}_r^T \mathbf{R}_{0,r}^{-1} \boldsymbol{\theta}_r / 2\sigma^2)}{\vartheta_s}\right)
\end{aligned} \tag{81}$$

where, r is the number of elements of the weight vector $\boldsymbol{\theta}$ for which the latent variable $Z_k = 1$. In the above it is to be noted that $p(\mathbf{Z})$ and $p(\sigma^2)$ are constants when ϑ_s is sampled. This can be observed from the DAG structure in Fig. 1. Thus given $\boldsymbol{\theta}, \mathbf{Z}$ and σ^2 and ϑ_s is sampled as, $\vartheta_s | \boldsymbol{\theta}, \mathbf{Z}, \sigma^2 \sim \mathcal{IG}\left(\alpha_\theta + \frac{r}{2}, \beta_\theta + \frac{\boldsymbol{\theta}_r^T \mathbf{R}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2}\right)$

The conditional distribution $p(\sigma^2 | \mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, \vartheta_s)$:

$$\begin{aligned}
p(\sigma^2 | \mathbf{Y}, \mathbf{Z}, \vartheta_s) &\propto \int p(\mathbf{Y}, \boldsymbol{\theta}, \mathbf{Z}, \vartheta_s, \sigma^2) d\boldsymbol{\theta} \\
&\propto \left(\int p(\mathbf{Y} | \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} | \mathbf{Z}, \vartheta_s, \sigma^2) d\boldsymbol{\theta} \right) p(\mathbf{Z}) p(\vartheta_s) p(\sigma^2) \\
&\propto \left(\int p(\mathbf{Y} | \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} | \mathbf{Z}, \vartheta_s, \sigma^2) d\boldsymbol{\theta} \right) p(\sigma^2)
\end{aligned} \tag{82}$$

From the DAG structure it can be understood that $p(\mathbf{Z})$ and $p(\vartheta_s)$ will act as a constant when σ^2 is sampled. Since the θ_k -values corresponding to the spike distribution does not contribute to the basis function selection, the remaining weights belonging to the slab denoted by $\boldsymbol{\theta}_r$ is used. Then one can expand the integrand $p(\mathbf{Y} | \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} | \mathbf{Z}, \vartheta_s, \sigma^2)$, as,

$$\begin{aligned}
&p(\mathbf{Y} | \boldsymbol{\theta}, \sigma^2) p(\boldsymbol{\theta} | \mathbf{Z}, \vartheta_s, \sigma^2) \\
&= \mathcal{N}(\mathbf{Y} | \mathbf{L}_r \boldsymbol{\theta}_r, \sigma^2 \mathbf{I}_{N \times N}) \mathcal{N}(\boldsymbol{\theta}_r | \mathbf{0}, \sigma^2 \boldsymbol{\vartheta}_s \mathbf{R}_{0,r}) \\
&= \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{Y} - \mathbf{D}_r \boldsymbol{\theta}_r)^T (\mathbf{Y} - \mathbf{D}_r \boldsymbol{\theta}_r)}{2\sigma^2}\right) \frac{(|\mathbf{R}_{0,r}^{-1}|)^{1/2}}{(2\pi\vartheta_s \sigma^2)^{r/2}} \exp\left(-\frac{\boldsymbol{\theta}_r^T \mathbf{R}_{0,r}^{-1} \boldsymbol{\theta}_r}{2\sigma^2 \vartheta_s}\right) \\
&= \frac{1}{(2\pi\sigma^2)^{N/2}} \frac{(|\mathbf{R}_{0,r}^{-1}|)^{1/2}}{(2\pi\vartheta_s \sigma^2)^{r/2}} \exp\left(-\frac{(\boldsymbol{\theta}_r - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta}_r - \boldsymbol{\mu})}{2\sigma^2}\right) \exp\left(-\frac{(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{2\sigma^2}\right)
\end{aligned} \tag{83}$$

where, $\boldsymbol{\Sigma}^{-1} = (\mathbf{L}_r^T \mathbf{L}_r + \vartheta_s^{-1} \mathbf{R}_{0,r}^{-1})$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma} \mathbf{L}_r^T \mathbf{Y}$. Upon integration of the the above expression about $\boldsymbol{\theta}_r$, one obtains:

$$\begin{aligned}
p(\sigma^2 | \mathbf{Y}, \mathbf{Z}, \vartheta_s) &\propto \frac{1}{(\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{2\sigma^2}\right) p(\sigma^2) \\
&\propto \frac{1}{(\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{2\sigma^2}\right) (\sigma^2)^{-\alpha_\sigma - 1} \exp\left(-\frac{\beta_\sigma}{\sigma^2}\right) \\
&\propto (\sigma^2)^{-(\alpha_\sigma + 0.5N) - 1} \exp\left(-\frac{\beta_\sigma + \frac{1}{2}(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})}{\sigma^2}\right)
\end{aligned} \tag{84}$$

The random variable σ^2 is sampled as, $\sigma^2 | \mathbf{Y}, \mathbf{Z}, \vartheta_s \sim \mathcal{IG}\left(\alpha_\sigma + \frac{N}{2}, \beta_\sigma + \frac{1}{2}(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})\right)$.

The conditional distribution $p(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{Z}, \vartheta_s, \sigma^2)$: The elements of the weight vector $\boldsymbol{\theta}$ corresponding to the spike distribution i.e. the weights for which the latent variable $z_k = 0; k = 1 \dots K$, are assigned the value 0. The conditional distribution for the

weights belonging to the slab distribution denoted by θ_r is derived as follows:

$$\begin{aligned}
p(\theta_r | \mathbf{Y}, \vartheta_s, \sigma^2) &\propto \mathcal{N}(\mathbf{Y} | \mathbf{L}_r \theta_r, \sigma^2 \mathbf{I}_{N \times N}) \mathcal{N}(\theta_r | \mathbf{0}, \sigma^2 \vartheta_s \mathbf{R}_{0,r}) \\
&\propto \exp\left(-\frac{(\mathbf{Y} - \mathbf{L}_r \theta_r)^T (\mathbf{Y} - \mathbf{L}_r \theta_r)}{2\sigma^2}\right) \exp\left(-\frac{\theta_r^T \mathbf{R}_{0,r}^{-1} \theta_r}{2\sigma^2 \vartheta_s}\right) \\
&\propto \exp\left(-\frac{\mathbf{Y}^T \mathbf{Y} + \theta_r^T \Sigma^{-1} \theta_r - 2\theta_r^T \Sigma^{-1} \mu}{2\sigma^2}\right) \\
&\propto \exp\left(-\frac{(\theta_r - \mu)^T \Sigma^{-1} (\theta_r - \mu)}{2\sigma^2}\right)
\end{aligned} \tag{85}$$

where, $\Sigma^{-1} = (\mathbf{L}_r^T \mathbf{L}_r + \vartheta_s^{-1} \mathbf{R}_{0,r}^{-1})$ and $\mu = \Sigma \mathbf{L}_r^T \mathbf{Y}$. Then, the random values of θ_r can be sampled as $\theta_r | \mathbf{Y}, \vartheta_s, \sigma^2 \sim \mathcal{N}(\mu, \sigma^2 \Sigma)$.

The conditional distribution $p(Z | \theta, p_0, \vartheta_s, \sigma^2)$: The latent variable Z_k corresponding to the k^{th} element of the weight vector are assigned the value 0 or 1 independently. The conditional probability distributions of Z_k are estimated by comparing the probabilities that the k^{th} latent variable $Z_k = 1$ will assume a value 1 or 0, given the values of ϑ_s , p_0 and remaining values of the latent vector \mathbf{Z}_{-k} . Here, the term \mathbf{Z}_{-k} denotes the latent vector \mathbf{Z} whose k^{th} element is removed. Let u_k denotes the probability with which the k^{th} latent variable Z_k takes a value 1. The probability u_k is then found as,

$$\begin{aligned}
u_k &= \frac{p(Z_k = 1 | \mathbf{Y}, \mathbf{Z}_{-k}, \vartheta_s, p_0)}{p(Z_k = 1 | \mathbf{Y}, \mathbf{Z}_{-k}, \vartheta_s, p_0) + p(Z_k = 0 | \mathbf{Y}, \mathbf{Z}_{-k}, \vartheta_s, p_0)} \\
&= \frac{p(\mathbf{Y} | Z_k = 1, \mathbf{Z}_{-k}, \vartheta_s) p(Z_k = 1 | p_0)}{p(\mathbf{Y} | Z_k = 1, \mathbf{Z}_{-k}, \vartheta_s) p(Z_k = 1 | p_0) + p(\mathbf{Y} | Z_k = 0, \mathbf{Z}_{-k}, \vartheta_s) p(Z_k = 0 | p_0)} \\
&= \frac{p(\mathbf{Y} | Z_k = 1, \mathbf{Z}_{-k}, \vartheta_s) p_0}{p(\mathbf{Y} | Z_k = 1, \mathbf{Z}_{-k}, \vartheta_s) p_0 + p(\mathbf{Y} | Z_k = 0, \mathbf{Z}_{-k}, \vartheta_s) (1 - p_0)} \\
&= \frac{p_0}{p_0 + \lambda_k (1 - p_0)}
\end{aligned} \tag{86}$$

where, $\lambda_k = \frac{p(\mathbf{Y} | Z_k = 0, \mathbf{Z}_{-k}, \vartheta_s)}{p(\mathbf{Y} | Z_k = 1, \mathbf{Z}_{-k}, \vartheta_s)}$. The marginal likelihood function $p(\mathbf{Y} | \mathbf{Z}, \vartheta_s)$ was derived by integrating out the random variables θ and σ^2 from the original likelihood function, which follows,

$$p(\mathbf{Y} | \mathbf{Z}, \vartheta_s) = \int p(\mathbf{Y} | \mathbf{Z}, \vartheta_s, \sigma^2) p(\sigma^2) d\sigma^2 \tag{87}$$

where, $p(\mathbf{Y} | \mathbf{Z}, \vartheta_s, \sigma^2)$ is obtained by marginalizing the likelihood function with respect to the variable θ . Considering only the weights whose corresponding latent variables $Z_k = 0$, denoted by θ_r , the probability $p(\mathbf{Y} | \mathbf{Z}, \vartheta_s, \sigma^2)$ is obtained as,

$$\begin{aligned}
p(\mathbf{Y} | \mathbf{Z}, \vartheta_s, \sigma^2) &= \int p(\mathbf{Y}, \theta | \mathbf{Z}, \vartheta_s, \sigma^2) d\theta \\
&= \int p(\mathbf{Y} | \theta_r, \sigma^2) p(\theta_r | \vartheta_s, \sigma^2) d\theta_r \\
&= \int \mathcal{N}(\mathbf{Y} | \mathbf{L}_r \theta_r, \sigma^2 \mathbf{I}_{N \times N}) \mathcal{N}(\theta_r | \mathbf{0}, \sigma^2 \vartheta_s \mathbf{R}_{0,r}) d\theta_r \\
&= \int \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{Y} - \mathbf{L}_r \theta_r)^T (\mathbf{Y} - \mathbf{L}_r \theta_r)}{2\sigma^2}\right) \frac{(|\mathbf{R}_{0,r}^{-1}|)^{1/2}}{(2\pi\vartheta_s\sigma^2)^{r/2}} \exp\left(-\frac{\theta_r^T \mathbf{R}_{0,r}^{-1} \theta_r}{2\sigma^2 \vartheta_s}\right) d\theta_r \\
&= \frac{1}{(2\pi\sigma^2)^{N/2}} \frac{(|\mathbf{R}_{0,r}^{-1}|)^{1/2} (|\Sigma^{-1}|)^{1/2}}{(\vartheta_s)^{r/2}} \exp\left(-\frac{(\mathbf{Y}^T \mathbf{Y} - \mu^T \Sigma \mu)}{2\Sigma^2}\right)
\end{aligned} \tag{88}$$

where, $\Sigma = (\mathbf{L}_r^T \mathbf{L}_r + \vartheta_s^{-1} \mathbf{R}_{0,r}^{-1})^{-1}$ and $\mu = \Sigma \mathbf{L}_r^T \mathbf{Y}$. The operator $|.|$ denotes the determinant. With the above result then, $p(\mathbf{Y} | \mathbf{Z}, \vartheta_s)$

can be derived as,

$$\begin{aligned}
p(\mathbf{Y} | \mathbf{Z}, \vartheta_s) &= \frac{(|\mathbf{R}_{0,r}^{-1}|)^{1/2} (|\boldsymbol{\Sigma}^{-1}|)^{1/2}}{(2\pi)^{N/2} (\vartheta_s)^{r/2}} \int \frac{1}{(\sigma^2)^{N/2}} \exp\left(-\frac{(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma} \boldsymbol{\mu})}{2\sigma^2}\right) \mathcal{IG}(\alpha_\sigma, \beta_\sigma) d\sigma^2 \\
&= \frac{(|\mathbf{R}_{0,r}^{-1}|)^{1/2} (|\boldsymbol{\Sigma}^{-1}|)^{1/2}}{(2\pi)^{N/2} (\vartheta_s)^{r/2}} \frac{(\beta_\sigma)^{\alpha_\sigma}}{\Gamma(\alpha_\sigma)} \int \frac{1}{(\sigma^2)^{\alpha_\sigma + N/2 + 1} \exp\left(-\frac{\beta_\sigma + \frac{1}{2}(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma} \boldsymbol{\mu})}{\sigma^2}\right)} d\sigma^2 \\
&= \frac{(|\mathbf{R}_{0,r}^{-1}|)^{1/2} (|\boldsymbol{\Sigma}^{-1}|)^{1/2}}{(2\pi)^{N/2} (\vartheta_s)^{r/2}} \frac{(b_\sigma)^{\alpha_\sigma}}{\Gamma(\alpha_\sigma)} \frac{\Gamma\left(\alpha_\sigma + \frac{N}{2}\right)}{\left(\beta_\sigma + \frac{1}{2}(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma} \boldsymbol{\mu})\right)^{\left(\alpha_\sigma + \frac{N}{2}\right)}}
\end{aligned} \tag{89}$$

where, the operator $\Gamma(\cdot)$ denotes the Gamma function. To summarise, the random variables Z_k can be computed from the Bernoulli distribution with the parameter u_k as: $Z_k; k = 1, \dots, K \mid \mathbf{Y}, \vartheta_s, p_0 \sim \text{Bern}(u_k)$. Where, the marginalised likelihood function is obtained as,

$$p(\mathbf{Y} | \mathbf{Z}, \vartheta_s) = \begin{cases} \frac{\Gamma\left(\alpha_\sigma + \frac{N}{2}\right)}{(2\pi)^{N/2} \frac{(\beta_\sigma)^{\alpha_\sigma}}{\Gamma(\alpha_\sigma)}} \frac{1}{\left(\beta_\sigma + \frac{1}{2}(\mathbf{Y}^T \mathbf{y})\right)^{\left(\alpha_\sigma + \frac{N}{2}\right)}} & , \text{when all } \{Z_k; k = 1, \dots, K\} = 0 \\ \frac{\Gamma\left(\alpha_\sigma + \frac{N}{2}\right)}{(2\pi)^{N/2} (\vartheta_s)^{r/2}} \frac{(\beta_\sigma)^{\alpha_\sigma}}{\Gamma(\alpha_\sigma)} \frac{\left(|\mathbf{R}_{0,r}^{-1}|^{1/2} (|\boldsymbol{\Sigma}^{-1}|)^{1/2}\right)}{\left(\beta_\sigma + \frac{1}{2}(\mathbf{Y}^T \mathbf{Y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma} \boldsymbol{\mu})\right)^{\left(\alpha_\sigma + \frac{N}{2}\right)}} & , \text{otherwise} \end{cases} \tag{90}$$