

Bayesian autoencoders for data-driven discovery of coordinates, governing equations and fundamental constants

L. Mars Gao¹, J. Nathan Kutz^{2*}

¹ Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, WA

² Department of Applied Mathematics and Electrical and Computer Engineering, University of Washington, Seattle, WA

Abstract

Recent progress in autoencoder-based sparse identification of nonlinear dynamics (SINDy) under ℓ_1 constraints allows joint discoveries of governing equations and latent coordinate systems from spatio-temporal data, including simulated video frames. However, it is challenging for ℓ_1 -based sparse inference to perform correct identification for real data due to the noisy measurements and often limited sample sizes. To address the data-driven discovery of physics in the low-data and high-noise regimes, we propose Bayesian SINDy autoencoders, which incorporate a hierarchical Bayesian sparsifying prior: Spike-and-slab Gaussian Lasso. Bayesian SINDy autoencoder enables the joint discovery of governing equations and coordinate systems with a theoretically guaranteed uncertainty estimate. To resolve the challenging computational tractability of the Bayesian hierarchical setting, we adapt an adaptive empirical Bayesian method with Stochastic gradient Langevin dynamics (SGLD) which gives a computationally tractable way of Bayesian posterior sampling within our framework. Bayesian SINDy autoencoder achieves better physics discovery with lower data and fewer training epochs, along with valid uncertainty quantification suggested by the experimental studies. The Bayesian SINDy autoencoder can be applied to real video data, with accurate physics discovery which correctly identifies the governing equation and provides a close estimate for standard physics constants like gravity g , for example, in videos of a pendulum.

Keywords— model discovery, dynamical systems, machine learning, bayesian deep learning, bayesian sparse inference, autoencoder

1 Introduction

Calculus-based models fundamentally relate the rates of change of quantities of interest in time and space through differential and partial differential equations. From population models to turbulence, physics and engineering principles are rooted in such governing equations. In the modern era of big data, there is a growing demand to transform rich spatio-temporal data into descriptive physical models in an automated, data-driven fashion. Video data, for example, contain optical snapshots of an observed dynamical system described by some specific governing equations. To understand the underlying physics, it is important not only to identify the equations, but to discover dependent state variables (sparse representations) directly from the video frames. To discover such an underlying sparse representation, it is essential to find a correct coordinate transformation (latent space or manifold) to compress the data into a low-dimensional space. Principal Component Analysis (PCA) is often applied to obtain a low-dimensional subspace with linearity constraints. For nonlinear transformations, Autoencoders with neural networks are frequently applied as a nonlinear extension of PCA [5, 39]. To identify the governing equation given the latent sparse representations, one can apply data-driven methods via sparse regression [51, 74, 75]. The sparse regression-based model discovery enables a computationally efficient way of identifying governing equation with convergence guarantees [104]. In this case of physics discovery from videos, Champion et al. [15] propose SINDy autoencoders which can jointly discover on

*Corresponding author

coordinates and equations for synthetic video data. More recently, Chen et al. [19] introduce the Neural State Variables learning from video data, which enables automated discovery of the latent state variables from the underlying dynamical system. [19].

The automated discovery of coordinates and equations in real video is significantly more challenging compared to these prior works [16, 19]. First, the temporal derivatives are not directly available from the video data. To resolve the missing temporal derivatives, one has to find an approximation for the temporal derivatives from discrete video frames, which is frequently numerically unstable and will inject a high level of noise. Indeed, lighting effects alone can greatly compromise derivative estimates. Second, the sample size of the real video dataset may not be sufficiently large in comparison with synthetic video [15]. These pragmatic constraints hinder the real video data from having a statistically sufficient sample size for learning. Thus automated discovery from real video data is much more challenging due to being in the low-data and high-noise limit.

In this low-data and high-noise regime, Bayesian sparse regression methods generally have significant advantages from both a theoretical and practical perspective [31, 64, 79]. The spike-and-slab prior [44] (e.g. Bernoulli-Gaussian, Bernoulli-Laplace [3]) with hierarchical Bayesian settings has a proven success for both sparse variable selection and uncertainty quantification [71]. In data-driven model discovery, Bayesian SINDy [40] with the spike-and-slab prior has also shown an advantage for correct model identification in the case study of Lynx-hare population [38] under the very low-data limit. However, a significant limitation of Bayesian methods comes from its high computational cost that hinders both speed and scalability. The MCMC-based hierarchical Bayesian model sampling requires a considerably long run because the underlying stochastic binary search grows exponentially with the number of parameters. Additionally, it is very costly to compute the full gradient of the entire video dataset for the MCMC sampling. Therefore, even if the Bayesian methods are much more powerful in the low-data regime, it is highly nontrivial to design a feasible Bayesian solution for the discovery of governing equations and coordinate systems given the computationally intractability.

In this paper, we propose Bayesian SINDy autoencoders that extends SINDy Autoencoder [15] into a Bayesian learning framework. Bayesian SINDy autoencoders can perform a joint discovery for governing equations and coordinate systems in a computationally tractable manner. Specifically, we apply the Spike-and-slab Gaussian-Laplace (SSGL) prior to the intermediate SINDy module to accelerate the sparse identification process of governing equation discovery. To resolve the computational burden arising from the hierarchical Bayesian model, we consolidate the Bayesian sampling procedure via Stochastic Gradient Langevin Dynamics (SGLD) with an adaptive empirical Bayesian variable selection method using Expectation-maximization. Instead of computing the full batch gradient, SGLD evaluates mini-batch gradients with injected random Gaussian noise, which is theoretically valid to generate Langevin-based proposal distribution [17, 95]. The mini-batch gradient learning naturally fits into the training of deep neural networks and relaxes the scalability issue at the same time. On the other hand, the adaptive Bayesian Expectation-maximization variable selection (EMVS) performs variable selection by optimizing the latent inclusion probability of each variable, which avoids the previously required lengthy binary stochastic search. Adapting these two ideas into the SINDy module, we can simultaneously perform an accelerated governing equation identification under the low-data, high-noise limit with a trustworthy uncertainty quantification through the power of Bayesian estimation.

In our numerical experiments, Bayesian SINDy autoencoders achieve accelerated discovery in governing equations and coordinate systems under the SSGL prior. In addition to the model discovery in synthetic datasets, from a real video data on a single moving pendulum, we obtain correct governing equation discovery with a close estimate of the gravity constant $\hat{g} = -9.876$

with only 390 data snapshots. With precise understanding of the underlying physics, Bayesian SINDy autoencoder enables explainable, trustworthy, and robust video predictions. In summary, the contribution of this paper is threefold:

1. We propose the Bayesian SINDy autoencoder with Spike-and-slab Gaussian-Laplace prior to accelerated sparse inference under low-data and high-noise environments.
2. We utilize Stochastic Gradient Langevin Dynamics to perform posterior sampling in Bayesian SINDy autoencoder with valid uncertainty estimations.
3. We conduct extensive experiments with successful joint discoveries on coordinate systems and governing equations for both synthetic and real datasets. Remarkably, we achieve a correct physics discovery with a 14 seconds recording on our experiment on a pendulum.

2 Background

The current work is built upon two primary mathematical innovations: (i) sparse regression used in the SINDy algorithm, including learning latent representations, and (ii) Bayesian learning with sparse priors. A quick review of each is given in order to better inform the reader how they are combined into our Bayesian SINDy autoencoder framework.

2.1 Sparse identification of nonlinear dynamics

We review the *sparse identification of nonlinear dynamics* (SINDy) [10] algorithm, which utilizes sparse regression to identify the latent dynamical system from snapshot data. SINDy takes snapshot data $\mathbf{x}(t) \in \mathbb{R}^n$ and aims to discover the underlying dynamical system

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t)). \quad (1)$$

The snapshots are collected by measurements at time $t \in [t_1, t_m]$, and the function \mathbf{f} characterizes the dynamics. Assuming the temporal derivatives of the snapshot are available from data, SINDy forms data matrices in the following way:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}(t_1) \\ \mathbf{x}(t_2) \\ \vdots \\ \mathbf{x}(t_m) \end{pmatrix}, \quad \dot{\mathbf{X}} = \begin{pmatrix} \dot{\mathbf{x}}_1(t_1) \\ \dot{\mathbf{x}}_1(t_2) \\ \vdots \\ \dot{\mathbf{x}}_1(t_m) \end{pmatrix},$$

with $\mathbf{X}, \dot{\mathbf{X}} \in \mathbb{R}^{m \times n}$. The candidate function library is constructed by p candidate model term θ_j 's that $\Theta(\mathbf{X}) = [\theta_1(\mathbf{X}) \cdots \theta_p(\mathbf{X})] \in \mathbb{R}^{m \times p}$. A common choice of candidate functions are polynomials in \mathbf{x} targeting common canonical models of dynamical systems [34]. The Fourier library is also very common with $\sin(\cdot)$ and $\cos(\cdot)$ terms. In summary, we build a model between \mathbf{X} and $\dot{\mathbf{X}}$ that

$$\dot{\mathbf{X}} = \Theta(\mathbf{X})\Xi$$

where the unknown matrix $\Xi = (\xi_1 \ \xi_2 \ \cdots \ \xi_n) \in \mathbb{R}^{p \times n}$ is the set of coefficients. The sparse inference on Ξ enables sparse identification of the dynamical system \mathbf{f} . For high-dimensional systems, the goal is to jointly identify a low-dimensional state $\mathbf{z} = \varphi(\mathbf{x})$ with dynamics $\dot{\mathbf{z}} = \mathbf{g}(\mathbf{z})$. The standard SINDy approach uses a sequentially thresholded least squares algorithm to perform

sparse inference [10], which is a proxy for ℓ_0 optimization [107] with convergence guarantees [104]. Small et al [81] and Yao and Bollt [102] previously formulated the dynamical system identification without sparsity constraints. These methods provide a computationally efficient counterpart to other model discovery frameworks [78, 91].

In an alternative approach, Hirsh et al. [40] utilize Bayesian sparse regression techniques in model discovery, which have improved the robustness of SINDy under high-noise and low-data settings. Bayesian sparse inference models the sparsity with various probabilistic priors $p(\Xi)$ like the Spike-and-slab [44], regularized horseshoe [13, 14], Laplace priors [67, 88] and so on. After defining the Bayesian likelihood function with sparsifying priors, Bayesian SINDy generates posterior distribution via the No-U-Turn MCMC sampler [42]. Since MCMC for sparse inference can be extremely computationally demanding, [71, 72] propose coordinate descent sparse inference methods via Spike-and-slab prior, targeting the mode detection. As an approximation to Bayesian inference, Fasel et al. [29] combine ensembling techniques via bootstrapping to perform uncertainty estimation and accelerated variable selection for system identification.

SINDy has been widely applied to model discovery for many scientific scenarios including fluid dynamics [12, 22, 33, 55–57], nonlinear optics [82], turbulence closures [6, 7, 77], ocean closures [103], chemical reaction [43], plasma dynamics [2, 21, 48], structural modeling [52], and for model predictive control [46]. There are many extensions of SINDy, including the identification of partial differential equations [74, 75], multiscale physics [16], parametrically dependent dynamical models [73], time-dependent PDEs [18], switching dynamical systems [60], rational function nonlinearities [45, 61], control inputs [46], constraints on symmetries [56], control for stability [47], control for robustness [2, 63, 63, 69, 70, 76], stochastic dynamical systems [9, 11], and multidimensional approximation on tensors [30].

A related and important extension to the SINDy framework is the SINDy autoencoder [15], which embeds SINDy into the training process of deep autoencoders. The SINDy autoencoder achieves remarkable performances on high-dimensional synthetic data to jointly discover coordinate systems and governing equations. Due to over-parametrization and non-convexity, deep neural networks do not generally have explainable guarantees for inference and predictions. Dynamical system learning and forecasting is typically an extrapolatory problem by nature. Therefore, the interpretability of models is very important to understand. In this case, SINDy autoencoder is satisfactory due to the transparency of the learning process. Specifically, the encoder and decoder focus on the specific task of learning only a coordinate transformation, with the SINDy layer targeting the inference of latent dynamical systems. The SINDy autoencoder can not only identify the governing equation from high-dimensional data but can also perform trustworthy predictions based on future dynamics. The limitations of the SINDy autoencoder are exactly what our Bayesian framework addresses.

2.2 Bayesian and sparse deep learning

Bayesian deep learning achieves outstanding success in various machine learning tasks like computer vision [35], physics-informed modeling [101], and complex dynamical system control [94]. From previous works, the main advantages of Bayesian deep learning come from two parts. First, the Bayesian framework can unify uncertainty quantification in deep learning. This includes the uncertainty from the neural network parameters, task-specific parameters, and exchanging information [1, 90], which is applicable to computer vision [49], spatiotemporal forecasting [99], weather forecasting [89], and so on [108]. Second, the Bayesian framework allows theoretically grounded ensemble neural networks via model averaging [96]. Bayesian neural networks apply weights distribution to neurons [8]. It is important to avoid computational requirements from Bayesian

when applying it to deep learning. To perform Bayesian sampling in deep learning, Welling and Teh establish Stochastic Gradient Langevin Dynamics as an MCMC sampler in mini-batch settings [95]. Neal introduces MCMC via Hamiltonian dynamics [66], and Chen et al. extend into deep learning via Stochastic gradient Hamiltonian Monte Carlo (SGHMC) [20, 59]. Other techniques include Nosé–Hoover thermostat [27], replica-exchange SGHMC [23], cyclic SGLD [105], contour SGLD [24], preconditioned SGLD [92], and adaptively weighted SGLD [25]. Variational methods could also help to perform approximate inference for Bayesian deep learning [8, 37].

Sparse deep learning is of emerging interest due to its lowered computational cost and improved interpretability. The discussion on sparse neural networks dates back decades [28, 65]. Glorot et al. use Rectifier Neurons with sparsity constraints to obtain sparse representations [32]. Liu et al. extend to Convolutional neural networks under sparse settings [54]. Then, several follow-up works aimed to train sparse neural networks efficiently [58, 83]. By adapting Bayesian sparse inference methods, the sparse inference process could be accelerated with valid uncertainty quantification via various prior options [26, 85, 93]. Variational methods are also applied for efficient Bayesian inference [4]. There have also been theoretical discussion of Bayesian sparse deep learning [68, 84, 86]. Overall, sparse deep learning has the potential to make computationally tractable the construction of a Bayesian model.

3 Bayesian SINDy Autoencoders

This section presents the SINDy autoencoder with a Bayesian learning process that incorporates sparsifying priors and posterior sampling. We first introduce the Sparse identification of nonlinear dynamics (SINDy) autoencoder in Sec. 3.1. Then, we propose a Bayesian learning framework that includes SINDy autoencoder in the likelihood function and specifies various setups for different sparsifying priors. Finally, we discuss Stochastic Gradient Langevin Dynamics (SGLD) to generate posterior samples from the Bayesian learning model.

3.1 Likelihood setting for Bayesian SINDy Autoencoders

The SINDy autoencoder enables a joint discovery of sparse dynamical models and coordinates. Figure 1 (a) provides an overview of an autoencoder. The input data $x(t) \in \mathbb{R}^d$ is mapped by an encoder function $f_{\theta_1}(\cdot)$ to a latent space $z(t) \in \mathbb{R}^{d_z}$, $d_z < d$. This latent space $z(t)$ contains sufficient information to recover $x(t)$ via a decoder function $g_{\theta_2}(\cdot)$.

SINDy autoencoder combines SINDy with autoencoders by constraining the latent space governed by a sparse dynamical system. The encoder function $f_{\theta_1}(\cdot)$ performs coordinate transformation to map the high-dimensional inputs into an appropriate latent subspace. The latent space $z(t) = f_{\theta_1}(x(t))$ has an associated sparse dynamical model governed by

$$\frac{d}{dt}z(t) = \Phi(z(t)) = \Theta(z(t))\Xi, \quad (2)$$

where $\Theta(z) = [\theta_1(z), \theta_2(z), \dots, \theta_p(z)]$ is a library of candidate basis functions, and a set of coefficients $\Xi = [\xi_1, \xi_2, \dots, \xi_p]$.

A statistical understanding of the model formulates the SINDy autoencoder as a parametric model \mathcal{M}_θ where $\theta = \{\theta_1, \theta_2, \Xi\}$ contains all parameters. The likelihood of this model $p(\mathcal{D}|\theta)$ is defined as

$$p(\mathcal{D}|\theta) \propto \exp\left(-\|\mathbf{x} - g_{\theta_2}(\mathbf{z})\|_2^2 + \lambda_1 \|\dot{\mathbf{x}} - (\nabla_{\mathbf{z}}g_{\theta_2}(\mathbf{z})) (\Theta(\mathbf{z}^T)\Xi)\|_2^2 + \lambda_2 \|(\nabla_{\mathbf{x}}\mathbf{z})\dot{\mathbf{x}} - \Theta(\mathbf{z}^T)\Xi\|_2^2\right). \quad (3)$$

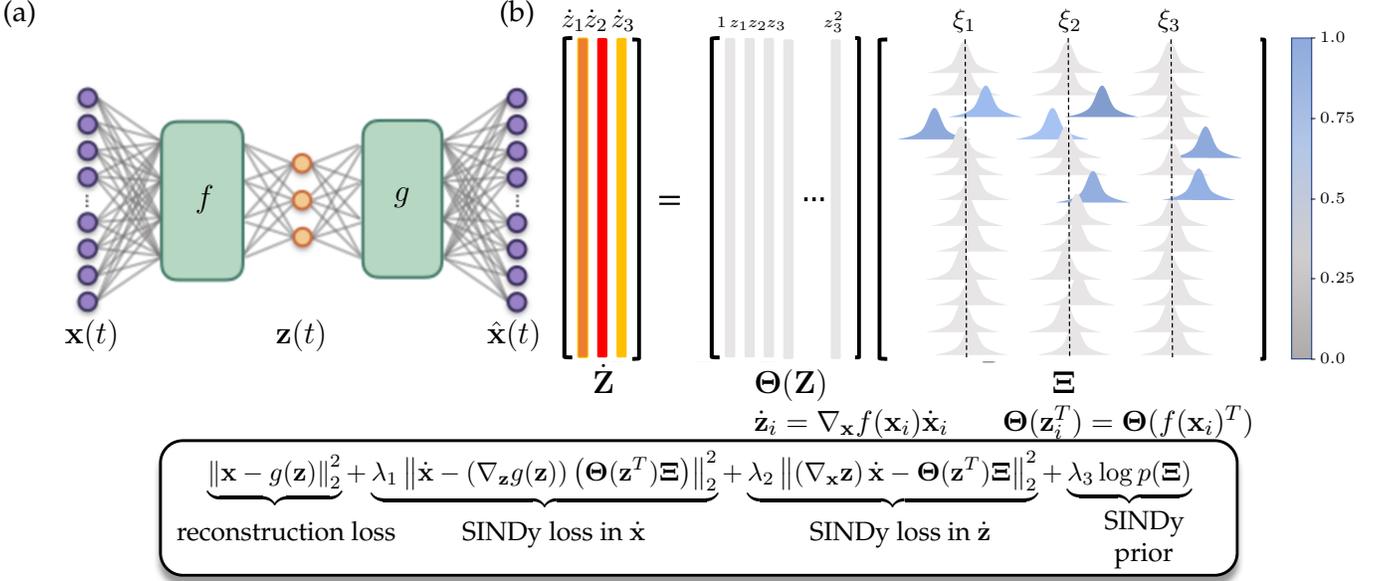


Figure 1: General structure of the Bayesian SINDy autoencoder architecture. (a) An autoencoder architecture is used to discover intrinsic coordinates \mathbf{z} from high-dimensional input data \mathbf{x} . The encoder $f(\mathbf{x})$ transforms the high-dimensional input into the a low-dimensional space \mathbf{z} , and the decoder $g(\mathbf{z})$ reconstructs \mathbf{x} from the low-dimensional subspace. (b) A Bayesian SINDy model infers the underlying dynamics of the latent dimension with uncertainty estimation. The active terms are identified by Ξ . The color represents the inclusion probability $[0.0, 1.0]$ from grey (non-active) to blue (active). The loss function encourages the network to minimize both the reconstruction error and the SINDy loss in \mathbf{z} and \mathbf{x} . The Bayesian prior works as a regularization term on Ξ for sparse inference.

The log-likelihood of this statistical model is similar to the setting in [15]. In order to promote sparsity on Ξ , we consider two sets of priors in what follows.

Laplace prior. The Laplace prior can be understood as a Bayesian LASSO [67, 88]. We define the Laplace prior such that

$$\Xi_j \sim \mathcal{L}(0, v_0), \quad (4)$$

where $\mathcal{L}(\cdot, \cdot)$ denotes the Laplace distribution defined as $f(\Xi_i; 0, v_0) = \frac{1}{2v_0} \exp\left(-\frac{|\Xi_i|}{v_0}\right)$. We can see the equivalence of the Laplace prior and the LASSO in the negative log-likelihood, where $\frac{1}{v_0} \|\Xi_i\|_1$ is included as a regularizer.

Spike-and-slab Gaussian-Laplace (SSGL) prior. We define the SSGL prior as

$$\Xi_j | \gamma_j \sim (1 - \gamma_j) \mathcal{L}(0, \sigma v_0) + \gamma_j \mathcal{N}(0, \sigma^2 v_1), \quad (5)$$

where γ_j is a binary variable, $\Xi \in \mathbb{R}^p$, $\sigma, v_0, v_1 \in \mathbb{R}$, $\mathcal{L}(\cdot, \cdot)$ denotes a Laplace distribution and $\mathcal{N}(\cdot, \cdot)$ denotes a Normal distribution. We assign a Bernoulli prior to $\gamma \sim \text{Ber}(\delta)$, $\delta \in [0, 1]$. The prior for θ_1, θ_2 is specified with a Gaussian, which is equivalent to an ℓ_2 regularization implementation-wise. For simplicity, we set δ, σ, v_0, v_1 as tunable hyperparameters.

Prior selection. In general, there is no optimal Bayesian prior for all statistical models since every prior has its own advantages and drawbacks. Therefore, the prior selection typically requires an assessment of both theory and experimental outcomes. If one selects the prior of Ξ to be a Laplace distribution, the setting will be identical to the basic SINDy autoencoder model, which is equivalent to adding a ℓ_1 regularization term (c.f. Eqn. (7) and Figure 1 in [15]). Even if the Laplace prior has a benefit in computation, its performance suffers for small sample sizes and large observation noises. Different from the Laplace prior, if one selects the SSGL prior for Ξ , it typically requires a slightly increased cost in computation. However, the SSGL prior typically has dominating performances for cases with very large noise and limited sample size, which is preferred in our case to learn the actual video data.

Bayesian formulation. Using the Bayes formula, we can construct the posterior distribution from the likelihood function and prior that

$$\pi(\theta, \gamma | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta_1) p(\theta_2) p(\Xi | \gamma) p(\gamma). \quad (6)$$

We aim to depict the posterior distribution $p(\Xi | \mathcal{D})$ via the joint distribution $p(\theta | \mathcal{D})$ under the sparsifying SSGL prior. The approximation of $p(\Xi | \mathcal{D})$ is accessible from posterior samples of $p(\theta | \mathcal{D})$ when dropping θ_1, θ_2 . In the setting of deep neural networks, we can sample the posterior distribution via Stochastic Gradient methods using mini-batches. The mini-batch setting not only naturally fits into the training of deep neural networks but also could accelerate the Bayesian posterior sampling process.

3.2 Stochastic Gradient Langevin Dynamics

To perform posterior sampling in mini-batch settings, Stochastic Gradient Langevin Dynamics (SGLD) is a popular method that combines stochastic optimization and Langevin Dynamics [95]. Denote the learning rate at epoch t by $\epsilon^{(t)}$ which decreases to zero, and the dataset $\mathcal{D} = \{d_i\}_{i=1}^N$. The mini-batch setting estimates the gradient $\nabla_{\theta} L(\theta)$ from a subset (batch) $\mathcal{B} = \{d_i\}_{i=1}^n$. The injected noise from mini-batches facilitates the generation of posterior samples while reducing the high computational cost in computation for full-batch gradients. We follow the classical SGLD setting where

$$\Delta\theta_{t+1} = \frac{\epsilon^{(t)}}{2} \left(\nabla \log p(\theta_t) + \frac{N}{n} \sum_{i=1}^n \nabla \log p(X_i | \theta_t) \right) + \eta_t, \quad \eta_t \sim \mathcal{N}(0, \epsilon^{(t)}). \quad (7)$$

Here, $p(\theta_t)$ denotes the prior specified in Eqn. (4), (5) and $p(X_i | \theta_t)$ denotes the data likelihood. Prior works study the asymptotic convergence of SGLD to the target distribution, which validates SGLD for posterior sampling in theory [41, 87, 106]. An advantageous property of SGLD in posterior sampling is that with decaying step size $\epsilon^{(t)}$, SGLD automatically transfers from a stochastic optimization algorithm to a posterior sampling procedure (c.f. Sec. 4.1 [95]). The Metropolis-Hasting correction can be ignored since the rejection rate for sampling goes to zero asymptotically, resulting from $\epsilon^{(t)} \rightarrow 0$ when $t \rightarrow \infty$. The discretization error similarly decreases as ϵ goes to zero.

Cyclical SGLD. The cyclical SGLD method [105] consists of exploration and sampling stages via a cyclical step-size schedule for $\epsilon^{(t)}$. In the training of deep autoencoders, the optimization process is highly non-convex with a very complex loss landscape. The cyclical step-size schedule helps to explore the parameter space when ϵ is large as well as sample local mode when ϵ is small. It is also possible to understand cyclical SGLD from a parallel SGLD perspective [23]. We apply this idea in our experiments to perform better inference.

3.3 Empirical Bayes Variable Selection in SINDy Autoencoder

The empirical bayesian method infers prior hyperparameters from data. In this case, we aim to optimize γ (ignoring the uncertainty) and sample $\theta|\mathcal{D}$. The posterior distribution of $\Xi|\mathcal{D}$ can be derived from posterior samples of $\theta|\mathcal{D}$.

Using the mini-batch setting, the posterior distribution follows

$$\pi(\theta, \gamma|\mathcal{B}) \propto p(\mathcal{B}|\theta) \frac{N}{n} p(\theta_1) p(\theta_2) p(\Xi|\gamma) p(\gamma). \quad (8)$$

The term $p(\mathcal{B}|\theta)$ can be evaluated from the loss of SINDy autnecoder for the current mini-batch \mathcal{B} ; the term $p(\theta_1), p(\theta_2)$ can be computed from the ℓ_2 loss; the term $p(\Xi|\gamma)$ can be computed from Eqn. (5); and the term $p(\gamma)$ can be known from the Bernoulli prior setting.

As a binary variable, γ is difficult to be optimized for due to non-continuity and non-convexity. An important trick to perform the optimization on γ is to alternatively optimize the adaptive posterior mean $\mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{D}} [\pi(\theta, \gamma|\mathcal{D})]$ which treats γ as a latent variable. Following previous works with similar settings [23, 71], we could indirectly evaluate $\pi(\theta, \gamma|\mathcal{D})$ by a strict lower bound $Q(\cdot|\cdot)$ that

$$Q(\theta|\theta^{(k)}) = \mathbb{E}_{\mathcal{B}} \left[\mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{D}} [\log \pi(\theta, \gamma|\mathcal{B})] \right] \leq \log \mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{D}} [\mathbb{E}_{\mathcal{B}} [\pi(\theta, \gamma|\mathcal{B})]]. \quad (9)$$

The inequality holds by Fubini's theorem and Jensen's inequality (c.f. Eqn. (7) in [23]).

The variable $Q(\theta|\theta^{(k)})$ can be decomposed into

$$Q(\theta|\theta^{(k)}) = \frac{N}{n} \log \pi(\mathcal{B}|\theta) - \|\theta_1\| - \|\theta_2\| - \sum_{i \in |\Xi|} \left[\frac{|\Xi_i| \kappa_{i0}}{\sigma} + \frac{\Xi_i^2 \kappa_{i1}}{2\sigma^2} \right] + \sum_{i \in |\Xi|} \log \left(\frac{\delta}{1-\delta} \right) \rho_i + C, \quad (10)$$

where $\kappa_{i0} = \mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{D}} \left[\frac{1}{v_0(1-\gamma_i)} \right]$, $\kappa_{i1} = \mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{D}} \left[\frac{1}{\gamma_i} \right]$, $\rho_i = \mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{D}} [\gamma_i]$, and $C \in \mathbb{R}$ is a constant. Notice here γ is treated as a latent variable, and we only consider the expectation given the conditional distribution of γ given $\theta^{(k)}, \mathcal{D}$. The ρ_i could be considered as an inclusion probability estimate of Ξ_i . In this way, ρ softens γ from a binary variable into a continuous one.

The term κ performs an elastic net-like approach which adaptively optimizes the ℓ_1 and ℓ_2 coefficients. Suppose Ξ_i is identified with a high probability to be a sparse variable (e.g., $\rho_i < 0.05$), the term κ_{i0} will be large, strengthening the sparsity constraints. Otherwise, if Ξ_i is identified as a non-sparse variable (e.g. $\rho_i > 0.95$), the term κ_{i1} will be small, and the ℓ_2 constraint will be dominated instead.

3.3.1 Stochastic Approximation from Expectation Maximization of ρ and κ

We could derive an asymptotically correct posterior distribution on $\pi(\theta, \kappa, \rho)$ following the steps in [23]:

1. Sample θ from $Q(\cdot)$ that

$$\theta^{(k+1)} = \theta^{(k)} + \eta^{(k)} \nabla_{\theta} Q(\rho^{(k)}, \kappa^{(k)}|\mathcal{B}^{(k)}) + \mathcal{N}(0, 2\eta^{(k)}) \quad (11)$$

2. Perform stochastic approximation to latent variables ρ, κ that

$$\rho^{(k+1)} = (1 - \omega^{(k+1)})\rho^{(k)} + \omega^{(k+1)}\tilde{\rho}^{(k+1)}, \quad (12)$$

$$\kappa^{(k+1)} = (1 - \omega^{(k+1)})\kappa^{(k)} + \omega^{(k+1)}\tilde{\kappa}^{(k+1)}. \quad (13)$$

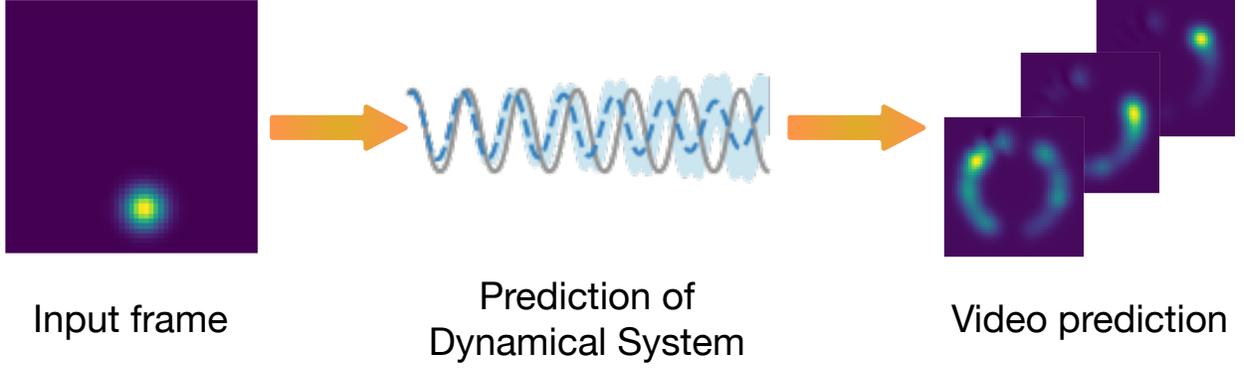


Figure 2: Add more things to this figure including: encoder (with its shape), decoder (with its shape), and Bayesian SINDy.

Here, $\tilde{\rho}_i^{(k+1)}$, $\tilde{\kappa}_i^{(k+1)}$ is the Expectation Maximization (EM) estimation of ρ , κ given iteration k [71]. The EM estimation of the inclusion probability ρ_i is

$$\tilde{\rho}_i^{(k+1)} = \mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{B}}[\gamma_i] = P(\gamma_i = 1 | \Xi^{(k)}) = \frac{a_i}{a_i + b_i}, \quad (14)$$

where $a_i = \pi(\Xi_i^{(k)} | \gamma_i = 1)P(\gamma_i = 1 | \delta)$ and $b_i = \pi(\Xi_i^{(k)} | \gamma_i = 0)P(\gamma_i = 0 | \delta)$. In our case, the term $\pi(\Xi_i^{(k)} | \gamma_i = 1)$ is the probability of $\Xi_i^{(k)}$ given the Gaussian prior distribution, and the term $\pi(\Xi_i^{(k)} | \gamma_i = 0)$ is the probability of $\Xi_i^{(k)}$ from the Laplace distribution. The latter term $P(\gamma_i = 1 | \delta) = \delta$ given Bernoulli prior in the previous setting. Following a similar process, for κ , we have

$$\tilde{\kappa}_{i0}^{(k+1)} = \mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{B}} \left[\frac{1}{v_0(1 - \gamma_i)} \right] = \frac{1 - \rho_i}{v_0}, \quad (15)$$

$$\tilde{\kappa}_{i1}^{(k+1)} = \mathbb{E}_{\gamma|\theta^{(k)}, \mathcal{B}} \left[\frac{1}{v_1 \gamma_i} \right] = \frac{\rho_i}{v_1}. \quad (16)$$

3.4 Prediction with Bayesian SINDy Autoencoder

Bayesian SINDy autoencoder has a trustworthy application in our application of video prediction whereby we learn the dynamics and coordinates of the latent space. Based on an inference process established in previous sections, the video prediction can precisely understand the underlying dynamical system, which allows accurate, robust, and interpretable future forecasting. Additionally, the Bayesian framework enables precise uncertainty quantification for video prediction from posterior samples. We visualize the prediction process in Fig. 2. In this case of a pendulum video, the uncertainty of video prediction grows with larger t , and gradually fails to predict by showing a ring-like prediction.

From the procedure in Sec. 3.3, we can generate samples from $p(\theta | \mathcal{D})$, which is the posterior distribution of neural network parameters θ given observed data \mathcal{D} . Suppose the posterior samples are $\Xi^1, \Xi^2, \dots, \Xi^m \sim p(\theta | \mathcal{D})$. The full posterior predictive distribution is defined as

$$\begin{aligned} p(\hat{x}(t) | x_0, \mathcal{D}) &= \int p(\hat{x}(t) | \theta, x_0) p(\theta | \mathcal{D}) d\theta \\ &= \int p(\hat{x}(t) | \theta_2, \hat{z}(t)) p(\hat{z}(t) | \Xi, \hat{z}(0)) p(\hat{z}(0) | \theta_1, x_0) p(\theta_1, \theta_2, \Xi | \mathcal{D}) d\theta \end{aligned} \quad (17)$$

Therefore, we can approximately generate samples from $p(\hat{x}(t)|x_0, \mathcal{D})$ using Monte Carlo estimation. For simplicity, we only consider the Maximum likelihood Estimation $\theta_1^{\text{MLE}}, \theta_2^{\text{MLE}}$. We have the following process:

1. From the input image x_0 , compute $z(0) = f_{\theta_1^{\text{MLE}}}(x_0)$.
2. From $z(0)$, using posterior samples of Ξ , generate $\hat{z}^{(i)}(t)$ samples via

$$z^{(i)}(t) = z(0) + \int_0^t \Theta(z(t'))\Xi^{(i)} dt'.$$

3. From samples of $z^{(i)}(t)$, generate $\hat{x}^{(i)}(t) = g_{\theta_2^{\text{MLE}}}(z^{(i)}(t))$.

4 Experiments

In the following subsections, we conduct four case studies on Bayesian SINDy autoencoder for governing equations and coordinate system discovery for video data. In Sec. 4.1, we study similar cases to those in [15] with synthetic high-dimensional data generated from the Lorenz system, reaction-diffusion, and a single pendulum. We explore the Laplace prior for chaotic Lorenz system and the SSGL prior for reaction-diffusion and single pendulum. In Sec. 4.2, we study real video data that consists of 390 video frames of a moving rod. This setting is particularly challenging due to the high dimensionality in video data, noisy observation, missing temporal derivatives, and prior setting for correct learning. In both experiments on synthetic and real video data, we observe the Bayesian SINDy autoencoder can accurately perform nonlinear identification of dynamical system under the correct setting of preprocessing, prior, and training parameters. All the experiments are implemented and run on a single NVIDIA GeForce RTX 2080 Ti.

4.1 Learning Physics from Synthetic Video Data

4.1.1 Chaotic Lorenz System via the Laplace prior

We consider the chaotic Lorenz system in the following

$$\begin{cases} \dot{z}_1 = -\sigma z_1 + \sigma z_2 \\ \dot{z}_2 = \rho z_1 - z_2 - z_1 z_3 \\ \dot{z}_3 = z_1 z_2 - \beta z_3, \end{cases} \quad (18)$$

where $z = [z_1, z_2, z_3] \in \mathbb{R}^3$ and σ, ρ, β are constants. The Lorenz system is very representative of the chaotic and nonlinear system, which is an ideal example of applying model discovery techniques. In the numerical simulation, we first set $\sigma = 10, \rho = 28, \beta = -2.7$. We only generate partial Lorenz via time range $t = [0, 5]$ with $\Delta t = 0.02$ for 1024 different Lorenz systems from random initial conditions. The initial condition follows a uniform distribution centered at $[0, 0, 25]$ with width $[36, 48, 41]$ respectively.

From the underlying dynamical system, we create a high-dimensional dataset via six fixed spatial models given by Legendre polynomials that $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_6 \in \mathbb{R}^{128}$. We transfer from the low-dimensional dynamical system into a high-dimensional dataset via the following rule:

$$\mathbf{x}(t) = \mathbf{u}_1 z_1(t) + \mathbf{u}_2 z_2(t) + \mathbf{u}_3 z_3(t) + \mathbf{u}_4 z_1(t)^3 + \mathbf{u}_5 z_2(t)^2 + \mathbf{u}_6 z_3(t)^3. \quad (19)$$

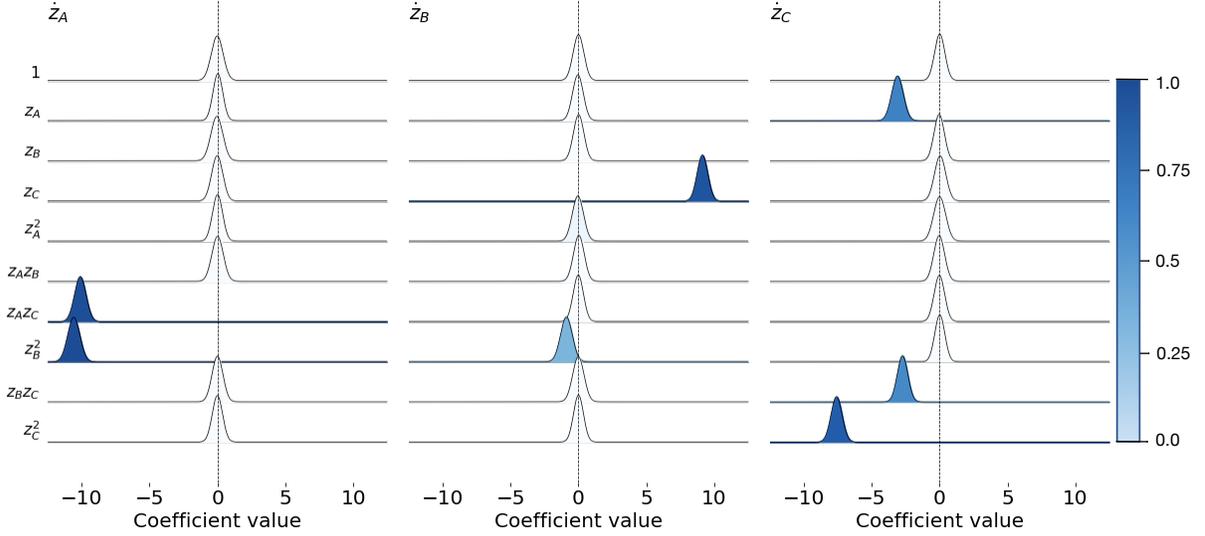


Figure 3: Bayesian estimation and uncertainty quantification visualization of SINDy coefficient for Lorenz system under Laplace prior. The color bar represents the inclusion probability estimate given the coefficient magnitude.

We set the autoencoder with latent dimension $d = 3$ corresponding to the latent system in the 3D coordinate system with z_A, z_B, z_C . We include polynomials with the highest order 3 composing a library $[1, z_A, z_B, z_C, z_A^2, z_A z_B, z_A z_C, z_B^2, z_B z_C, z_C^2, z_A^3, z_A^2 z_B, z_A^2 z_C, \dots, z_C^3]$. Via the autoencoder, we wish to identify the correct active terms as well as the value of the coefficients. The coefficient of Ξ is uniformly initialized from constant 1. The loss coefficients are $\lambda_1 = 0.0, \lambda_2 = 1 \times 10^{-4}$. For the encoder and decoder, we use the sigmoid activation function with widths $[64, 32]$. For optimization, we select Adam optimizer with learning rate 1×10^{-3} and the batch size to be 1024. For the Laplace prior setting, we set $\lambda_3 = 1 \times 10^{-5}$.

By training with 5,000 epochs following the setting from [15], in terms of the error metrics, the best test error of the decoder reconstruction achieves 2×10^{-5} of the fraction of the variance from the input. The fraction of unexplained variances are 2×10^{-4} for the reconstruction of \dot{z} , and 1.3×10^{-3} for the reconstruction of \dot{x} . Notice here decoder reconstruction is better compared to SINDy autoencoder without uncertainty quantification, but the reconstruction of \dot{z} is slightly worse compared to point estimation for SINDy autoencoder.

The coefficient estimate from the Bayesian SINDy autoencoder is shown in Fig. 3. From the figure, we can observe the uncertainty quantification in the parameter space. The outcome of this model identifies 7 active terms marked with deeper blues. It is known that the identification of Lorenz dynamics suffers from the symmetry in the coordinate system as described in [15]. The discovered governing equation could be equivalently transformed via the affine group transformation on the coefficients and a permutation group transformation on the latent variables. Therefore, the identification is still correct in Fig. 3 from the following transformations. (a) We inversely transform the permutation group via assigning z_C to z_1, z_B to z_2 , and z_A to z_3 . (b) We inversely perform affine transformation by $z_1 = 1.0, z_2 = -0.94z_2$, and $z_3 = 0.55z_3 - 2.81$. We demonstrate the effectiveness of the affine transformation process in Fig. 4 (c) that the discovered model (c.1) can be transformed

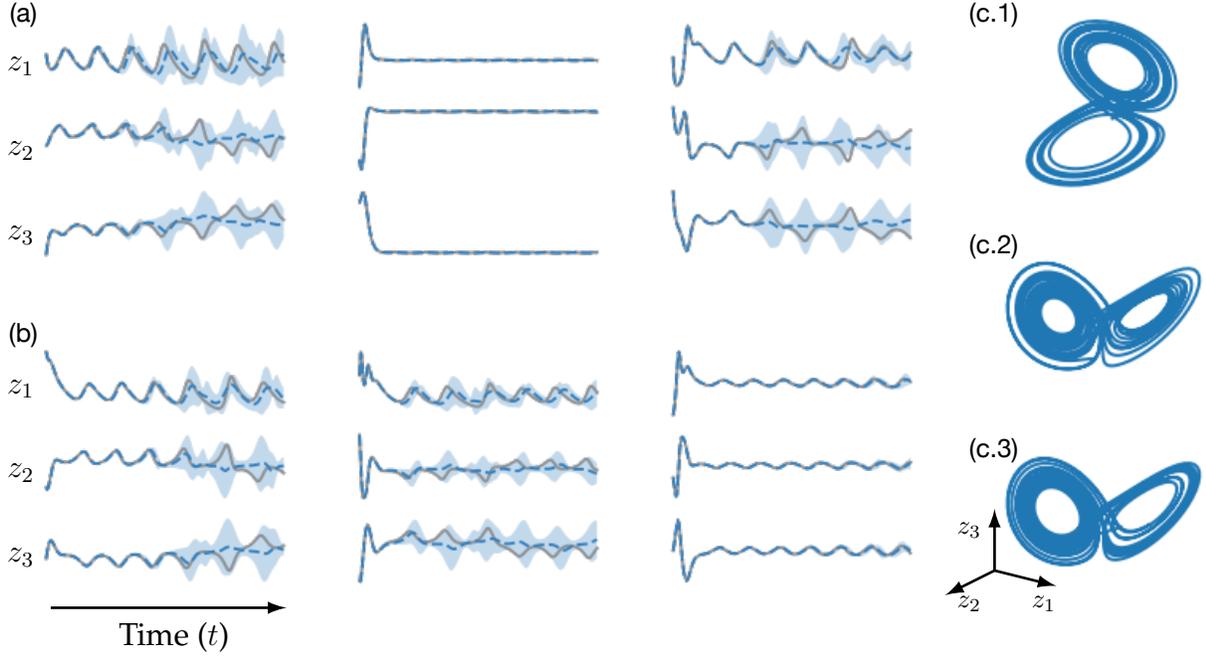


Figure 4: (a) Trajectory prediction for **in-distribution** data with uncertainty quantification generated by posterior samples from Bayesian inference on SINDy coefficient. (b) Trajectory prediction for **out-of-distribution** data with uncertainty quantification generated by posterior samples from Bayesian inference on SINDy coefficient. (c.1) Discovered Lorenz systems and (c.2) transformation to the standard space. (c.3) Lorenz system generated from the ground truth.

into (c.2), which is close to the ground truth model (c.3). The discovered governing equation from the Laplace prior is

$$\begin{cases} \dot{z}_1 = -10.09z_1 + 10.00z_2 \\ \dot{z}_2 = 27.09z_1 - 0.86z_2 - 5.35z_1z_3 \\ \dot{z}_3 = 5.35z_1z_2 - 2.71z_3, \end{cases} \quad (20)$$

Stepping further from the uncertainty quantification of coefficients, we can see the uncertainty quantification in the prediction space in Fig. 4 (a) (b). In general, the uncertainty coverage correctly captures the true trajectory when the model fails to predict correctly and mildly covers the trajectory when the model prediction is confident.

4.1.2 Reaction Diffusion via Spike-and-slab prior

Reaction-diffusion is governed by a partial differential equation (PDE) that has complex interactions between spatial and temporal dynamics. We define a lambda-omega reaction-diffusion system by

$$\begin{aligned} u_t &= (1 - (u^2 + v^2))u + \beta(u^2 + v^2)v + d_1(u_{xx} + u_{yy}) \\ v_t &= -\beta(u^2 + v^2)u + (1 - (u^2 + v^2))v + d_2(v_{xx} + v_{yy}), \end{aligned} \quad (21)$$

which $d_1 = d_2 = 0.1, \beta = 1$.

For the synthetic data, we first generate the latent dimensions $u(x, y, t)$ and $v(x, y, t)$ by 21 from $t = [0, 500]$ with time step $\Delta t = 0.05$. Then, we generate the snapshots of the dynamical

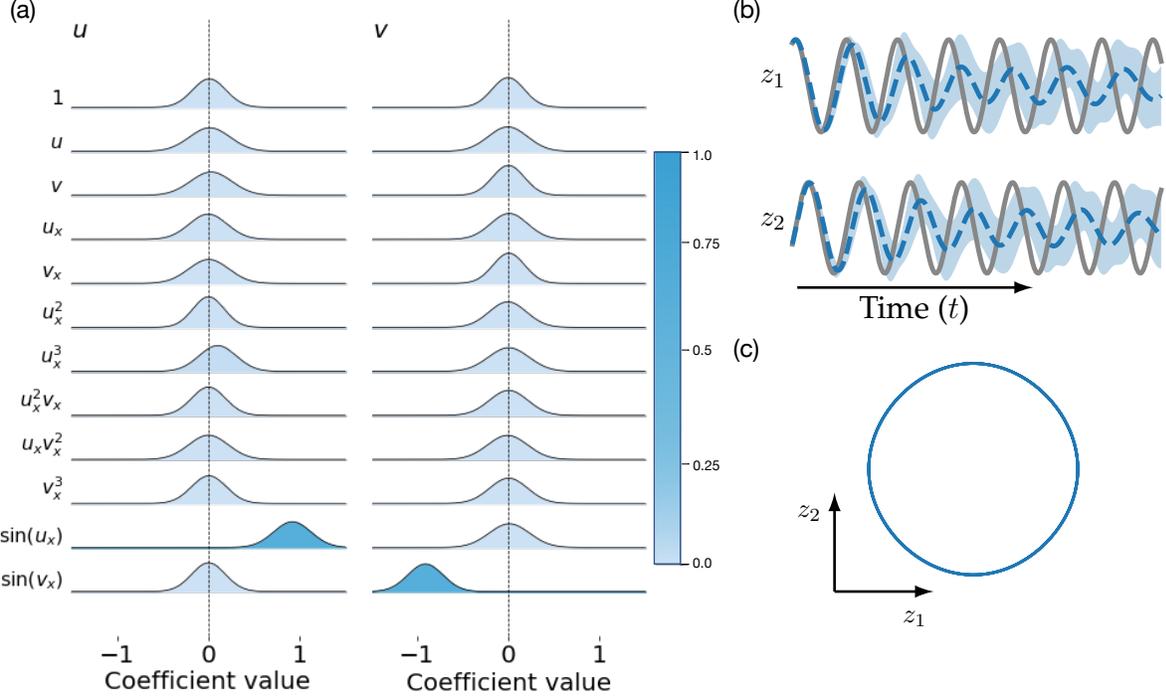


Figure 5: (a) Bayesian estimation and uncertainty quantification visualization of SINDy coefficient for reaction-diffusion under Spike-and-slab prior. (b) Visualization of in-distribution (top) and out-of-distribution (bottom) predicted dynamics. (c) The generated attractor from the mean of Bayesian inference.

system spatially in xy -domain from its latent dimensions into a video with shape (10000, 100, 100). We obtain a training dataset with size 9,000. We also generate a testing dataset with 1,000 samples.

We setup the autoencoder with latent dimension $d = 2$, targeting the u , v , and apply a first-order library of functions including $[1, u_x, v_x, u_x^2, u_x v_x, v_x^2, u_x^3, u_x^2 v_x, u_x v_x^2, v_x^3, \sin(u_x), \sin(v_x)]$. We hope to discover the two oscillating spatial modes for this nonlinear oscillatory dynamics. The coefficient of Ξ is randomly initialized from Gaussian $\mathcal{N}(0, 0.1)$. The loss coefficients are $\lambda_1 = 1.0 \times 10^{-2}$, $\lambda_2 = 1.0 \times 10^{-1}$, $\lambda_3 = 20.0$. For the encoder and decoder, we use sigmoid activation function with width [256]. For optimization, we select Adam optimizer with learning rate $1e^{-3}$, and batch size to be 1000. For the setting of SSSL prior, we set $\delta = 0.08$, $v_0 = 0.1$, $v_1 = 3.0$, $\omega^{(k)} = 0.02 \times (0.999)^k$.

We follow the same setting in [15]. The SSSL prior only requires 1,500 epochs for training while LASSO setting frequently needs more than 3,000 epochs for training. The best test error is 2.1×10^{-3} for decoder loss, 2.1×10^{-5} for the reconstruction of \dot{z} , and 1.7×10^{-4} for the reconstruction of \dot{x} . The fraction of unexplained variance of decoder reconstruction is 9.1×10^{-5} . For SINDy predictions, the fraction of unexplained variance of \dot{x} is 0.013. The fraction of unexplained variance of \dot{z} is 0.001. These results all improve from LASSO based SINDy Autoencoder [15]. The posterior samples from Bayesian SINDy are shown in Fig. 5. The discovered governing equation from the Bayesian SINDy autoencoder is

$$\begin{cases} \dot{u}_x = 0.91 \sin(u_x) \\ \dot{z}_2 = -0.91 \sin(v_x) \end{cases} \quad (22)$$

4.1.3 Nonlinear Pendulum via the Spike-and-slab prior

We consider simulated video of a nonlinear pendulum in pixel space with two spatial dimensions. The nonlinear pendulum is governed by the following second-order differential equation:

$$\ddot{z} = -\sin z. \quad (23)$$

We generate the synthetic dataset following the settings in [15]. The synthetic data first generates the latent dimension z as the angle of pendulum from $t = [0, 10]$ with time step $\Delta t = 0.02$. Then, we form a Gaussian ball around the mass point with angle z and length l . This process transfers the dynamical system from its latent dimension z into a video with shape $(500, 51, 51)$. By simulating this process with 100 times, we obtain a training dataset with size 50,000. We also generate a testing dataset with 5,000 samples.

We set up the autoencoder with latent dimension $d = 1$, targeting the angle of pendulum, and apply a second-order library of functions including $[1, z, \dot{z}, z^2, z\dot{z}, \dot{z}^2, z^3, z^2\dot{z}, z\dot{z}^2, \dot{z}^3, \sin(z), \sin(\dot{z})]$. We hope to infer the real coefficient that $\Xi = [0, 0, 0, 0, 0, 0, 0, 0, 0, -1, 0]$. The coefficient of Ξ is randomly initialized from Gaussian. The loss coefficients are $\lambda_1 = 5e^{-3}$, $\lambda_2 = 5e^{-5}$, $\lambda_3 = 8e^{-4}$. For the encoder and decoder, we use the sigmoid activation function with widths $[128, 64, 32]$. For optimization, we select Adam optimizer with a learning rate $1e^{-3}$, and batch size to be 1000. For the setting of SSGL prior, we set $\delta = 0.08$, $v_0 = 0.05$, $v_1 = 3.0$, $\omega^{(k)} = 0.05 \times (0.995)^k$.

We follow the same setting in [15]. The success discovery rate is 80% from 15 training instances, which improves from 50% success discovery rate in the LASSO setting. Additionally, Bayesian SINDy Autoencoder with the SSGL prior requires only 1,500 epochs for training, while the LASSO setting frequently needs 5,000 epochs for training. The best test error is 6.5×10^{-8} for decoder loss, 1.1×10^{-3} for the reconstruction of \ddot{x} , and 8.1×10^{-4} for the reconstruction of \ddot{z} . The best fraction of unexplained variance of decoder reconstruction is 4.5×10^{-4} . For SINDy predictions, the best fraction of unexplained variance of \ddot{x} and \ddot{z} reconstruction are 2.3×10^{-4} and 5.5×10^{-3} . The posterior samples from Bayesian SINDy are shown in Fig. 6 (a). The discovered governing equation from the Bayesian SINDy autoencoder is

$$\ddot{z} = -0.99 \sin(z). \quad (24)$$

The prediction of trajectory with uncertainty quantification is shown in Fig. 6 (b) with 6 different initial conditions that $\theta_0 = [\frac{2}{3}\pi, \frac{1}{2}\pi, \frac{1}{3}\pi, \frac{1}{4}\pi, \frac{1}{8}\pi, \frac{1}{16}\pi]$ from time $t = [0, 100]$. The error and uncertainty grow with longer time interval, generally starting from $t = 20$.

4.2 Learning dynamical system from real video data

Experimental setting and dataset description. The raw video data has 14 seconds recording of a moving rod. We process the raw video via the following steps.

- (1) We first move the RGB channels of a video to greyscale, normalizing to $[0, 1]$. This is step (a) to (b) in Fig. 7.
- (2) We estimate the background in this video. In our setting, we compute an averaged frame over the entire video as the background. In a wider context, one shall remove low-rank information via linear models [53, 98]. This is step (b) to (c) in Fig. 7.
- (3) Apply Gaussian filters with appropriate variance. The selection of variance varies from case to case. A good hyperparameter setting should keep the mix-max gap of the original image and could also smoothen the sharp edges of the objects. This is step (c) to (d) in Fig. 7.

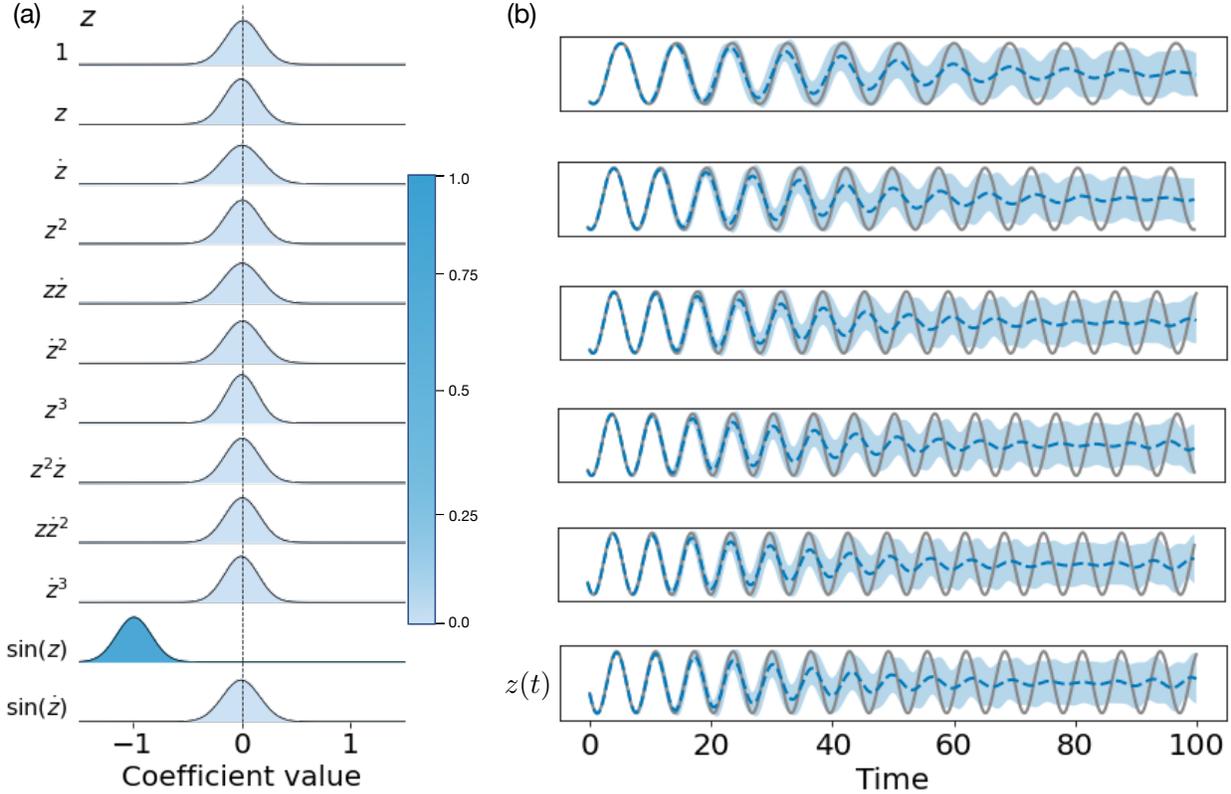


Figure 6: (a) Bayesian estimation and uncertainty quantification visualization of SINDy coefficient for pendulum data under Spike-and-slab prior. (b) Visualization of predicted dynamics with uncertainty quantification from 6 different initializations.

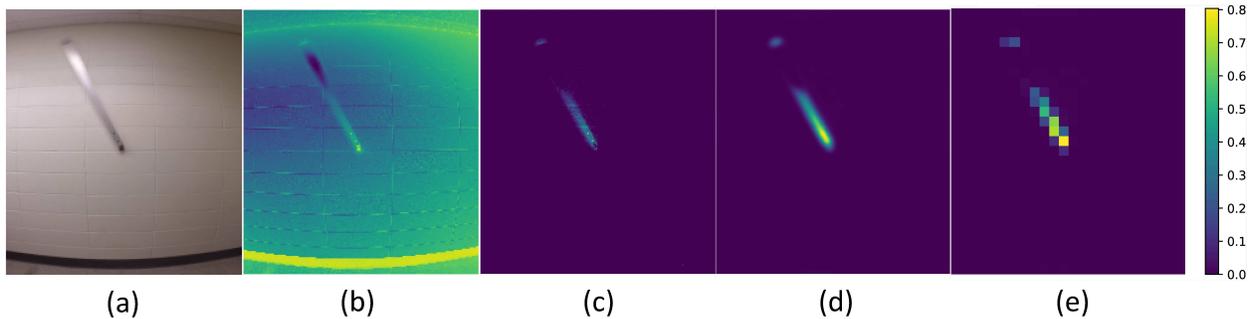


Figure 7: Preprocessing pipeline to high-dimensional video data to remove auxiliary information and reduce the size of video frames.

(4) The final step downsample the frames from (1080, 960) to (27, 24) via interpolation.

After these preprocessing steps, we obtain a training dataset with 390 samples and shape (27, 24).

Bayesian masked autoencoder setup We set up the autoencoder with latent dimension $d = 1$, with a second-order library of functions similar to synthetic video data. We remove the constant term for simplicity of the inferential process. The loss coefficients are $\lambda_1 = 5e^{-7}$, $\lambda_2 = 5e^{-8}$, and

$\lambda_3 = 1 \times 10^{-4}$. For the encoder and decoder, we use the sigmoid activation function with widths [64, 32, 16]. For optimization, we select Adam optimizer with a learning rate $1e^{-3}$, and batch size to be 10. Real video data is much harder compared to synthetic video data since it contains rich information, including object colors, background, and processing noises. At the same time, only 390 snapshot samples are available to train the Bayesian autoencoder. To resolve this problem, we involve masked autoencoder [36], which applies 50% random masking to the input data. The effect of masked autoencoder could be understood as data augmentation, which allows the encoder to focus more on the key information [100].

4.2.1 Results from the Bayesian discovery via SSSL prior

We discuss both Laplace and Spike-and-slab priors for learning real moving rod data. In this case with very low and noisy data, the Laplace prior identifies an incorrect term z over than the true target on $\sin(z)$. We provide detailed discussion of the Laplace prior in the Supplemental material (5). The SSSL prior could perform correct dynamical identification yet with a slightly large coefficient estimation. For the setting of SSSL prior, we set $\delta = 0.09$, $v_0 = 1.0$, $v_1 = 5.0$, $\omega^{(k)} = 0.01 \times (0.999)^k$. The success discovery rate is 100% from 10 training instances, which improves from 70% success discovery rate in the LASSO setting. Bayesian SINDy Autoencoder requires 1,500 epochs for training. The best test error is 0.767 for decoder loss, 0.305 for the reconstruction of \ddot{x} , and 0.0263 for the reconstruction of \ddot{z} . The fraction of the unexplained variance of decoder reconstruction is 0.003. For SINDy predictions, the fraction of unexplained variance of \ddot{x} and \ddot{z} reconstruction are 0.081 and 0.961.

The posterior samples from Bayesian SINDy are shown in Fig. 8 (a), and the prediction of trajectory with uncertainty quantification of latent dynamics is shown in the bottom figure of Fig. 8 (b). The discovered governing equation from the SSSL prior is

$$\ddot{z} = -16.06 \sin(z). \quad (25)$$

The latent dimension \ddot{z} with SINDy correctly identifies $\sin(z)$ as the active term and generates an uncertainty estimate. Utilizing the uncertainty estimation in the parameter space, we could perform uncertainty quantification on the prediction space as in Fig. 8 (b). The training outcome of Bayesian SINDy autoencoder are shown in Fig. 8, 9. First, from Fig. 8 (c), we note that the reconstruction of the Bayesian SINDy autoencoder is very promising. The top line of figures is input data, and the second line of figures is generated by the autoencoder. Even if the scale is slightly different, it is convincing that the overall reconstructions on x , \dot{x} and \ddot{x} are very good. This can be suggested by Fig. 9 that the log loss converges with more training epochs. Throughout these experiments, we uniformly set training epochs to be 1,500 and refinement epochs to be 1,500. This decision is suggested by the observation in a longer run (7,500 training epochs and 1,500 refinement epochs) that the log loss remains stable after 1000 epochs. For the reconstruction of \ddot{z} , having longer training epochs will harm the testing error due to overfitting. We note here the discovery by SSSL prior has dominant performance in testing loss comparing to SINDy autoencoder as shown in 5. This suggests the superiority of the SSSL prior in the real video data setting.

As suggested in Sec. 3.4, the Bayesian SINDy autoencoder enables a trustworthy solution in deep learning for video prediction. In Fig. 8 (b), we predict the following frames of video using the decoder and the simulation of latent dynamical systems. We observe that in the very beginning of future time frames, the prediction is very confident and accurate. Yet as we move to a longer time interval, the prediction becomes very uncertain, and we can observe this interesting phenomenon in Fig. 8 (b).

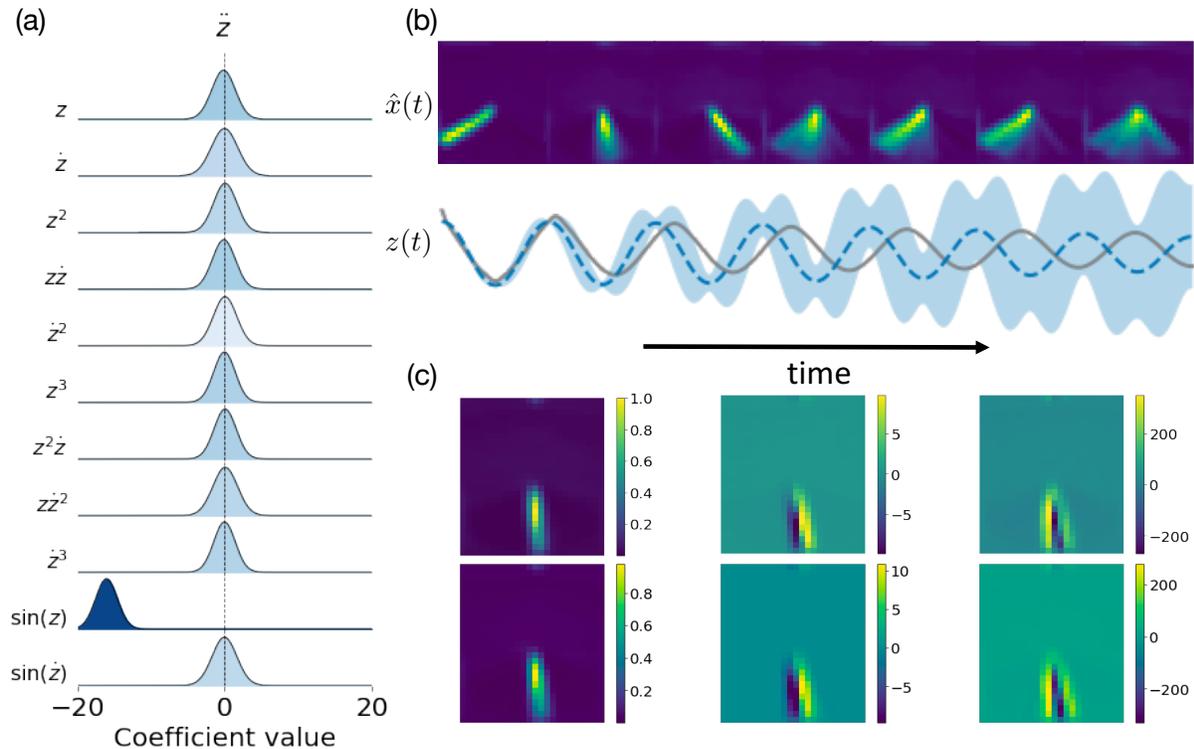


Figure 8: (a) Bayesian estimation and uncertainty quantification visualization of SINDy coefficient for real moving rod data under Spike-and-slab prior. (b) Visualization of prediction in both latent dynamics and pixel space with uncertainty quantification. (c) Reconstruction of Bayesian SINDy autoencoder for video frame, and temporal derivatives.

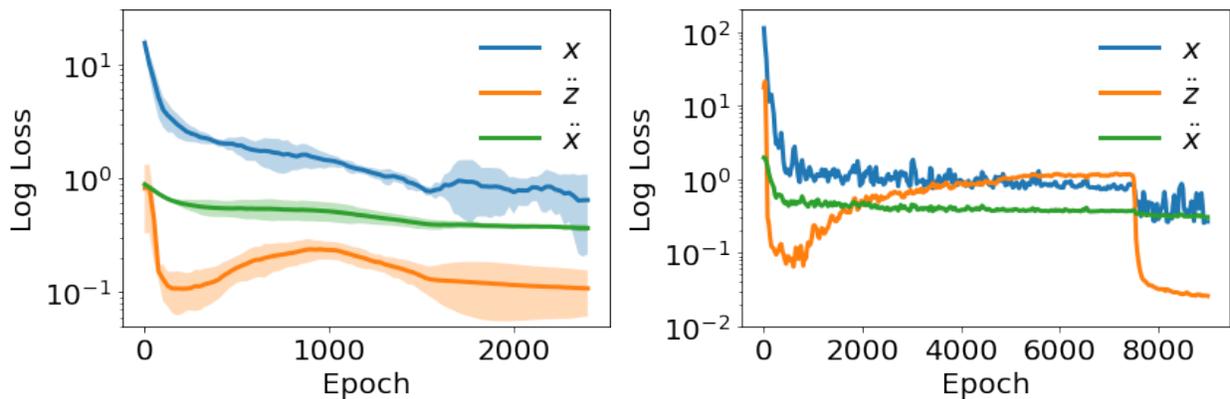


Figure 9: Left: by running the experiment 5 times with 2000 training epochs and 1500 refinement epochs, we plot the mean and standard deviation of testing loss under the log scale. Right: by running the experiment once with 7500 training epochs and 1500 refinement epochs, we plot the testing loss under the log scale.

It is essential to understand the latent dimension of the autoencoder in order to know the validity of the coordinate system discovery. Therefore, we create the following analysis on the latent dimension. We first manually label all angles θ from all video frames by humans, as shown

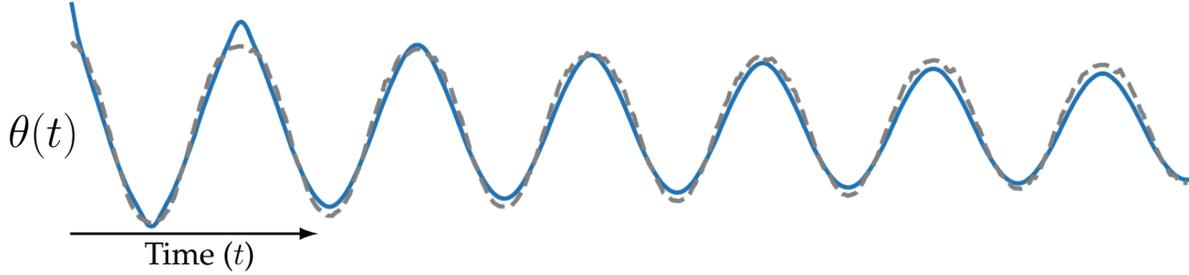


Figure 10: Latent dimension visualization after rescaling (blue curve) versus manually labelled moving rod angle by human (grey curve).

in the grey dashed curve in Fig. 10. The blue curve represents the rescaled latent dimension z , and we observe these two curves match mostly perfect to each other. We apply rescaling in this process due to the coordinate system for the angles θ are equivalent after arbitrary transformation in its scale. For example, the angles represented by radians and degrees are equivalent to each other under different scales.

On the discovery of the standard gravity constant g The SSSL prior accurately discovers the model with $\sin(z)$ term and reveals the standard gravity constant g . In our experimental setting, the length of the rod in our experiment is 123 cm with the initial condition at 75° . We have an estimate of g given these conditions that $\hat{g} = -9.876$. The estimation \hat{g} is slightly overconfident and biased which concentrates at $(-9.889, -9.865)$. The bias is more likely to be removed with an enlarged video datasets. It is also possible to use the Laplace prior discovery to estimate the standard gravity constant because the discovered dynamical system with the term z empirically has the form that $\ddot{z} = -\frac{\tilde{g}}{2L}z$ (c.f. [15] S2.3). This would result in an estimation of the standard gravity g to be $\tilde{g} = -8.733$. However, the estimation \tilde{g} is more biased, and lacks rigor in the context of the physics discovery.

Given this low-data and high-noise setting, the estimate from SSSL prior \hat{g} is very remarkable given the true value that $g \approx -9.807$. Even if the discovery of the coordinate system could be up to an arbitrary scaling, the frequency of the dynamical system could uniquely identify the coefficient of $\sin(z)$. This validates the discovery under random scaling group transformation in the latent dimension.

5 Discussion

In this paper, we design and implement a Bayesian SINDy autoencoder for automated coordinate and governing equation discovery from high-dimensional data. Using the Bayesian learning framework and sparsity-promoting priors, the proposed model identifies sparse dynamics in the latent dimension. Through experimental studies on synthetic Lorenz, a reaction-diffusion system, a pendulum, and real video of moving pendulum. We successfully perform model discovery under a learned coordinate system. In the small-data and high-noise regime for real video data, we identify the correct physical law with close estimation of the standard gravity constant g . Besides the model and physics discovery, in terms of video prediction, the Bayesian SINDy autoencoder provides uncertainty-aware future predictions with an exact understanding of its underlying physics, which enables a trustworthy alternative for deep learning-based video prediction.

From our perspective, the Bayesian SINDy autoencoder has great potential to enable the concept of ‘GoPro physics’ [97], given the initial accomplishment in real video data discovery. However,

the current framework does have several limitations, which hinder a more general application in various problems. The largest limitation comes from the current framework requiring a high manual workload on hyperparameter tuning. The Bayesian SINDy autoencoder has over ten hyperparameters, including the encoder-decoder size, loss function weights, and the sparsity-promoting prior. The current inferential process relies heavily on a reasonable hyper-parameter setting. There are several potential strategies to improve the hyper-parameter tuning part to improve the robustness of the training outcome.

1. The training process will be unstable and converge to a suboptimal solution if $\lambda_1, \lambda_2, \lambda_3$ is misspecified. To solve this problem, one can directly solve the multi-objective optimization for encoder-decoder architecture using methods like Multiple Gradient Descent Algorithm [80]. Such adaptations can not only work in this paper on Bayesian SINDy autoencoder, but will also be able to generalize into a wider context for physics-informed machine learning.
2. The hyperparameter setting for the sparsity level δ , spike-and-slab distribution variances ν_0, ν_1 , and environmental noise σ can be tackled using a full Bayesian inference setup as described in [26, 71]. However, this will exceedingly enlarge the computational requirement and the difficulty in implementation.
3. Under minor misspecification of $\lambda_1, \lambda_2, \lambda_3$, one may consider to utilize an ensemble of trained Bayesian SINDy autoencoder initialized from various starting points [96, 99]. This process could ultimately construct a posterior distribution of the SINDy coefficient Ξ , and we can perform subset selection via stability selection or inclusion probability thresholding [29, 62]. Another consideration would be applying variational autoencoders [50] to avoid the hard-thresholding procedure. Applying variational autoencoders could also help to perform better subset selection with Bayesian uncertainty estimation to the SINDy coefficient [50].
4. The global minima exists, representing the real governing equation with the best coordinate system, but the nonconvex nature of neural network training makes it difficult to be found. The nonconvexity in neural network loss space is an open problem with very limited theoretical guarantees. To illustrate the existence of the global minima, from the case study for the moving rod (pendulum) video, the discovery with the $\sin(z)$ term leads to the minimal aggregated loss (1.098) throughout all experimental trails. All other suboptimal discoveries have a larger aggregated loss (for example via the Laplace prior, the discovery of z has aggregated loss 1.559). It is also observable from the data reconstruction (like in Fig. 8) where the result from the Spike-and-slab prior is visually the best. However, there is no guarantee that the neural network training could always converge to the global minima. Again, applying the Bayesian Model Averaging initializing from various random initial conditions (as mentioned in the previous point) could be helpful, and this will be an interesting future direction.
5. The training of encoder-decoder structure requires fixing an approximately correct neural network depth and width. The current study on autoencoders still lacks a complete understanding of the structure. In general, an oversized encoder-decoder will result in an underdetermined system which is likely to result in an overfitted encoder-decoder. This will hinder the autoencoder from extracting the correct latent dimension as expected. Convolutional filters could be helpful in reducing the possibility of learning a wrong coordinate system by focusing more on the 2D features.

In conclusion, the Bayesian SINDy autoencoder is an improved option compared to the SINDy autoencoder especially in the small-data and high-noise limit. The Bayesian learning framework not

only accelerates the training but also has dominating performance in various cases. Additionally, the Bayesian SINDy autoencoder enables uncertainty estimation, which is essential for data-driven decision-making purposes. Experimentally, we demonstrate the effectiveness of the Bayesian SINDy autoencoder using both synthetic and real video data. With only 14 seconds of real moving pendulum video, the Bayesian SINDy autoencoder successfully discovers the coordinate system and governing equation and estimates the fundamental constant g . In future works, we hope to extend the current methodology into more complex scenarios with real video to explore further the power of Bayesian SINDy autoencoder in broad applications.

Acknowledgments

We acknowledge support from the National Science Foundation AI Institute in Dynamic Systems (grant number 2112085). JNK further acknowledges support from the Air Force Office of Scientific Research (FA9550-19-1-0011). We would also like to thank Dr. Wei Deng, Dr. Bethany Lusch, Prof. Simon Du, Prof. Lexing Ying, and Dr. Joseph Bakarji for their insightful suggestions and comments.

Appendix: Real video discovery via the Laplace prior

We present the result from the Laplace prior which is suboptimal in this case. In this case with very low and noisy data, the Laplace prior identifies an incorrect term z over than the true target on $\sin(z)$. However, the discovery on z is an elegant failure due to the z term also captures part of the dynamics via the identification on z .

For the setting of Laplace prior, we set $\lambda_3 = 1 \times 10^{-2}$. The success discovery rate is 70% from 10 training instances while the other 30% instances will lead to zero discovery, indicating no active indices. Bayesian SINDy Autoencoder requires 1,500 epochs for training. The best test error is 0.779 for decoder loss, 0.364 for the reconstruction of \ddot{x} , and 0.416 for the reconstruction of \ddot{z} . The fraction of unexplained variance of decoder reconstruction is 0.012. For SINDy predictions, the fraction of unexplained variance of \ddot{x} and \ddot{z} reconstruction are 0.210 and 0.941.

The posterior samples from Bayesian SINDy are shown in Fig. 11 (a) and the prediction of trajectory with uncertainty quantification of latent dynamics is shown in Fig. 11 (c). The discovered governing equation from the Laplace prior is

$$\ddot{z} = -7.10z. \tag{26}$$

The discovered governing equation with z slightly deviates from the true dynamics similarly as reported in [15], but still creates reasonable uncertainty estimate under model misspecification. The latent dimension in Fig. 11 (b) follows a similar process in Fig. 10. The latent dimension z still matches the human labels, but not as close as the learned result by the SSGL prior.

References

- [1] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021.

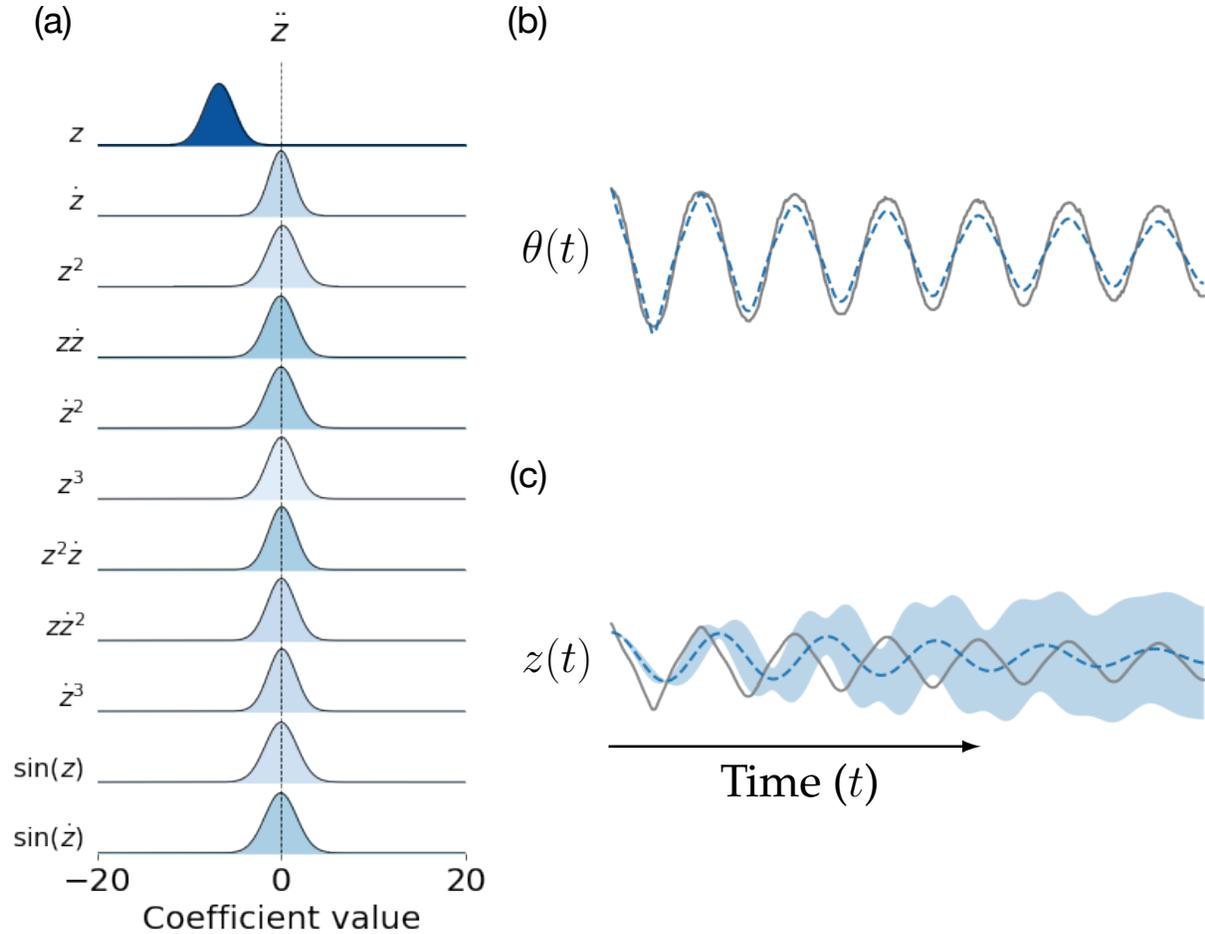


Figure 11: (a) Bayesian estimation and uncertainty quantification visualization of SINDy coefficient for real moving rod data under Laplace prior. (b) Latent dimension visualization after rescaling (blue curve) versus manually labelled moving rod angle by human (grey curve). (c) Visualization of predicted dynamics in latent space with uncertainty quantification.

- [2] E Paulo Alves and Frederico Fiuza. Data-driven discovery of reduced plasma physics models from fully kinetic simulations. *Physical Review Research*, 4(3):033192, 2022.
- [3] Arash Amini, Ulugbek S Kamilov, and Michael Unser. The analog formulation of sparsity implies infinite divisibility and rules out bernoulli-gaussian priors. In *2012 IEEE Information Theory Workshop*, pages 682–686. Ieee, 2012.
- [4] Jincheng Bai, Qifan Song, and Guang Cheng. Efficient variational inference for sparse deep learning with theoretical guarantee. *Advances in Neural Information Processing Systems*, 33:466–476, 2020.
- [5] Xuchan Bao, James Lucas, Sushant Sachdeva, and Roger B Grosse. Regularized linear autoencoders recover the principal components, eventually. *Advances in Neural Information Processing Systems*, 33:6971–6981, 2020.
- [6] Sarah Beetham and Jesse Capecelatro. Formulating turbulence closures using sparse regression with embedded form invariance. *Physical Review Fluids*, 5(8):084611, 2020.

- [7] Sarah Beetham, Rodney O Fox, and Jesse Capecelatro. Sparse identification of multiphase turbulence closures for coupled fluid–particle flows. *Journal of Fluid Mechanics*, 914, 2021.
- [8] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR, 2015.
- [9] Lorenzo Boninsegna, Feliks Nüske, and Cecilia Clementi. Sparse learning of stochastic dynamical equations. *The Journal of chemical physics*, 148(24):241723, 2018.
- [10] Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the national academy of sciences*, 113(15):3932–3937, 2016.
- [11] Jared L Callaham, J-C Loiseau, Georgios Rigas, and Steven L Brunton. Nonlinear stochastic modelling with langevin regression. *Proceedings of the Royal Society A*, 477(2250):20210092, 2021.
- [12] Jared L Callaham, Georgios Rigas, Jean-Christophe Loiseau, and Steven L Brunton. An empirical mean-field model of symmetry-breaking in a turbulent wake. *Science Advances*, 8(19):eabm4786, 2022.
- [13] Carlos M Carvalho, Nicholas G Polson, and James G Scott. Handling sparsity via the horseshoe. In *Artificial Intelligence and Statistics*, pages 73–80. PMLR, 2009.
- [14] Carlos M Carvalho, Nicholas G Polson, and James G Scott. The horseshoe estimator for sparse signals. *Biometrika*, 97(2):465–480, 2010.
- [15] Kathleen Champion, Bethany Lusch, J Nathan Kutz, and Steven L Brunton. Data-driven discovery of coordinates and governing equations. *Proceedings of the National Academy of Sciences*, 116(45):22445–22451, 2019.
- [16] Kathleen P Champion, Steven L Brunton, and J Nathan Kutz. Discovery of nonlinear multiscale systems: Sampling strategies and embeddings. *SIAM Journal on Applied Dynamical Systems*, 18(1):312–333, 2019.
- [17] Rohitash Chandra, Mahir Jain, Manavendra Maharana, and Pavel N Krivitsky. Revisiting bayesian autoencoders with mcmc. *IEEE Access*, 10:40482–40495, 2022.
- [18] Aoxue Chen and Guang Lin. Robust data-driven discovery of partial differential equations with time-dependent coefficients. *arXiv preprint arXiv:2102.01432*, 2021.
- [19] Boyuan Chen, Kuang Huang, Sunand Raghupathi, Ishaan Chandratreya, Qiang Du, and Hod Lipson. Discovering state variables hidden in experimental data. *arXiv preprint arXiv:2112.10755*, 2021.
- [20] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In *International conference on machine learning*, pages 1683–1691. PMLR, 2014.
- [21] Magnus Dam, Morten Brøns, Jens Juul Rasmussen, Volker Naulin, and Jan S Hesthaven. Sparse identification of a predator-prey system from simulation data of a convection model. *Physics of Plasmas*, 24(2):022310, 2017.

- [22] Nan Deng, Bernd R Noack, Marek Morzyński, and Luc R Pastur. Galerkin force model for transient and post-transient dynamics of the fluidic pinball. *Journal of Fluid Mechanics*, 918, 2021.
- [23] Wei Deng, Qi Feng, Liyao Gao, Faming Liang, and Guang Lin. Non-convex learning via replica exchange stochastic gradient mcmc. In *International Conference on Machine Learning*, pages 2474–2483. PMLR, 2020.
- [24] Wei Deng, Guang Lin, and Faming Liang. A contour stochastic gradient langevin dynamics algorithm for simulations of multi-modal distributions. *Advances in neural information processing systems*, 33:15725–15736, 2020.
- [25] Wei Deng, Guang Lin, and Faming Liang. An adaptively weighted stochastic gradient mcmc algorithm for monte carlo simulation and global optimization. *Statistics and Computing*, 32(4):1–24, 2022.
- [26] Wei Deng, Xiao Zhang, Faming Liang, and Guang Lin. An adaptive empirical bayesian method for sparse deep learning. *Advances in neural information processing systems*, 32, 2019.
- [27] Nan Ding, Youhan Fang, Ryan Babbush, Changyou Chen, Robert D Skeel, and Hartmut Neven. Bayesian sampling using stochastic gradient thermostats. *Advances in neural information processing systems*, 27, 2014.
- [28] D Elizondao, Emile Fiesler, and Jerzy Korczak. Non-ontogenic sparse neural networks. In *Proceedings of ICNN’95-International Conference on Neural Networks*, volume 1, pages 290–295. IEEE, 1995.
- [29] Urban Fasel, J Nathan Kutz, Bingni W Brunton, and Steven L Brunton. Ensemble-sindy: Robust sparse model discovery in the low-data, high-noise limit, with active learning and control. *Proceedings of the Royal Society A*, 478(2260):20210904, 2022.
- [30] Patrick Gelß, Stefan Klus, Jens Eisert, and Christof Schütte. Multidimensional approximation of nonlinear dynamical systems. *Journal of Computational and Nonlinear Dynamics*, 14(6), 2019.
- [31] Edward I George and Robert E McCulloch. Approaches for bayesian variable selection. *Statistica sinica*, pages 339–373, 1997.
- [32] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323. JMLR Workshop and Conference Proceedings, 2011.
- [33] Yifei Guan, Steven L Brunton, and Igor Novosselov. Sparse nonlinear models of chaotic electroconvection. *Royal Society Open Science*, 8(8):202367, 2021.
- [34] John Guckenheimer and Philip Holmes. *Nonlinear oscillations, dynamical systems, and bifurcations of vector fields*, volume 42. Springer Science & Business Media, 2013.
- [35] Fredrik K Gustafsson, Martin Danelljan, and Thomas B Schon. Evaluating scalable bayesian deep learning methods for robust computer vision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 318–319, 2020.
- [36] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.

- [37] José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *International conference on machine learning*, pages 1861–1869. PMLR, 2015.
- [38] Charles Gordon Hewitt. *The conservation of the wild life of Canada*. New York: C. Scribner, 1921.
- [39] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.
- [40] Seth M Hirsh, David A Barajas-Solano, and J Nathan Kutz. Sparsifying priors for bayesian uncertainty quantification in model discovery. *Royal Society Open Science*, 9(2):211823, 2022.
- [41] Matthew Hoffman and Yi-An Ma. Black-box variational inference as distilled langevin dynamics. In *Proceedings of the 37th International Conference on Machine Learning*, pages 4324–4341, 2020.
- [42] Matthew D Hoffman, Andrew Gelman, et al. The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *J. Mach. Learn. Res.*, 15(1):1593–1623, 2014.
- [43] Moritz Hoffmann, Christoph Fröhner, and Frank Noé. Reactive sindy: Discovering governing reactions from concentration data. *The Journal of chemical physics*, 150(2):025101, 2019.
- [44] Hemant Ishwaran and J Sunil Rao. Spike and slab variable selection: frequentist and bayesian strategies. *The Annals of Statistics*, 33(2):730–773, 2005.
- [45] Kadierdan Kaheman, J Nathan Kutz, and Steven L Brunton. Sindy-pi: a robust algorithm for parallel implicit sparse identification of nonlinear dynamics. *Proceedings of the Royal Society A*, 476(2242):20200279, 2020.
- [46] Eurika Kaiser, J Nathan Kutz, and Steven L Brunton. Sparse identification of nonlinear dynamics for model predictive control in the low-data limit. *Proceedings of the Royal Society A*, 474(2219):20180335, 2018.
- [47] Alan A Kaptanoglu, Jared L Callahan, Aleksandr Aravkin, Christopher J Hansen, and Steven L Brunton. Promoting global stability in data-driven models of quadratic nonlinear dynamics. *Physical Review Fluids*, 6(9):094401, 2021.
- [48] Alan A Kaptanoglu, Kyle D Morgan, Chris J Hansen, and Steven L Brunton. Physics-constrained, low-dimensional models for magnetohydrodynamics: First-principles and data-driven approaches. *Physical Review E*, 104(1):015206, 2021.
- [49] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.
- [50] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [51] J Nathan Kutz, Steven L Brunton, Bingni W Brunton, and Joshua L Proctor. *Dynamic mode decomposition: data-driven modeling of complex systems*. SIAM, 2016.
- [52] Zhilu Lai and Satish Nagarajaiah. Sparse structural system identification method for nonlinear dynamic systems with hysteresis/inelastic behavior. *Mechanical Systems and Signal Processing*, 117:813–842, 2019.

- [53] Liyuan Li, Weimin Huang, Irene Yu-Hua Gu, and Qi Tian. Statistical modeling of complex backgrounds for foreground object detection. *IEEE Transactions on Image Processing*, 13(11):1459–1472, 2004.
- [54] Baoyuan Liu, Min Wang, Hassan Foroosh, Marshall Tappen, and Marianna Pensky. Sparse convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 806–814, 2015.
- [55] Jean-Christophe Loiseau. Data-driven modeling of the chaotic thermal convection in an annular thermosyphon. *Theoretical and Computational Fluid Dynamics*, 34(4):339–365, 2020.
- [56] Jean-Christophe Loiseau and Steven L Brunton. Constrained sparse galerkin regression. *Journal of Fluid Mechanics*, 838:42–67, 2018.
- [57] Jean-Christophe Loiseau, Bernd R Noack, and Steven L Brunton. Sparse reduced-order modelling: sensor-based dynamics to full-state estimation. *Journal of Fluid Mechanics*, 844:459–490, 2018.
- [58] Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks through l_0 regularization. *arXiv preprint arXiv:1712.01312*, 2017.
- [59] Yi-An Ma, Tianqi Chen, and Emily Fox. A complete recipe for stochastic gradient mcmc. *Advances in neural information processing systems*, 28, 2015.
- [60] Niall M Mangan, Travis Askham, Steven L Brunton, J Nathan Kutz, and Joshua L Proctor. Model selection for hybrid dynamical systems via sparse regression. *Proceedings of the Royal Society A*, 475(2223):20180534, 2019.
- [61] Niall M Mangan, Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Inferring biological networks by sparse identification of nonlinear dynamics. *IEEE Transactions on Molecular, Biological and Multi-Scale Communications*, 2(1):52–63, 2016.
- [62] Nicolai Meinshausen and Peter Bühlmann. Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4):417–473, 2010.
- [63] Daniel A Messenger and David M Bortz. Weak sindy for partial differential equations. *Journal of Computational Physics*, 443:110525, 2021.
- [64] Toby J Mitchell and John J Beauchamp. Bayesian variable selection in linear regression. *Journal of the american statistical association*, 83(404):1023–1032, 1988.
- [65] Peter H Morgan. Differential evolution and sparse neural networks. *Expert Systems*, 25(4):394–413, 2008.
- [66] Radford M Neal et al. Mcmc using hamiltonian dynamics. *Handbook of markov chain monte carlo*, 2(11):2, 2011.
- [67] Trevor Park and George Casella. The bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.
- [68] Nicholas G Polson and Veronika Ročková. Posterior concentration for sparse deep learning. *Advances in Neural Information Processing Systems*, 31, 2018.

- [69] Patrick AK Reinbold, Daniel R Gurevich, and Roman O Grigoriev. Using noisy or incomplete data to discover models of spatiotemporal dynamics. *Physical Review E*, 101(1):010203, 2020.
- [70] Patrick AK Reinbold, Logan M Kageorge, Michael F Schatz, and Roman O Grigoriev. Robust learning from noisy, incomplete, high-dimensional experimental data via physically constrained symbolic regression. *Nature communications*, 12(1):1–8, 2021.
- [71] Veronika Ročková and Edward I George. Emvs: The em approach to bayesian variable selection. *Journal of the American Statistical Association*, 109(506):828–846, 2014.
- [72] Veronika Ročková and Edward I George. The spike-and-slab lasso. *Journal of the American Statistical Association*, 113(521):431–444, 2018.
- [73] Samuel Rudy, Alessandro Alla, Steven L Brunton, and J Nathan Kutz. Data-driven identification of parametric partial differential equations. *SIAM Journal on Applied Dynamical Systems*, 18(2):643–660, 2019.
- [74] Samuel H Rudy, Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Data-driven discovery of partial differential equations. *Science advances*, 3(4):e1602614, 2017.
- [75] Hayden Schaeffer. Learning partial differential equations via data discovery and sparse optimization. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 473(2197):20160446, 2017.
- [76] Hayden Schaeffer and Scott G McCalla. Sparse model selection via integral terms. *Physical Review E*, 96(2):023302, 2017.
- [77] Martin Schmelzer, Richard P Dwight, and Paola Cinnella. Discovery of algebraic reynolds-stress models using sparse symbolic regression. *Flow, Turbulence and Combustion*, 104(2):579–603, 2020.
- [78] Michael Schmidt and Hod Lipson. Distilling free-form natural laws from experimental data. *science*, 324(5923):81–85, 2009.
- [79] Steven L Scott and Hal R Varian. Predicting the present with bayesian structural time series. *International Journal of Mathematical Modelling and Numerical Optimisation*, 5(1-2):4–23, 2014.
- [80] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018.
- [81] Michael Small, Kevin Judd, and Alistair Mees. Modeling continuous processes from data. *Physical Review E*, 65(4):046704, 2002.
- [82] Mariia Sorokina, Stylianos Sygletos, and Sergei Turitsyn. Sparse identification for nonlinear optical communication systems: Sino method. *Optics express*, 24(26):30433–30443, 2016.
- [83] Suraj Srinivas, Akshayvarun Subramanya, and R Venkatesh Babu. Training sparse neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 138–145, 2017.
- [84] Yan Sun, Qifan Song, and Faming Liang. Consistent sparse deep learning: Theory and computation. *Journal of the American Statistical Association*, pages 1–15, 2021.

- [85] Yan Sun, Qifan Song, and Faming Liang. Learning sparse deep neural networks with a spike-and-slab prior. *Statistics & Probability Letters*, 180:109246, 2022.
- [86] Yan Sun, Wenjun Xiong, and Faming Liang. Sparse deep learning: A new framework immune to local traps and miscalibration. *Advances in Neural Information Processing Systems*, 34:22301–22312, 2021.
- [87] Yee Whye Teh, Alexandre H Thiery, and Sebastian J Vollmer. Consistency and fluctuations for stochastic gradient langevin dynamics. *Journal of Machine Learning Research*, 17, 2016.
- [88] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [89] Bin Wang, Jie Lu, Zheng Yan, Huaishao Luo, Tianrui Li, Yu Zheng, and Guangquan Zhang. Deep uncertainty quantification: A machine learning approach for weather forecasting. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2087–2095, 2019.
- [90] Hao Wang and Dit-Yan Yeung. A survey on bayesian deep learning. *ACM Computing Surveys (CSUR)*, 53(5):1–37, 2020.
- [91] Wen-Xu Wang, Rui Yang, Ying-Cheng Lai, Vassilios Kovanis, and Celso Grebogi. Predicting catastrophes in nonlinear dynamical systems by compressive sensing. *Physical review letters*, 106(15):154101, 2011.
- [92] Yating Wang, Wei Deng, and Guang Lin. Bayesian sparse learning with preconditioned stochastic gradient mcmc and its applications. *Journal of Computational Physics*, 432:110134, 2021.
- [93] Yuexi Wang and Veronika Rocková. Uncertainty quantification for sparse deep learning. In *International Conference on Artificial Intelligence and Statistics*, pages 298–308. PMLR, 2020.
- [94] Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. *Advances in neural information processing systems*, 28, 2015.
- [95] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011.
- [96] Andrew G Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. *Advances in neural information processing systems*, 33:4697–4708, 2020.
- [97] Charlie Wood. Powerful ‘machine scientists’ distill the laws of physics from raw data. *Quanta Magazine*, May 2022.
- [98] John Wright and Yi Ma. *High-dimensional data analysis with low-dimensional models: Principles, computation, and applications*. Cambridge University Press, 2022.
- [99] Dongxia Wu, Liyao Gao, Xinyue Xiong, Matteo Chinazzi, Alessandro Vespignani, Yi-An Ma, and Rose Yu. Quantifying uncertainty in deep spatiotemporal forecasting. *arXiv preprint arXiv:2105.11982*, 2021.

- [100] Haohang Xu, Shuangrui Ding, Xiaopeng Zhang, Hongkai Xiong, and Qi Tian. Masked autoencoders are robust data augmentors. *arXiv preprint arXiv:2206.04846*, 2022.
- [101] Liu Yang, Xuhui Meng, and George Em Karniadakis. B-pinns: Bayesian physics-informed neural networks for forward and inverse pde problems with noisy data. *Journal of Computational Physics*, 425:109913, 2021.
- [102] Chen Yao and Erik M Bollt. Modeling and nonlinear parameter estimation with kronecker product representation for coupled oscillators and spatiotemporal systems. *Physica D: Nonlinear Phenomena*, 227(1):78–99, 2007.
- [103] Laure Zanna and Thomas Bolton. Data-driven equation discovery of ocean mesoscale closures. *Geophysical Research Letters*, 47(17):e2020GL088376, 2020.
- [104] Linan Zhang and Hayden Schaeffer. On the convergence of the sindy algorithm. *Multiscale Modeling & Simulation*, 17(3):948–972, 2019.
- [105] Ruqi Zhang, Chunyuan Li, Jianyi Zhang, Changyou Chen, and Andrew Gordon Wilson. Cyclical stochastic gradient mcmc for bayesian deep learning. *arXiv preprint arXiv:1902.03932*, 2019.
- [106] Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradient langevin dynamics. In *Conference on Learning Theory*, pages 1980–2022. PMLR, 2017.
- [107] Peng Zheng, Travis Askham, Steven L Brunton, J Nathan Kutz, and Aleksandr Y Aravkin. A unified framework for sparse relaxed regularized regression: Sr3. *IEEE Access*, 7:1404–1423, 2018.
- [108] Yin hao Zhu and Nicholas Zabaras. Bayesian deep convolutional encoder–decoder networks for surrogate modeling and uncertainty quantification. *Journal of Computational Physics*, 366:415–447, 2018.