

A study on quantifying effective training of DLDMD

A Thesis
Presented to the
Faculty of
San Diego State University

In Partial Fulfillment
of the Requirements for the Degree
Master of Science in Applied Mathematics
with a Concentration in
Dynamical Systems

by
Joseph Aaron Gene Diaz

July 18, 2022

SAN DIEGO STATE UNIVERSITY

The Undersigned Faculty Committee Approves the
Thesis of Joseph Aaron Gene Diaz:

A study on quantifying effective training of DLDMD

Christopher Curtis, Chair
Department of Mathematics and Statistics

Jérôme Gilles
Department of Mathematics and Statistics

Saeed Manshadi
Department of Electrical Engineering

Approval Date

Copyright © 2022
by
Joseph Aaron Gene Diaz

DEDICATION

Everyone who's helped me get this far: You know who you are.

ABSTRACT OF THE THESIS

A study on quantifying effective training of DLDMD

by

Joseph Aaron Gene Diaz

Master of Science in Applied Mathematics with a Concentration in Dynamical Systems

San Diego State University, 2022

Write this once everything else is written.

TABLE OF CONTENTS

	PAGE
ABSTRACT	v
LIST OF TABLES	vii
LIST OF FIGURES	viii
CHAPTER	
1 INTRODUCTION	1
1.1 What is E-DMD?	1
1.2 What are Neural Networks?	3
1.3 DLDMD and it's limitations	5
1.4 Statement of Dilemma	6
2 KULLBACK-LEIBLER DIVERGENCE	8
2.1 Basic Definitions	9
2.2 How to build distributions (Non-parametric statistics)	10
2.3 Measuring entropy flow	11
2.4 Implementation Details	12
3 RESULTS	14
3.1 The Duffing Oscillator	14
3.2 The Van der Pol Oscillator	17
4 DISCUSSION	26
4.1 Non-unitary densities	26
4.2 Statistical issues with edge layers	26
4.3 Future work	27
BIBLIOGRAPHY	28

LIST OF TABLES

	PAGE
3.1 Hyperparamters of the DLDMD algorithm held constant for each training run on Duffing.	16
3.2 ReLU activation function vs ELU activation function.	18

LIST OF FIGURES

	PAGE
2.1 Entropy and density function change with respect to σ^2 . The base of the logarithm shown here is e , but the base itself is irrelevant due to the shared properties of all logarithms: that they increase monotonically and unbounded. As the variance increases we see higher entropy in the left plot and a flattening of the probability density in the right plot.	9
2.2 KL Divergence example with Normal distributions. As expected, the divergence decreases to 0 as $\sigma^2 \rightarrow 1$ from either direction.	10
2.3 Example of divergences for layer Dec 2 of Van der Pol model with latent dimension 8 with linear fitting and confidence bands.	13
3.1 Duffing phase space.	15
3.2 Duffing average slopes with standard deviations.	17
3.3 Duffing loss plots, smoothed to accentuate trends.	18
3.4 Duffing Average bandwidth.	18
3.5 Duffing slope linear fits by layer.	19
3.6 Duffing slope linear fits by layer.	20
3.7 Van der Pol phase space.	21
3.8 Van der Pol average slopes with standard deviations.	22
3.9 Van der Pol loss plots, smoothed to accentuate trend.	23
3.10 Van der Pol Average bandwidth.	23
3.11 Van der Pol slope linear fits by layer.	24
3.12 Van der Pol slope linear fits by layer.	25
4.1 Histograms of computed areas from KDE.	27

CHAPTER 1

INTRODUCTION

In the study of dynamical systems a central problem is how to derive models from measured data to facilitate the prediction of future states. Many approaches and techniques exist in the literature, from deriving sets of governing equations via application of the simple physical principles to the statistically meaningful Principal Component Analysis and other modal decompositions. The main goal of many of these methods is to come up with a generalized framework so that the dynamics can be extracted from the data for the sake of control and prediction. While there are genuine advantages to the more concrete and scientifically-sound methods of deriving models from first principles, it is often very difficult if not outright untenable.

As such, data-driven methods have rapidly become very useful approaches for coming up with rough and ready models.

Say more once you've thought of it.

1.1 What is E-DMD?

As it says on the tin, we seek a predictive model for a time series $\{\mathbf{y}_j\}_{j=1}^{N^T+1}$, which are the measurements of a uniformly sampled unknown dynamical system of the form

$$\frac{d\mathbf{y}}{dt} = f(\mathbf{y}(t)), \quad \mathbf{y}(0) = \mathbf{x} \in \mathcal{M} \subseteq \mathbb{R}^{N_s}$$

As established in the introduction, there is an old and venerable literature dedicated to deriving system using classical methods. What this literature will admit is that such a process is not trivial and even difficult or impossible to do in practice; particularly for nonlinear phenomena worth investigating. We want something that is more easily generalized and algorithmic.

One method that can be leveraged for this is via the Koopman Operator \mathcal{K} . If we denote $\boldsymbol{\varphi}(t; \mathbf{x}) = \mathbf{y}(t)$ to be the flow map affiliated with the initial condition \mathbf{x} and denote $g : \mathcal{M} \rightarrow \mathbb{C}$ to be a square integrable observable of the system then, according to [7], there exists a linear representation of the flow map given by \mathcal{K} ;

$$\mathcal{K}^t g(\mathbf{x}) = g(\boldsymbol{\varphi}(t; \mathbf{x})),$$

which means that \mathcal{K} is a linear time-advancement operator of the dynamics. This might seem like it trivializes the problem, given that we now have a linear system, but it does

not. The Koopman operator is an *infinite* dimensional operator, so we've traded a potentially nonlinear problem for an infinite dimensional linear one.

In order to actually use this new representation, it suffices to find the eigenvalues $\{\lambda_\ell\}$, and their affiliated eigenfunctions $\{\phi_\ell\}$, of the Koopman Operator such that

$$\mathcal{K}^t \phi_\ell = \exp(t\lambda_\ell) \phi_\ell \implies g(\mathbf{x}) = \sum_{\ell \in \mathbb{N}} c_\ell \phi_\ell(\mathbf{x}),$$

where we can essentially construct a modal decomposition of g . From here advancing the dynamics to time t is equivalent to writing

$$\mathcal{K}^t g(\mathbf{x}) = \sum_{\ell \in \mathbb{N}} a_\ell \phi_\ell(\varphi(t, \mathbf{x}))$$

The most useful property of this framing is that we have a global linearization of the flow, with a major caveat: Generally, finding these eigenvalues and eigenfunctions is impossible. This led to the development of the already mentioned and much much lauded Dynamic Mode Decomposition (DMD) and its extensions, which seeks to numerically approximate a finite number of these modes. We'll focus on the particular extension that [1] implements: Extended DMD (E-DMD) [10]. With E-DMD, we suppose that

$$g(\mathbf{x}) = \sum_{\ell=1}^{N_O} a_\ell \psi_\ell(\mathbf{x})$$

which is to say that the observable g exists in a finite dimensional subspace of $L^2(\mathcal{M})$, applying the Koopman operator implies that

$$\begin{aligned} \mathcal{K}^{\delta t} g(\mathbf{x}) &= \sum_{\ell=1}^{N_O} a_\ell \phi_\ell(\varphi(\delta t, \mathbf{x})) \\ &= \sum_{\ell=1}^{N_O} \phi_\ell(\mathbf{x}) (\mathbf{K}_O^T \mathbf{a})_\ell + r(\mathbf{x}; \mathbf{K}_O) \end{aligned}$$

for discrete time step δt . \mathbf{K}_O is the $N_O \times N_O$ matrix that minimizes

$$\mathbf{K}_O = \underset{K}{\operatorname{argmin}} ||\Psi_+ - K\Psi_-||_F^2$$

and $r(\mathbf{x}, \mathbf{K}_O)$ is a residual that represents the total error due to DMD. If the ansatz that g lives in a finite space holds, then r is identically 0. We define Ψ_\pm to be

$$\Psi_- = (\Psi_1 \ \Psi_2 \ \cdots \ \Psi_{N_T}), \quad \Psi_+ = (\Psi_2 \ \Psi_3 \ \cdots \ \Psi_{N_T+1})$$

where $\{\Psi_j\}$ is an observable of the time series of interest $\{\mathbf{y}_j\}$. What the expression for \mathbf{K}_O tells us is that, we are trying to find a one-step mapping from each data point to the next. Practically speaking, \mathbf{K}_O is found using an SVD, with which we can write

$$\Psi_- = U\Sigma W^\dagger \implies \mathbf{K}_O = \Psi_+ W \Sigma^{-P} U^\dagger$$

where $-P$ denotes the Moore-Penrose pseudo-inverse and \dagger denotes the conjugate transpose. This gives us an expression for r in terms of the observables Ψ :

$$E_r(\mathbf{K}_O) = \|\Psi_+(I - WW^\dagger)\|_F$$

Finding the eigenvalues, eigenfunctions, and Koopman modes comes down to an eigen-decomposition, from which the dynamics can be approximated as

$$y(t; \mathbf{x}) \approx V \exp(t\Lambda) V^{-1} \Psi(\mathbf{x})$$

where $\mathbf{K}_O = VTV^{-1}$, $\Lambda_{\ell\ell} = \ln(T_{\ell\ell})/\delta t$ is a diagonal matrix and Ψ is the representation of the initial condition in terms of the observables.

Expand on everything else above for a more complete picture of E-DMD.

1.2 What are Neural Networks?

Before we move on, a brief digression on Neural Networks (NNs) is appropriate. Ever since scientists discovered the relatively simple interaction between neurons and axons in the human brain, they have been enamored with the ability to create a learning computer with the similar ability to become more adept at a task with training and practice. In the same way that art imitates life, the most straight-forward attempts have been to construct artificial neural networks that pantomime the biological ones that we carry around in our heads; and while these pale imitations have not been developed to rival the human brain, they have led to some major achievements in automation. In much the same way that a person is “trained” to perform a task by repeating it with feedback and adjusting their performance as they go, an artificial NN uses data, a loss function, and an optimization algorithm to change the state of the NN to one which can better accomplish the task. The loss function and optimizer are to the feedback and behavioral adjustment what the data is to the experience.

The oldest known example of NNs were the perceptrons of the 1950s, which could perform binary classifications of linearly separable data. Multilayer versions of this simple architecture allowed for classification of larger set of classes, but the researcher of the time arrived at the limits of their computing environments rather

quickly and a brief dark age in the study and application of neural networks ensued. Interest was renewed in the 1980s when many of the hardware limitations of previous decades were ameliorated and new algorithms for optimizing had been pioneered.

Expand once you've got more to say.

In modern mathematical study and application, the question of function representation is an ever-present one; specifically with respect to approximating quantities that can be given via series. While truncations and use of the likes of fourier series are useful and effective in some narrow applications, one means of representation that has shown much promise in recent years due to the greater availability of computing power is through the use of Neural Networks (NNs). In 1989, [6] showed that, with some appropriate parameters, a NN can be used to approximate any continuous function from one vector space to another to any arbitrary precision required.

From this we can consider ML as a tool of regression. Say we have a set of training data $d_{tr} = \{(x_j, y_j)\}_{j=1}^{N_{tr}}$ and we believe that there is some function $f(x)$ such that

$$y_j = f(x_j) + \varepsilon_j,$$

where ε_j is assumed to be a sample drawn from $\mathcal{N}(0, \bar{\sigma}^2)$. Generally f is not known; so we must, using the training data, build an estimator $f_e(x; d_{tr})$ and measure the performance of the estimator by examining the quantity

$$\mathbb{E}_{d_{tr}} [(f(x) + \varepsilon - f_e(x; d_{tr}))^2] = \text{bias}_{d_{tr}}^2 + \mathbb{V}(f_e) + \bar{\sigma}^2$$

where

$$\text{bias}_{d_{tr}} = \mathbb{E}_{d_{tr}} [f_e] - f(x)$$

and

$$\mathbb{V}(f_e) = \mathbb{E}_{d_{tr}} [(\mathbb{E}_{d_{tr}} [f_e] - f_e(x; d_{tr}))^2]$$

With this framing, ML methods can be cast as attempts to find an estimator relative to a given training that balances the competing dilemmas of being too biased or allowing too much variance. From the point of view of non-parametric statistics, we would need to utilize other techniques[5] to determine the bias and variance of any given estimator. For our purposes, we will instead consider how metrics from information theory can be used to assess the change in information as data is fed through a NN and a model is trained.

1.3 DLDM and it's limitations

Incomplete picture, add more detail later. Consider re-skimming Jay's paper.

The key innovation of [1] is to use a neural network to come up with the collection of observables on $\{\mathbf{y}_j\}$ that allow for the best prediction of future system states. This is implemented by defining an encoder $\mathcal{E} : N_S \rightarrow N_O$ and decoder $\mathcal{D} : N_O \rightarrow N_S$ composed of Dense layers such that

$$(\mathcal{D} \circ \mathcal{E})(\mathbf{x}) = \mathbf{x}$$

We choose $N_O \geq N_S$ and an appropriate loss function so that \mathcal{E} and \mathcal{D} give a richer space of observables, called the latent space, for EDMD to use when advancing the dynamics. The implementation of NNs for this purpose requires a method of tuning to allow \mathcal{E} and \mathcal{D} to learn the best representations possible. As such, a loss function that correctly identifies and prioritizes the desired properties is a necessary condition for the DLDM to function as needed. A natural choice considering these constraints is given by

$$\mathcal{L} = \alpha_1 \mathcal{L}_{\text{recon}} + \alpha_2 \mathcal{L}_{\text{dmd}} + \alpha_3 \mathcal{L}_{\text{pred}} + \alpha_4 \|\mathbf{W}_g\|_2^2 \quad (1.1)$$

where

$$\mathcal{L}_{\text{recon}} = \frac{1}{N_T + 1} \sum_{j=1}^{N_T+1} \|\mathbf{y}_j - (\mathcal{D} \circ \mathcal{E})(\mathbf{y}_j)\|_2, \quad (1.2)$$

$$\mathcal{L}_{\text{dmd}} = E_r(\mathbf{K}_O), \quad (1.3)$$

$$\mathcal{L}_{\text{pred}} = \frac{1}{N_T} \sum_{j=1}^{N_T} \|\mathbf{y}_{j+1} - \mathcal{D}(VT^j V^{-1} \mathcal{E}(x))\|_2, \quad (1.4)$$

Each component guides the machine to a particular outcome:

1. $\mathcal{L}_{\text{recon}}$ is the Mean Squared Error (MSE) of each time step with respect to the reconstruction from the composition of \mathcal{E} and \mathcal{D} . This component ensures that, under training, the network effectively acts as a near identity transformation for data that is fed into it. This quality allows the DMD advanced trajectories to be recovered from the higher dimensional latent space back to the original dimension of the data.
2. \mathcal{L}_{DMD} is the error associated with DMD. Consequently, this component is the one that is most responsible for finding the optimal set of observables for DMD. \mathcal{E} immerses the data into a higher dimensional latent space, in effect acting as our set of observable; minimizing this gives us greater flexibility in the latent space.
3. $\mathcal{L}_{\text{pred}}$ is the (MSE) for each forward time prediction due to DMD and immersion/submersion due to \mathcal{E}/\mathcal{D} . In addition to balancing the last conditions,

this condition ensures that the DMD step in the latent dimension is consistent with the next time-step from the time series after encoding, advancing, and decoding.

4. \mathbf{W}_g is the vectorized quantity that represents all weights in both \mathcal{E} and \mathcal{D} . This is really only a regularization condition to keep the coefficients of the weight matrices from blowing up in value as the model trains, which can be a concern for ML models.
5. α_1 , α_2 , α_3 , and α_4 are 4 positive constants that allow us to assign a weighting to each component of the loss. This allows the loss function to be dynamically weighted to prioritize some conditions over others. In [1], $\alpha_1 = \alpha_2 = \alpha_3 = 1$ and $\alpha_4 < 10^{-10}$.

The marriage of E-DMD with Dense NNs and some meaningful choices of loss function as an innovation was made by [1, 3, 8], and has proved to be quite a robust method for learning dynamics from data. In researching this thesis, many attempts were made to test the limits of DLDMD. [1] was primarily interested in recreating the phase space behavior of a handful of well known nonlinear oscillators such as the nonlinear harmonic oscillator, the Duffing oscillator, and the Van der Pol oscillator. Additionally, chaotic systems like Lorenz 63 and the Rössler system were examined but chaos proved to be quite a challenge for the algorithm to overcome. A study into the addition of gaussian noise led to the discovery of the machine being quite robust to data due to the characteristics of it's loss function and additional noise mitigation was accomplished by implementing convolutional layers instead of traditional dense layers.

1.4 Statement of Dilemma

Having established DLDMD and having discussed testing it's limitations via training new models and *visually* qualifying each run as a success or failure, a question is begged: Can we, in an empirical and quantitative fashion, determine whether good training is taking place or otherwise deduce whether our model is well-posed given the hyperparameters? A machine learning model that must learn how to accomplish even a simple task can take quite a long time and plenty of computing resources to complete it's training and only then can the results be checked to determine whether productive training was taking place as opposed to the optimizer becoming stuck in a local minimum or there being no global minimum at all.

How can this be quantified? A useful starting point would be to examine the weights of the matrices that make up the layers of a model. For Dense layers, passing a vector of data $\mathbf{x} \in \mathbb{R}^d$ through a layer L can be written as

$$L(\mathbf{x}; A, \sigma, b) = \sigma(A\mathbf{x} + b)$$

for some matrix $A \in \mathbb{R}^{D \times d}$, vector $b \in \mathbb{R}^D$ and (typically nonlinear) activation function $\sigma : \mathbb{R}^D \rightarrow \mathbb{R}^D$. For a given layer, σ is a fixed hyperparameter of the machine, but both A and b have entries that are tuned by the optimizer as training takes place. Training takes place on a per epoch basis for most models, so considering the configuration of these matrices epoch to epoch can help us determine whether they are converging to a set of elements that accomplish the desired task or not. The next state for each layer is dependent on the previous, so a reasonable framing would be as one of a discrete dynamical system of the form

$$Q_{n+1} = P(Q_n) \quad (\text{Come up with better symbols later})$$

where P is the training procedure and Q_n represents the configuration of the network's weights at each epoch n . Typically, the initial state Q_0 is “random” in the sense that the weights are selected from a probability distribution. In the interest of examining convergences, one might consider the following limit

$$\lim ||Q_{n+1} - Q_n||_2$$

for some L^2 norm on the space that Q_n inhabits. While this “one-step” Cauchy convergence will tell us whether the machine is approaching a particular configuration point-wise, this might not actually tell us much of anything else. The changes between epochs is, in some sense, stochastic due to not knowing how the optimizer will update Q between epochs. As such, we propose a more statistical approach wherein we examine how the information content is changing from epoch to epoch. By the noted stochastic nature of the evolution, we consider the entries of A and b for each layer to be drawn from a continuous probability distribution X and examine how it changes epoch to epoch.

CHAPTER 2

KULLBACK-LEIBLER DIVERGENCE

In the field of information theory, it is often useful to frame information as being represented by a random variable and the values it can take on. For example, consider a normally distributed random variable $X \sim \mathcal{N}(\mu, \sigma^2)$ where μ is the mean or expected value $\mathbb{E}[X]$ of X and σ^2 is the variance $\mathbb{V}[X]$ of X . This is one of the most well understood probability distributions due to countless years of study into its properties, as such we know that if we were to draw a sample for X that we're much more likely to get something close to μ than we are to get something far away from μ . For example; if $\mu = 0$ and $\sigma^2 = 1$, then upon drawing a sample the value 1 is more likely to appear than -37. The level of certainty of the outcome of the draw begs the question of how that certainty can be quantified and conversely, we could instead consider the uncertainty of any given probability distribution. A measure of the uncertainty was proposed in 1948 and this measure is known as *informational entropy* [9] or, simply, entropy. For a continuous distribution with density function $f(x)$, the entropy is given by

$$h[f] = \mathbb{E}[-\log(f(x))] = - \int_X f(x) \log(f(x)) \, dx \quad (2.1)$$

For the normal distributed X defined above, we have that probability density function

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) \quad (2.2)$$

and the entropy integral for this evaluates to

$$h[f] = \frac{1}{2} (\log(2\pi\sigma^2) + 1) \quad (2.3)$$

which only depends on the variance. This expression tells us that the entropy of a normally distributed random variable increases as the variance does, graphically this is associated with the flattening of the bell curve which can be interpreted as the values around the mean being less likely to appear relative to values far from the mean. As such, the result of any given sample is less certain; Figure 2.1 demonstrates this.

When training an ML model, the expectation is that the machine will converge to some configuration that allows it to accomplish the task it was designed for. For the sake of

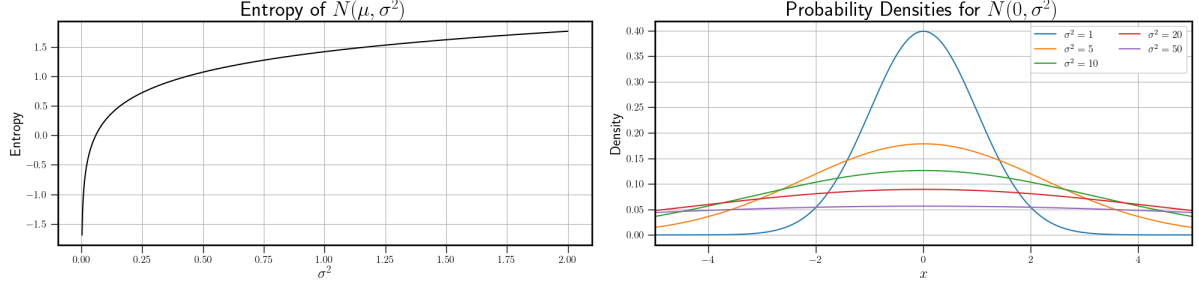


Figure 2.1. Entropy and density function change with respect to σ^2 . The base of the logarithm shown here is e , but the base itself is irrelevant due to the shared properties of all logarithms: that they increase monotonically and unbounded. As the variance increases we see higher entropy in the left plot and a flattening of the probability density in the right plot.

2.1 Basic Definitions

The KL Divergence (KLD) is a statistical distance between a pair of probability distributions which measures how different a distribution P is from a reference distribution Q . “A simple interpretation of the divergence of P from Q is the expected excess surprise from using Q as a model when the actual distribution is P .” If P and Q are discrete distributions defined on the probability space \mathcal{X} , then the KLD of P with respect to Q is given by

$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log (P(x)/Q(x)) \quad (2.4)$$

In other words, it is the expectation of the logarithmic difference between the probabilities P and Q , where the expectation is taken using the probabilities P . As a “distance”, the KLD is non-negative and 0 when $P = Q$ almost everywhere. Unlike more standard metrics, it is not symmetric and does not satisfy the triangle equality. In order to deal with something that is, at least, symmetric we’ll consider using the symmetric KL Divergence (SKLD) which is defined as

$$D_{SKL}(P, Q) = \frac{1}{2} (D_{KL}(P \parallel Q) + D_{KL}(Q \parallel P)) \quad (2.5)$$

Which is the average of the KLD for P with respect to Q and for Q with respect to P . In our cases, we’ll be dealing with a continuous random variable with probability density functions p and q , for which the corresponding KLD formula is

$$D_{KL}(P \parallel Q) = \int_{\mathcal{X}} p(x) \log (p(x)/q(x)) \, dx \quad (2.6)$$

As an analogue of the discrete case, the same simple interpretations about the formula apply. The reason we can consider this a “distance” is that the KLD is non-negative,

$D_{KL}(P \parallel Q) \geq 0$, which is known as Gibbs inequality and only equals 0 when $P = Q$ almost everywhere. As such, the distance between P and Q is 0 when they are the “same” distribution and some positive number if they differ. For example, consider the random variables $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ with corresponding density functions given by equation 2.2; the divergence between the distributions is

$$D_{KL}(p \parallel q) = \frac{1}{2} \left(\ln \left(\frac{\sigma_2^2}{\sigma_1^2} \right) + \frac{\sigma_1^2}{\sigma_2^2} + \frac{(\mu_1 - \mu_2)^2}{\sigma_2^2} - 1 \right), \quad (2.7)$$

Letting $\mu_1 = \mu_2$, $\sigma_1^2 = 1$, and allowing σ_2^2 to vary, we can write

$$D_{KL}(p \parallel q) = \frac{1}{2} \left(\ln \sigma_2^2 + \frac{1}{\sigma_2^2} - 1 \right), \quad (2.8)$$

Figure 2.2 illustrates the KLD between the example normal distributions.

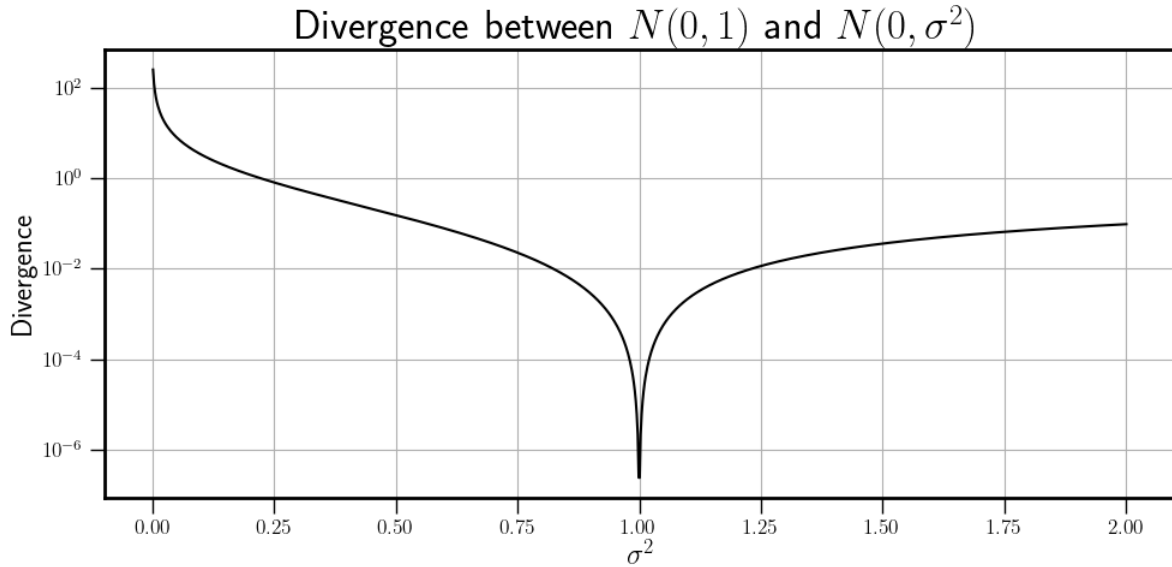


Figure 2.2. KL Divergence example with Normal distributions. As expected, the divergence decreases to 0 as $\sigma^2 \rightarrow 1$ from either direction.

2.2 How to build distributions (Non-parametric statistics)

The question remains about how we obtain the probability distribution of the random variable that represents the evolution of each layer of an ML model from measured data. Many techniques exist, but we’ll be using Kernel Density Estimation (KDE). The most basic idea behind KDE is using a different bin centering method and a smoothing factor, called a kernel, to approximate a probability density f from

measured data in the form of a histogram. This is accomplished with the *kernel density estimator* for f given by:

$$\hat{f}(x; K, h) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

where K is the chosen kernel, h is what's called the bandwidth parameter, and $\{x_i\}_{i=1}^n$ is the data that generates our histogram. Per [4], the choice of kernel does not provide much statistical significance; but the choice of the bandwidth parameter h is very crucial for finding a density estimate that approximates the underlying density function appropriately. Choosing an optimal h depends on the given data and the common method, Silverman's rule of thumb, works on the assumption that the density function of the data is unimodal and close to normal. For general data, on which no assumptions of normality are made, a more general method is found in the Improved Sheather-Jones algorithm (ISJ)[2]. The ISJ seeks to ...

2.3 Measuring entropy flow

At each epoch n of training for each layer l in the encoder \mathcal{E} and decoder \mathcal{D} networks of the DLDMD algorithm, we generate an affiliated matrix of weights $\mathbf{W}_{l,\mathcal{E}}^{(n)}$ and $\mathbf{W}_{l,\mathcal{D}}^{(n)}$, respectively. Should any given training cause the machine to converge to a fixed state, then for each layer we should have that

$$\mathbf{W}_{l,\mathcal{E}}^{(n)}, \mathbf{W}_{l,\mathcal{D}}^{(n)} \rightarrow \mathbf{W}_{l,\mathcal{E}}^*, \mathbf{W}_{l,\mathcal{D}}^* \text{ as } n \rightarrow \infty \quad (2.9)$$

Consequently, we can characterize any set of weights via the formulas

$$\mathbf{W}_{l,\mathcal{E}}^{(n)} = \mathbf{W}_{l,\mathcal{E}}^* + \mathbf{W}_{l,\mathcal{E}}^{(n),f}, \quad \mathbf{W}_{l,\mathcal{D}}^{(n)} = \mathbf{W}_{l,\mathcal{D}}^* + \mathbf{W}_{l,\mathcal{D}}^{(n),f}, \quad (2.10)$$

where the fluctuation matrices $\mathbf{W}_{l,\mathcal{E}}^{(n),f}$ and $\mathbf{W}_{l,\mathcal{D}}^{(n),f}$ are defined implicitly. These fluctuations are what we're most interested in studying. Using first-order differencing, we can study the affiliated detrended matrices $\delta\mathbf{W}_{l,\mathcal{E}}^{(n)}$ where

$$\begin{aligned} \delta\mathbf{W}_{l,\mathcal{E}}^{(n)} &= \mathbf{W}_{l,\mathcal{E}}^{(n+1)} - \mathbf{W}_{l,\mathcal{E}}^{(n)} \\ &= \mathbf{W}_{l,\mathcal{E}}^{(n+1),f} - \mathbf{W}_{l,\mathcal{E}}^{(n),f} \end{aligned}$$

The definitions for the decoder weights are identical, and thus omitted. We can consider $\delta\mathbf{W}_{l,\mathcal{E}}^{(n)}$ to be the deviation from the steady state and with it in hand we flatten it's data and use KDE to generate an affiliated empirical probability distribution $p_{l,\mathcal{E}}^{(n)}(w)$. With these distributions we now find the KLD between consecutive distributions to generate the following sequence:

$$D = \left\{ D_{KL} \left(p_{l,\mathcal{E}}^{(n+1)} \parallel p_{l,\mathcal{E}}^{(n)} \right) \right\}_{n=1}^{N_E-2} \quad (2.11)$$

which we can interpret to be the change in information between consecutive fluctuations from the steady state.

2.4 Implementation Details

The devil is in the details with respect to how this analysis is carried out and we will go over the important ones here. After the training of each model for N_E epochs, the data that we have to play with is

$$\mathbf{W}_{\mathcal{E}} = \left\{ \mathbf{W}_{l,\mathcal{E}}^{(n)} \right\}_{n=1, l=1}^{N_E, N_l}, \mathbf{W}_{\mathcal{D}} = \left\{ \mathbf{W}_{l,\mathcal{D}}^{(n)} \right\}_{n=1, l=1}^{N_E, N_l}, \quad (2.12)$$

which are the sets of weights of the layers of the encoder and decoder. Note that both sub-networks have the same number of layers and only differ by the dimension of their inputs and outputs. The following procedures are identical for both sets, so we'll only be referencing the encoder set from here on out. For each layer l , the set of detrended matrices are computed

$$\delta \mathbf{W}_{l,\mathcal{E}} = \left\{ \mathbf{W}_{l,\mathcal{E}}^{(n+1),f} - \mathbf{W}_{l,\mathcal{E}}^{(n),f} \right\}_{n=1}^{N_E-1}, \quad (2.13)$$

these matrices are flattened and their entries are used as the data for KDE. The KDE implementation used is that of KDEpy (<https://pypi.org/project/KDEpy/>, **make this a URL later**), which implements several algorithms for KDE in python including the bandwidth computation of ISJ. The kernel used is the Epanechnikov kernel, defined as

$$K(x) = \begin{cases} \frac{3(1-x^2)}{4} & |x| < 1 \\ 0 & \text{else} \end{cases}, \quad (2.14)$$

which is optimal in the mean square error sense. However, the choice of kernel is not a statistically significant one[4] and we choose this kernel for the convenience of it having finite support and non-vanishing first and second derivative on that support. The ISJ algorithm is used to determine the appropriate bandwidth for each set of data. These bandwidths are also being saved to see if the optimal bandwidth with respect to each data set might also contain some instructive information for quantification. From this, we obtain the set of density estimates for each detrended matrix

$$p_{l,\mathcal{E}} = \text{KDE}(\delta \mathbf{W}_{l,\mathcal{E}}) = \left\{ p_{l,\mathcal{E}}^{(n)}(w) \right\}_{n=1}^{N_E-1} \quad (2.15)$$

where w is the weight value itself and $p_{l,\mathcal{E}}(w)$ is the approximate density associated with that weight value. For the sake of our discretization of this data, we use 10001 linearly spaced points to represent the weights of each distribution from the smallest weight to

the largest. The number of discrete points is also of some importance for the next step; given these density approximations, we can now compute what we're truly interested in: the relative entropy of consecutive densities and the associated change in information:

$$D_l = \left\{ D_{SKL} \left(p_{l,\mathcal{E}}^{(n+1)} \parallel p_{l,\mathcal{E}}^{(n)} \right) \right\}_{n=1}^{N_E-2} \quad (2.16)$$

Unlike in the first mention of entropy flow, we instead use the SKLD instead of the KLD. This consideration simplifies our analysis by disregarding the asymmetric nature of KLD and consequences of this choice will be discussed in Chapter 4. Given the discretized nature of the densities, we turn to numerical integration to compute the integral in Equation 2.6. For the sake of an apples to apples comparison, we make the assumption that all densities have the same support and use the 10001 discretized points with Simpson's rule for the integral approximation. Each set of these divergences corresponds to a layer in the full network, and we can now use the divergences to characterize how the information is changing with respect to the steady state. Figure 2.3 shows an example of the divergence over the whole training.

We're interested in what these divergences tell us about information flow in the model's weights, and as such we seek a fitting of the data. Any overall trends are more important to our analysis than lone data points, so to avoid transient behavior skewing any fittings we ignore the first 20% of the divergences computed. The values span several orders of magnitude, so a simple fitting would certainly fail and we instead attempt a linear fit on the base 10 logarithm of the data. This corresponds to an exponential fit in the original domain of the divergences with slope of the linear trend being the approximate growth/decay rate of the information flow from the machines steady state. The slope intercept is not so important for this reason, being only a constant that modifies the order of magnitude without changing the trends, so for the time being we ignore it. In the interest of trying to identify a classifier of good training, we then find the average and variance of the slopes of the linear fits of the layers' divergence data. The backbone of our results is to compare these averages for each latent dimension to the loss curves and phase portraits from the training for each dimension.

Figure 2.3. Example of divergences for layer Dec 2 of Van der Pol model with latent dimension 8 with linear fitting and confidence bands.

CHAPTER 3

RESULTS

From experimentation by [1], we know that some sets of hyperparameters yield good reconstruction of the structure of the phase space of each of the dynamical systems in question and that some do not. For the sake of testing our method of measuring information change, we elect to focus on changing only the lifting or latent dimension of the observables we wish to use for the DMD algorithm. After conducting multiple trainings with latent dimensions ranging from $N_O = 2$ to $N_O = 10$ for the Duffing and Van der Pol oscillators, the algorithm detailed in Chapter 2 was implemented. The initial hope of our analysis is the data points that our algorithm generates will show clear trends that will allow us to easily differentiate the optimal parameter choice for the model from suboptimal or even bad parameter choices. Given the model's task of learning an imbedding and find optimal observables in the latent space, we additionally surmised that we would see the most variability in the model in what we call the edge layers: Enc in, Enc out, Dec in and Dec out. The other layers in the model certainly accomplish something as the model trains, but the dimension changes that help E-DMD along are accomplished at the edge layers, so we expect them to evolve differently than the hidden layers.

Though not explicitly part of the algorithm, we'll examine how the bandwidths generated by the ISJ method might also give us insight into optimal parameter choice as well. As a surrogate of the bin size of the affiliated histogram of the data, it can give us an idea of how small a window each possible weight can exist in.

3.1 The Duffing Oscillator

The Duffing oscillator is a non-linear second-order differential equation used to model certain damped and driven oscillators, whose general form is given by

$$\ddot{x} + \delta\dot{x} + \alpha x + \beta x^3 = \gamma \cos(\omega t) \quad (3.1)$$

There are many sets of coefficients for which many types of behavior are exhibited by the system; but for our purposes we'll only consider the following configuration:

$$\ddot{x} - x + x^3 = 0 \quad (3.2)$$

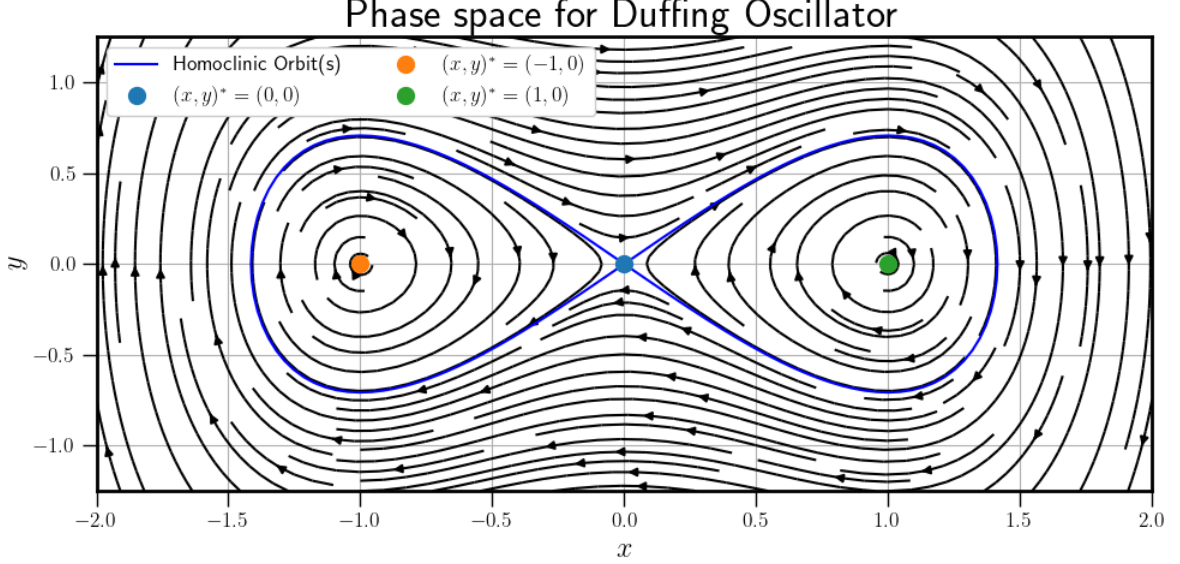


Figure 3.1. Duffing phase space.

which when turned into a system of first-order equations by introducing $y = \dot{x}$ is

$$\begin{cases} \dot{x} = y \\ \dot{y} = x - x^3 \end{cases} \quad (3.3)$$

The interesting attributes of this system are that it has 3 fixed points, 2 of which are non-linear centers with the third being a saddle point with homoclinic connections that encompass the other 2 fixed points. Figure 3.1 shows a graphical representation of the phase space of our Duffing oscillator. The homoclinic connection itself provides an interesting challenge topologically in that it divides the phase space into distinct regions that affect the end behavior of the system depending on the initial condition. Table 3.1 describes the hyperparameters used for the Duffing Oscillator that were left unchanged for each run; note that the distribution given as $\mathcal{N}(\mu, \sigma^2, n)$ is the truncated normal distribution, which is a normal distribution in which all values drawn that are more than n standard deviations from the mean are discarded.

A few notes on the trainings themselves are worth mentioning here. No learning rate scheduling is used, so the learning rate with the Adam optimizer is held constant for all epochs. The activation function Exponential Linear Unit (ELU) is used instead of the more standard Rectified Linear Unit (ReLU) due an error that occurs in the numerics of DLDMD for latent dimension choices greater than 3. ELU is a parameterized activation function whose definition is given in Table 3.2 and is closely related to ReLU; in fact, if the parameter α is set to 0 then ELU and ReLU are identical. For the sake of having similar performance for ELU with respect to ReLU, we

Table 3.1. Hyperparamters of the DLDMD algorithm held constant for each training run on Duffing.

Hyperparameter	Value
Initial Conditions	15000
Training proportion	66.66%
Validation proportion	19.98%
Testing proportion	13.36%
Optimizer	Adam
t_0	0
Δt	0.05
t_f	20
N_E	1000
Enc/Dec layers	3
Enc/Dec hidden activation	ELU
Enc/Dec output activation	identity
Neurons/layer	128
Weight initializer	$\mathcal{N}(0, .1, 2)$
Bias Initializer	0
Learning Rate	10^{-3}
a_1, a_2, a_3	1
a_4	10^{-14}

choose $\alpha = .01$. For reasons still under investigation, the RELU activation function causes the Singular Value Decomposition to fail if the latent dimension is 3 or greater and the algorithm subsequently ends with an error. After experimentation, it was found that using ELU instead circumvents this problem without significantly affecting model training. As a final note, the training for this model (and Van der Pol) across all latent dimensions was done on the same data set with the same subsets chosen to be the training, test, and validation data. We do this to reduce the number of variables down to the latent dimension alone.

Having established this, we'll consider the results of our analysis. Figure 3.2 shows the average of the slopes of the linear fits across all layers for each epoch. What we can see plainly is that for latent dimensions 2 through 5, the average slope is small in magnitude with a small to medium standard deviation relative to the scale with latent dimension 6 showing not only a much larger average slope, but as also much more variance. We may contrast this with the losses in Figure 3.3 for each latent dimension, which seems to show that latent dimension 5 should be the optimal choice with all other choices being more or less equivalent. This does not match the results of the phase space representation, which show roughly that dimension 3 is best. If we

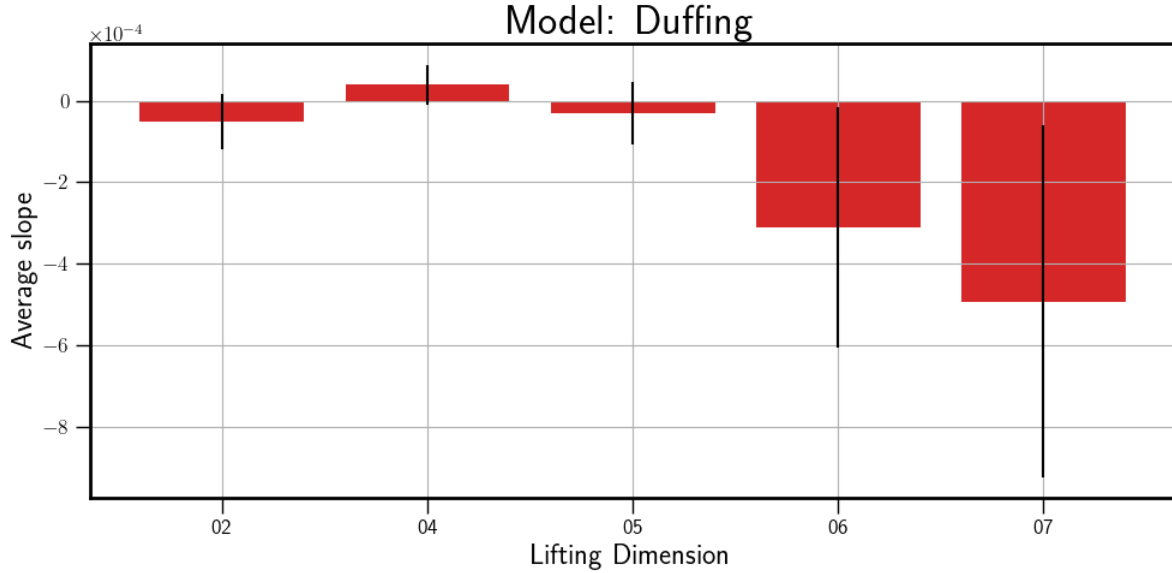


Figure 3.2. Duffing average slopes with standard deviations.

examine Figure 3.4, which shows the average bandwidth from the ISJ algorithm across all layer for a given latent dimension per epoch, we can see that 3 has a slightly different trend and separates cleanly from each of the other choices. While this behavior tracks with our understanding that this dimension is best and should be easily differentiated from the rest, it is inconsistent the fact that the other bandwidth curves are banded together but don't have much in common with respect to Figures 3.2 and 3.3 for both the fitted slopes and loss curves.

For a more granular view of the individual fittings for each layer of the model, we have Figures 3.5 and 3.6. Much to our expectation, the plots for the edge layers have slightly different behavior than for the hidden layers; generally, the linear trends have a smaller slope and y-intercept further from 0. We can interpret this as being indicative of having an information change trend that is further from the theoretical steady-state than for the other layers of the model. We've come to associate this with the kind of expressivity that is desirable for ML models to have in that they are able to cope with the kind of variability that real world data could exhibit. The ultimate conclusion we're able to draw is that DLDMD seems to adequately train and predict future states of Duffing for latent dimensions smaller than 6.

3.2 The Van der Pol Oscillator

The Van der Pol oscillator is a non-conservative oscillator with non-linear damping, and models a more specific kind of circuit than Duffing does. It's general

Model: Duffing

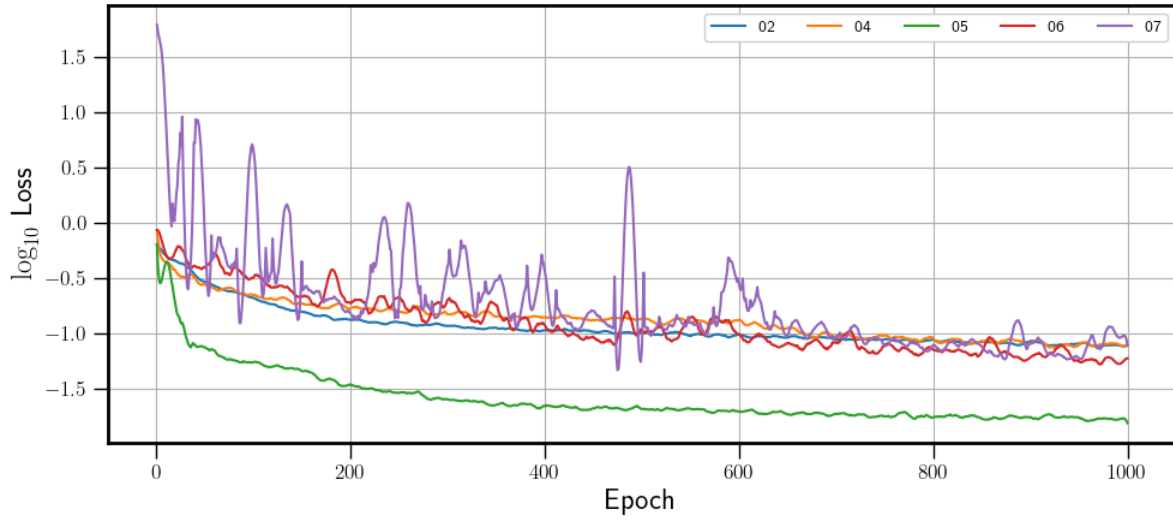


Figure 3.3. Duffing loss plots, smoothed to accentuate trends.

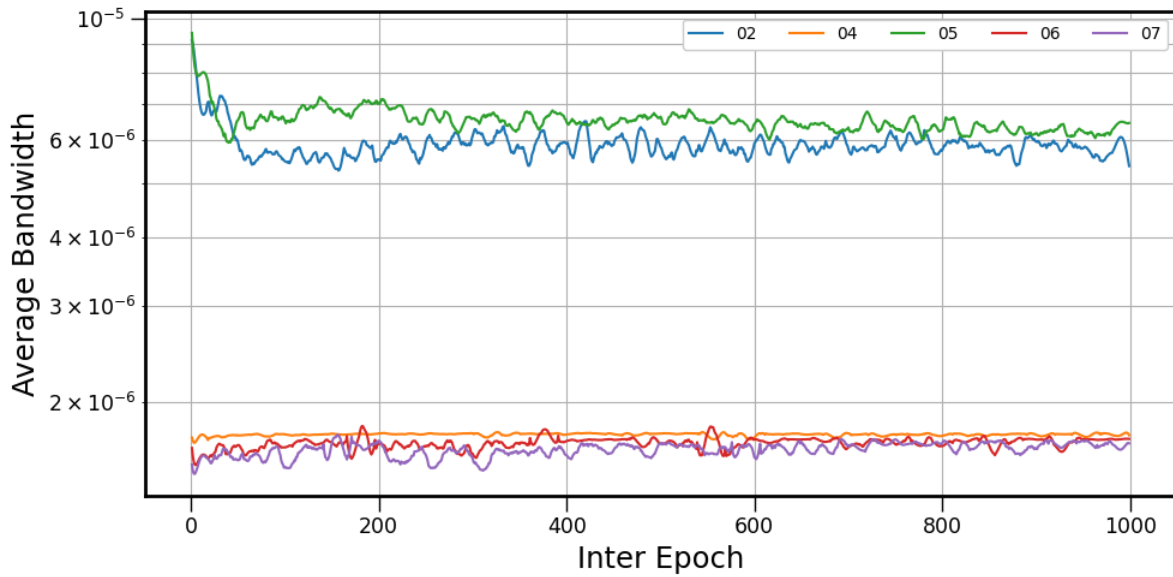


Figure 3.4. Duffing Average bandwidth.

Table 3.2. ReLU activation function vs ELU activation function.

ReLU	ELU
$\sigma(x) = \begin{cases} x, & x > 0 \\ 0, & \text{else} \end{cases}$	$\sigma(x) = \begin{cases} x, & x > 0 \\ \alpha(e^x - 1) & \text{else} \end{cases}$

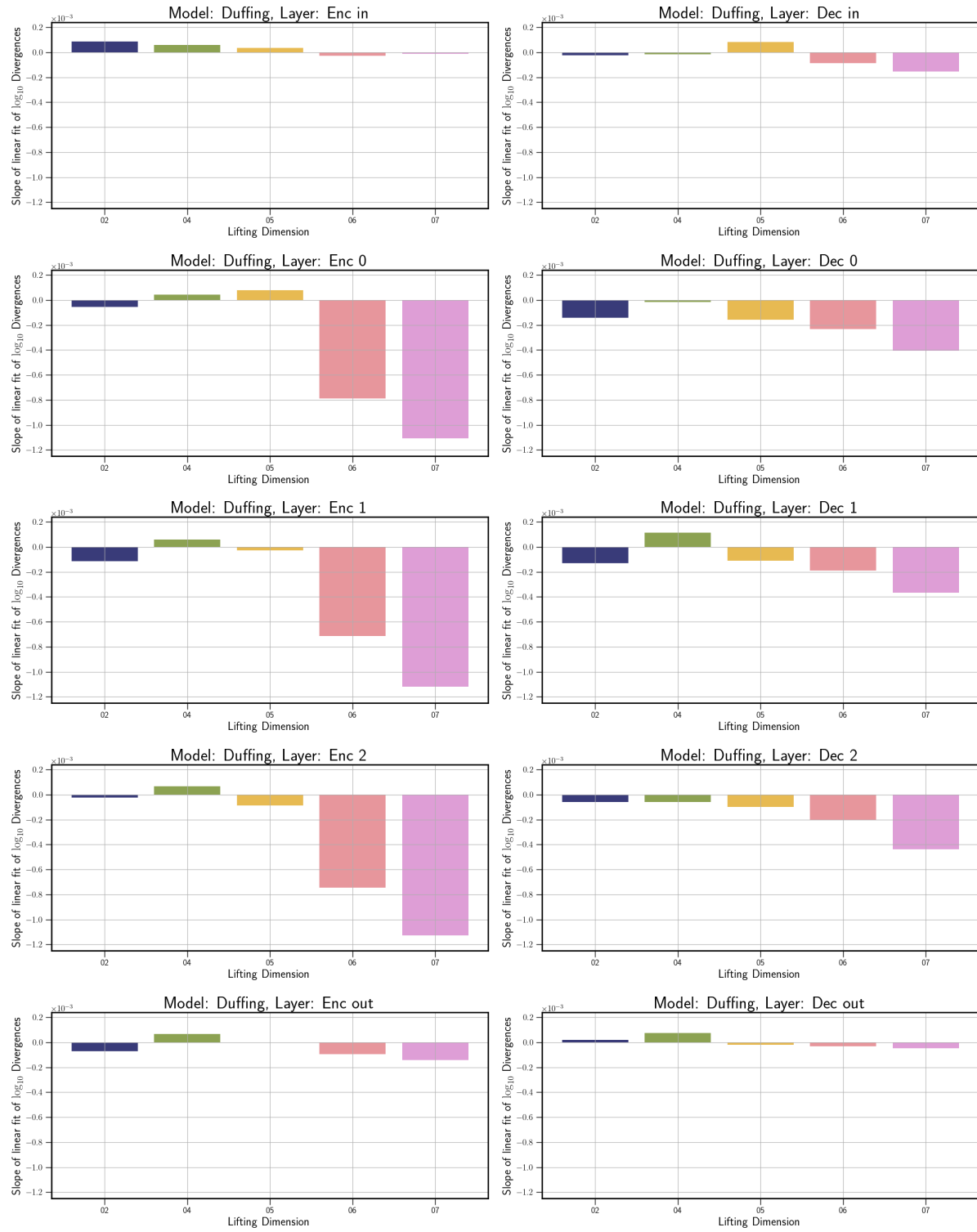


Figure 3.5. Duffing slope linear fits by layer.

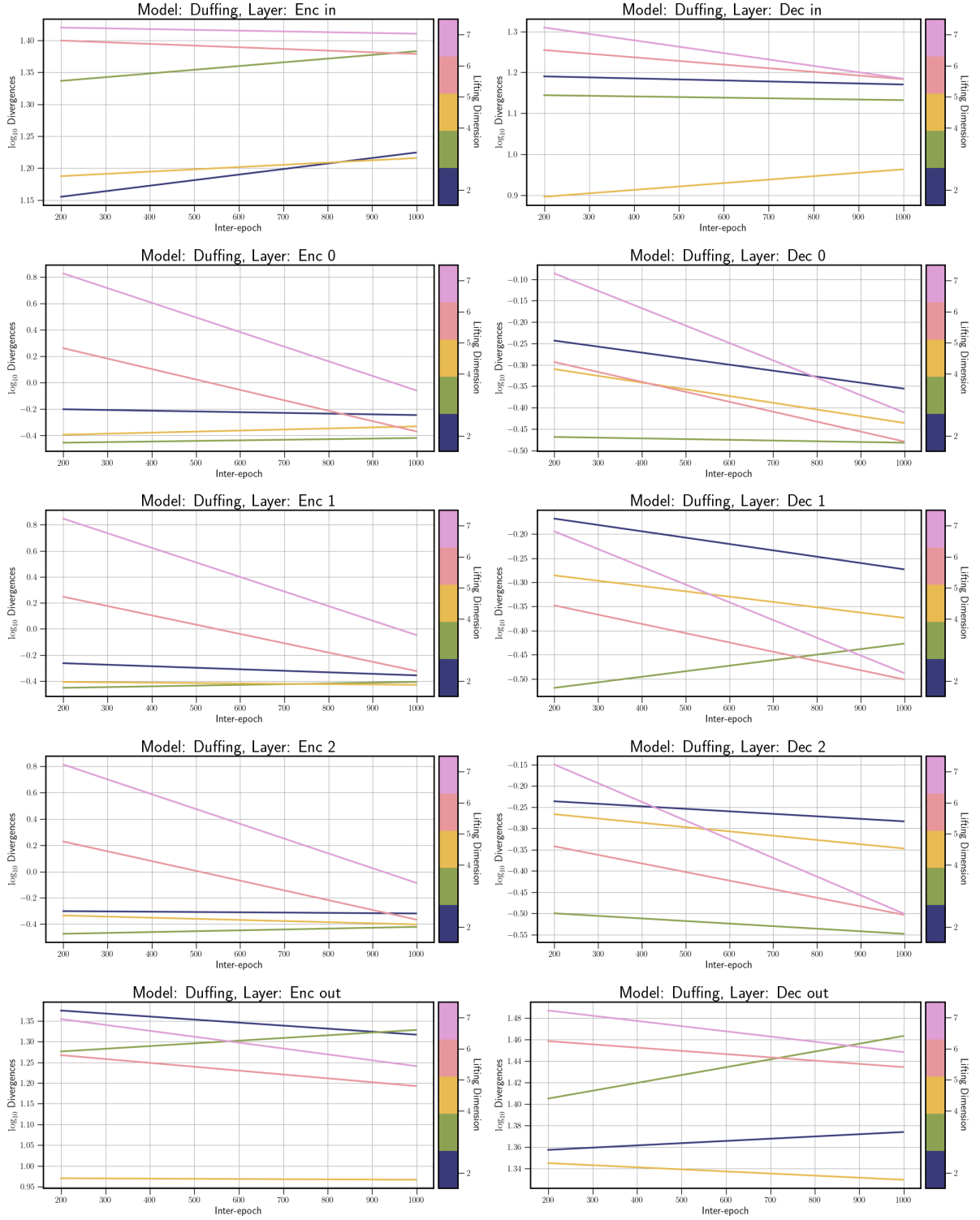


Figure 3.6. Duffing slope linear fits by layer.

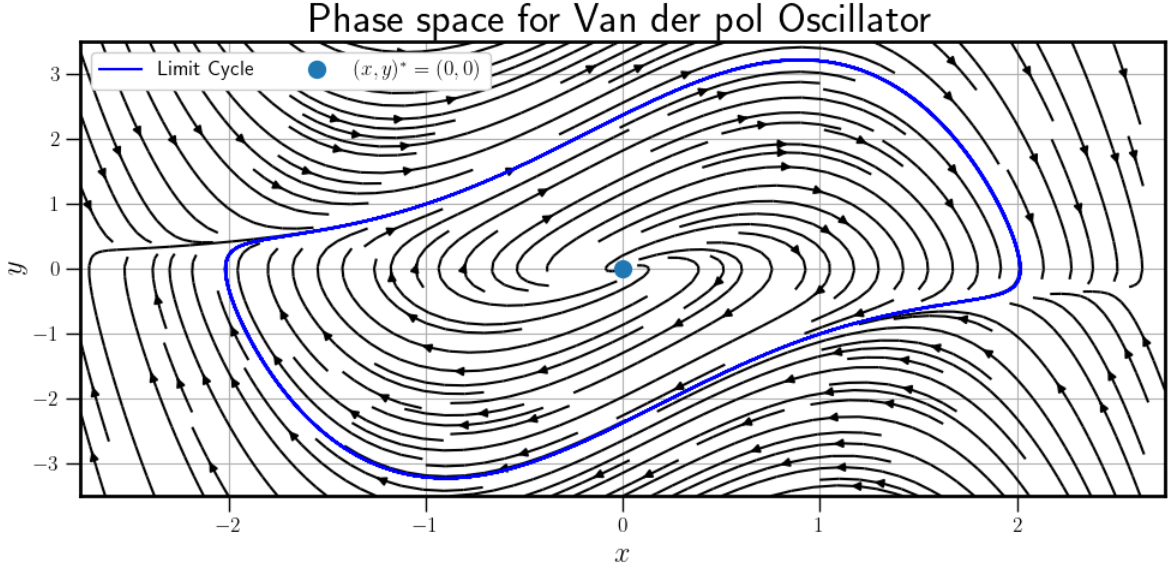


Figure 3.7. Van der Pol phase space.

form is given by

$$\ddot{x} - \mu(1 - x^2)\dot{x} + x = 0 \quad (3.4)$$

Much like Duffing, varying μ can give different and interesting phase space dynamics. For our purposes we will fix it to $\mu = 1.5$. With this, the equivalent system of first order equations is

$$\begin{cases} \dot{x} = y \\ \dot{y} = 1.5(1 - x^2)y - x \end{cases} \quad (3.5)$$

Again, by introducing $y = \dot{x}$. Unlike Duffing, this system has a single unstable fixed point at the origin and all nontrivial orbits are attracted to a limit cycle that surrounds the origin. Figure 3.7 shows a graphical representation of the phase space of our Van der Pol oscillator. Much like Duffing, we do have a curve that separates the phase space into regions of different behavior in the limit cycle; but since all orbits approach the limit cycle there should be much variability with respect to initial conditions. We would expect training on this system to be better overall due to the fact that all trajectories are essentially going to same place unlike the more complex behavior of Duffing. The hyperparameters for the Van der Pol training are identical to those of Duffing with the exception of the activation function; which is, as mentioned in Section 3.1, changed to ReLU. This difference is minor enough that a Van der Pol hyperparameters table is omitted.

Figure 3.8, shows the average and standard deviation of the slopes for the Van der Pol training per latent dimension choice. Relative to Figure 3.2 for Duffing, we have

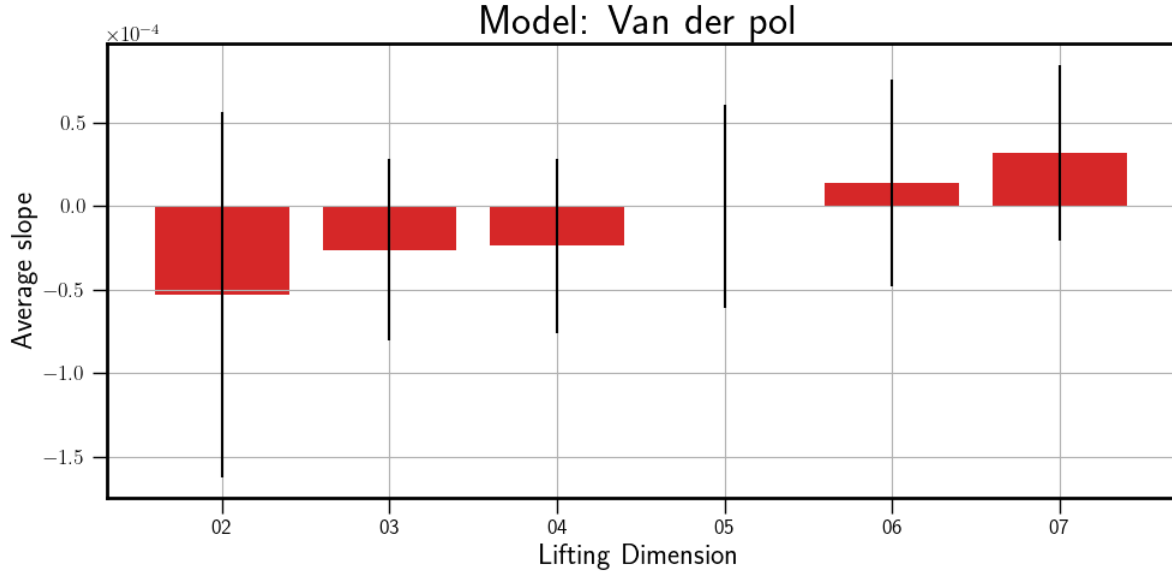


Figure 3.8. Van der Pol average slopes with standard deviations.

larger average slopes and more variety with respect to the standard deviations for each dimension. Based on the phase plane results, we should expect 8 to correspond to the best training so it seems that a small in magnitude slope is indicative of good training and this is consistent for some of the loss curves in Figure 3.9. The slope of 5 is essentially 0 and we note that it has a loss curve almost indistinguishable from that of 8; and we can make a similar observation about 7 as well. However, 10 also has a small slope but the loss curve is markedly different from the optimal choices and seems to match that of 3 which has a very large slope comparatively. One quality the optimal trainings seem to have in common is that they have smaller standard deviations than the trainings with worse loss curves.

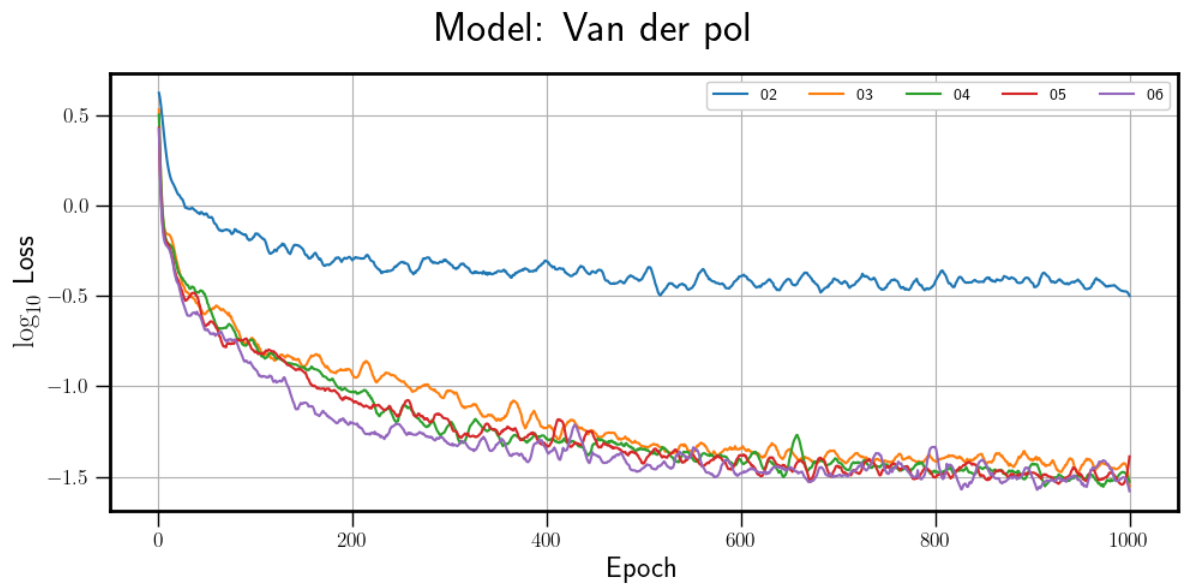


Figure 3.9. Van der Pol loss plots, smoothed to accentuate trend.

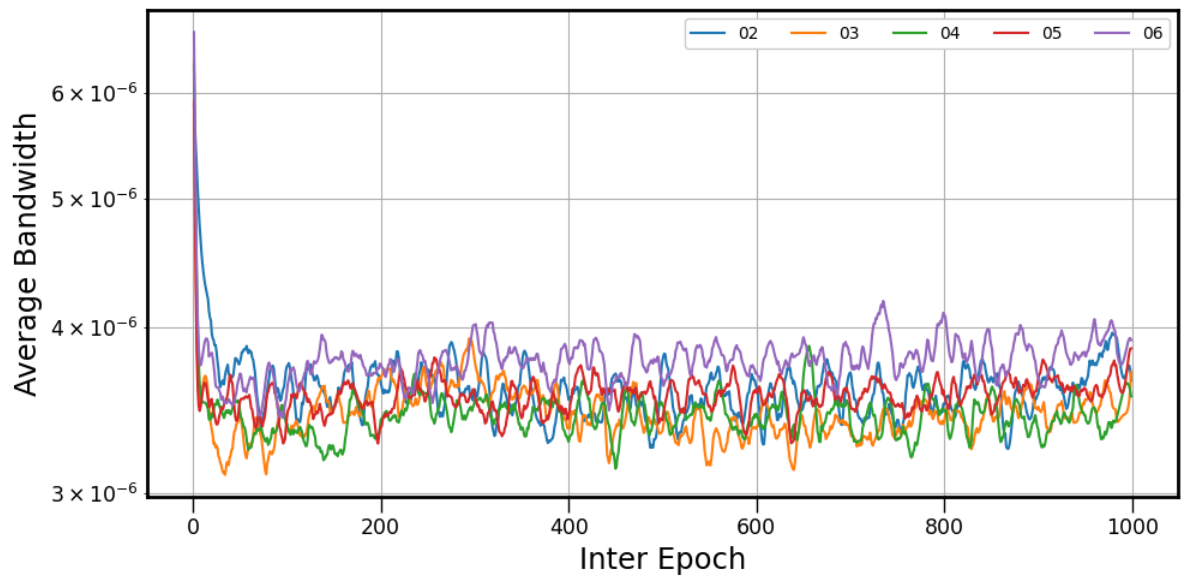


Figure 3.10. Van der Pol Average bandwidth.

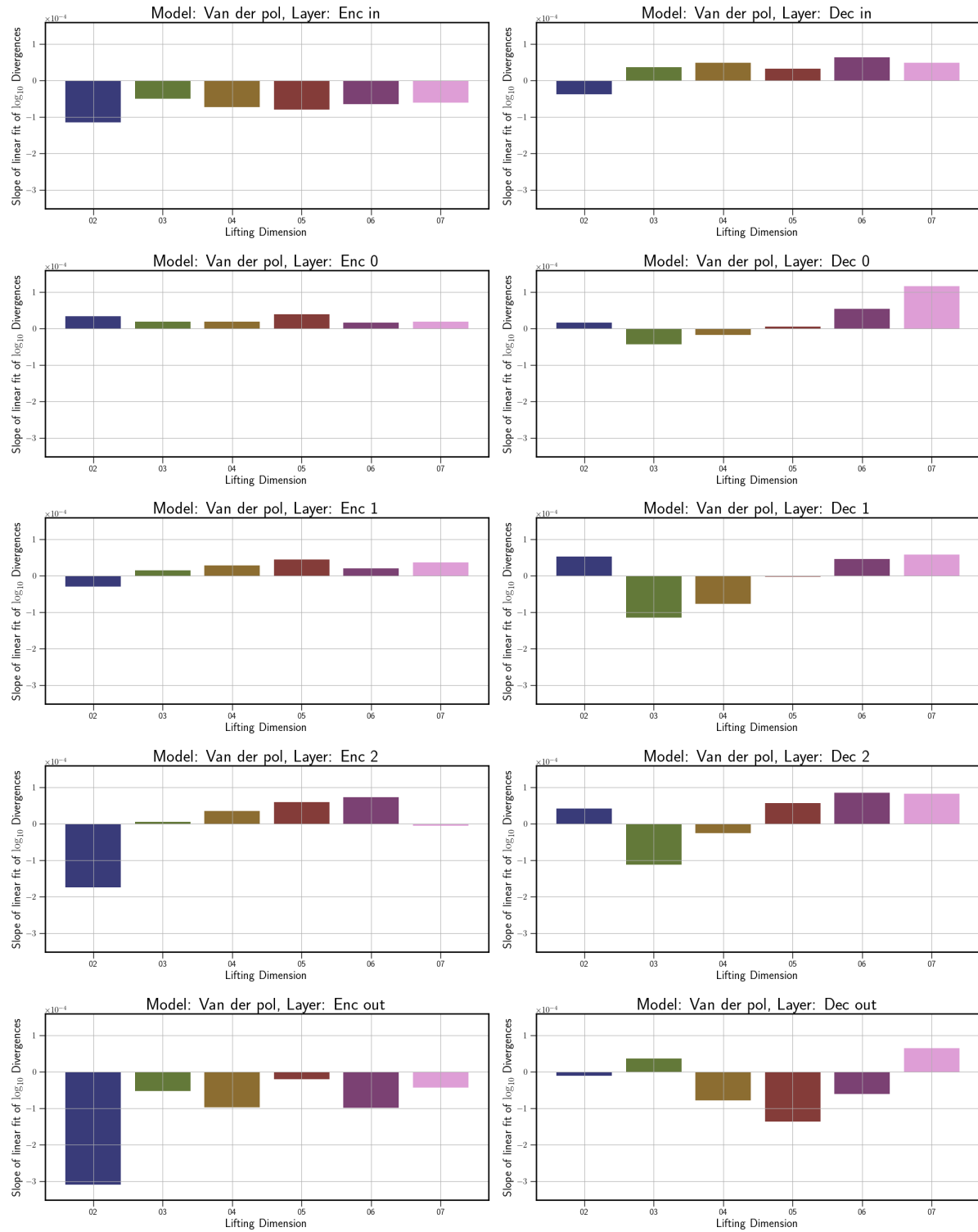


Figure 3.11. Van der Pol slope linear fits by layer.

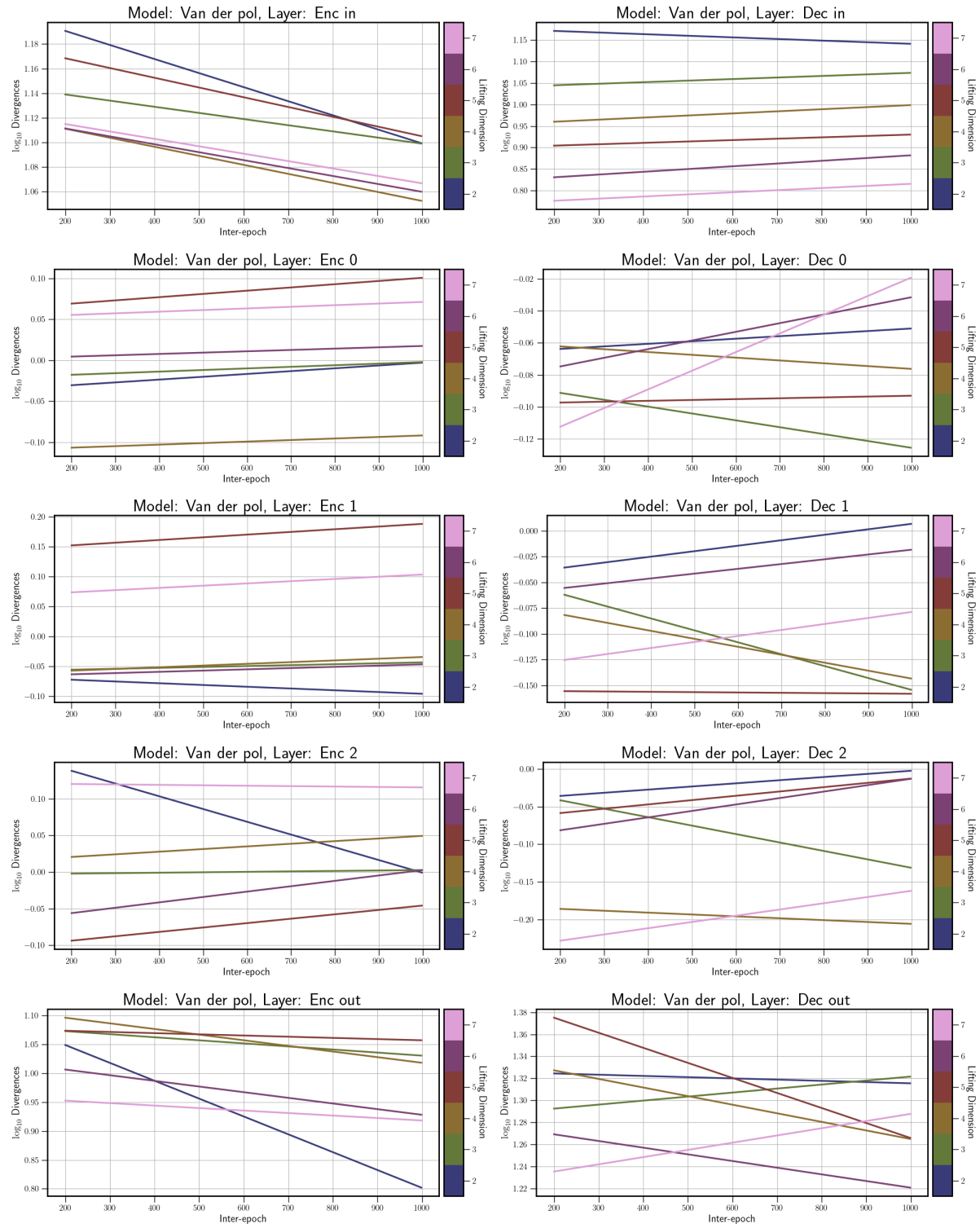


Figure 3.12. Van der Pol slope linear fits by layer.

CHAPTER 4

DISCUSSION

While we believe that we’ve struck upon some useful results with respect to finding different behavior and trends in some of the different values we’ve computed, but it must be made clear that some assumptions were made on our part that might not have been valid ones to make given what we’ve learned since the outset. Some of which is detailed in the coming sections. Perhaps the biggest assumption made is that we have sets of optimal parameters at all for the DLDMD algorithm. There hasn’t been much study into this algorithm yet, and the only measure of its validity is given by [1] in their results and our replications which admittedly rely on the “eye-ball” metric. Predicting future states in multidimensional nonlinear dynamical systems

4.1 Non-unitary densities

When creating the empirical densities in Chapter 3, often we end up deriving density functions that don’t core property of density functions: The sum of their probabilities does not integrate to 1. KDE was doubtless designed for less pathological data than the kind we seem to be extracting from the weight matrices of the model. This perhaps conflicts with the assumption that we should be able to extract a density from the weight matrices at all. This issue was rectified (though, you might say subverted) by numerically integrating the non-unitary densities and dividing by their area to force them into having the desired area. Figure 4.1 shows histograms of the numerically computed areas of the densities yielded by KDE. Peaks are visible near 1 for each model, but the spread around this average is not insignificant.

4.2 Statistical issues with edge layers

Considering that the weight matrices are used in the model for linear transformations the hidden layers have approximately D^2 entries while the edge layers have about $D \cdot d$ entries, where D is the number of neurons in the layer and d is the number of physical dimensions of the systems or the latent dimension. In either case, $d \ll D$ so from a statistical standpoint the edge layers don’t have as much data for the analysis. As we’ll elaborate on in each section of this chapter, we do notice marginally different behavior in the results for the edge layers as we do for the hidden layers and

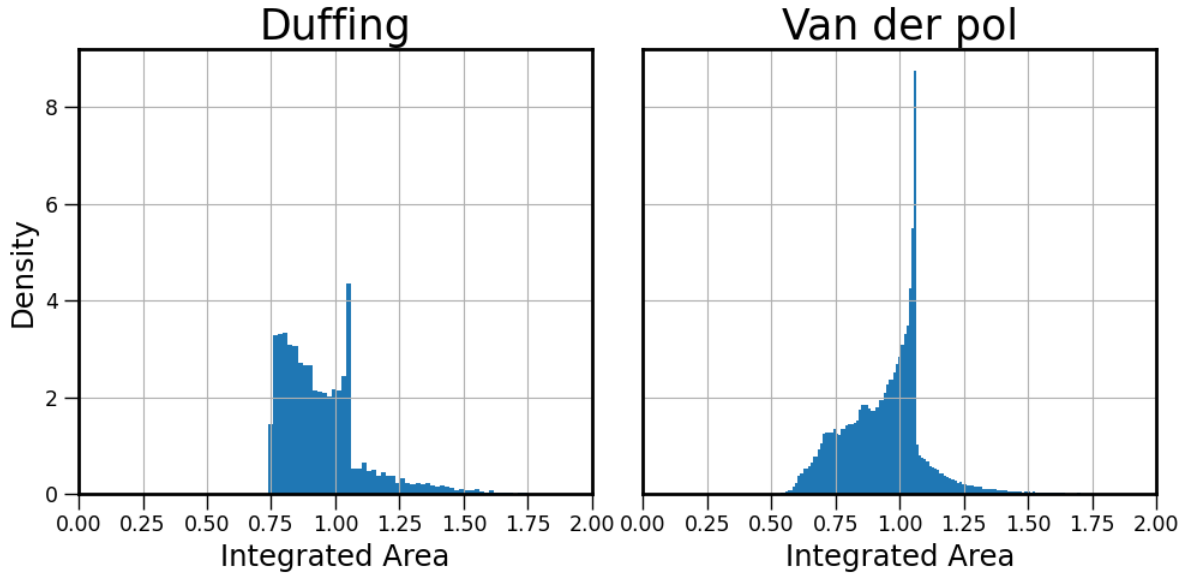


Figure 4.1. Histograms of computed areas from KDE.

it's as of yet unclear whether this is evidence of the validity of our hypothesis or a consequence of having less data to fully characterize their evolution.

4.3 Future work

In the interest of future exploration of the ideas in this thesis, another ML model that would be excellent to examine would be a classifier of some kind. Such a problem would also allow us to use alternative metrics such as model accuracy (which is much easier to quantify than the for the dynamical systems we investigate) along with our analysis to get a better idea of the limits of this discovery of good training through the frame of information change. Such classifiers are all over the modern literature of ML and would be of great interest and a method for determining good training during training would be of immense utility.

BIBLIOGRAPHY

- [1] D. ALFORD-LAGO, C. W. CURTIS, A. T. IHLER, AND O. ISSAN, *Deep Learning Enhanced Dynamic Mode Decomposition*, (2022).
- [2] Z. I. BOTEV, J. F. GROTHOWSKI, AND D. P. KROESE, *Kernel density estimation via diffusion*, The Annals of Statistics, 38 (2010), pp. 2916 – 2957.
- [3] S. BRUNTON, B. R. NOACK, AND P. KOUMOUTSAKOS, *Machine learning for fluid mechanics*, Ann. Rev. Fluid Mech., 52 (2020), pp. 477–508.
- [4] V. A. EPANECHNIKOV, *Non-parametric estimation of a multivariate probability density*, Theory of Probability & Its Applications, 14 (1969), pp. 153–158.
- [5] M. HOLLANDER, D. WOLFE, AND E. CHICKEN, *Nonparametric Statistical Methods*, Wiley Series in Probability and Statistics, Wiley, 2013.
- [6] K. HORNIK, M. STINCHCOMBE, AND H. WHITE, *Multilayer Feedforward Networks are Universal Approximators*, Neural Networks, (1989).
- [7] B. KOOPMAN, *Hamiltonian systems and transformations in Hilbert space*, Proc. Nat. Acad. Sci., 17 (1931), pp. 315–318.
- [8] B. LUSCH, J. N. KUTZ, AND S. L. BRUNTON, *Deep learning for universal linear embeddings of nonlinear dynamics*, Nature Comm., 9 (2018), p. 4950.
- [9] C. E. SHANNON, *A mathematical theory of communication*, The Bell System Technical Journal, 27 (1948), pp. 379–423.
- [10] M. WILLIAMS, I. G. KEVREKIDIS, AND C. W. ROWLEY, *A data-driven approximation of the Koopman operator: extending dynamic mode decomposition*, J. Nonlin. Sci., 25 (2015), pp. 1307–1346.