# ABSTRACT

LAGERGREN, JOHN HEIKKI. Data-Driven Discovery and Augmentation of Mathematical Models. (Under the direction of Kevin B. Flores.)

Mathematical modeling is a tool by which mathematicians and scientists can reason about the various mechanisms that drive physical and theoretical processes. A well-calibrated mechanistic model provides both *interpretability* about the underlying physical/biological laws that govern a system and *generalizability* which affords the ability to predict unseen dynamics beyond a set of observation data. There are many commonly encountered challenges that confound the capacity to conceive and calibrate these models. Many of these arise from poor understanding of the relevant dynamics, practical unidentifiability of model parameters, and the presence of uncertainties in the observed data, model parameters, and the model itself. This thesis presents a set of methodologies that leverage techniques from deep learning and data science to reduce the effects or presence of these challenges. Three cases are considered: (i) augmenting traditional mathematical modeling strategies with data-driven techniques when the governing set of equations is known, (ii) discovering mathematical models from data using deep learning and sparse regression when the governing equations are unknown, and (iii) learning the nonlinear terms comprising the governing system of equations with an end-to-end deep learning framework when only basic conservation laws are known. In (i), state space reconstruction is used to replace mechanistic equations, reducing the complexity of the inverse problem, and resulting in more accurate time-series forecasting and uncertainty quantification. In (ii), neural networks are shown to denoise and approximate the partial derivatives of data more accurately than state-of-the-art methods in the presence of biologically realistic forms and levels of noise, resulting in more robust equation discovery. In (iii), the end-to-end deep learning framework, nicknamed biologically-informed neural networks, is used to discover a previously unknown mechanism in cell biology from real-world experimental data. The results presented in this work are a step toward enabling the application of the mechanistic modeling strategy to a wider spectrum of intractable data sets arising from complex systems.

Data-Driven Discovery and Augmentation of Mathematical Models

by
John Heikki Lagergren

A dissertation submitted to the Graduate Faculty of
North Carolina State University
in partial fulfillment of the
requirements for the Degree of
Doctor of Philosophy

Applied Mathematics

Raleigh, North Carolina

2020

APPROVED BY:

_____                    _____
Hien T. Tran                                                    Ralph C. Smith


_____                    _____
Adriana San Miguel Delgadillo                        Kevin B. Flores
                                                                      Chair of Advisory Committee

# DEDICATION

To my wife, Ashlie Lagergren, and my parents, Eeva Valentine and James Lagergren.

# BIOGRAPHY

John Lagergren graduated summa cum laude with a B.S. in Computational Applied Mathematics with a minor in Physics from East Tennessee State University in December, 2015. He earned an M.S. in Applied Mathematics from North Carolina State University (NCSU) in December, 2018, and a Ph.D. in Applied Mathematics from NCSU in December, 2020. During his graduate studies, John was awarded the NCSU Graduate Fellowship in 2016-2017, the SAMSI (Statistical and Applied Mathematical Sciences Institute) Graduate Fellowship in 2018-2019, the Winton-Rose Award for Research Excellence in 2020, and had two internships at Applied Research Associates in 2019 and 2020.

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

CHAPTER

1

# INTRODUCTION

## 1.1 Problem statement

Limited understanding of the governing dynamics is an inherent issue when developing a mathematical model that approximates a dynamical system [Ban14a]. Ideally, a mathematical model is developed to achieve a balance between model complexity and the ability to parameterize the model using available data, with the ultimate goal of maximizing predictive accuracy for out-of-sample data. However, even in cases of full observability (i.e., when every variable in the model is known and observed), accurate parameterization may still be a challenge due to parameter identifiability-related issues. If parameters are not structurally identifiable with respect to an observed set of data, one can attempt to re-parameterize the model and then estimate aggregated parameters [Mes14]. However, a drawback with this technique is that if the primary goal is to accurately estimate parameters, the aggregated parameters may not be physically meaningful or interpretable. An alternative approach is to use subset selection techniques to find identifiable combinations of parameters and then fix the non-indentifiable parameters to constant values [CA09; Ban15b]. Yet, one then encounters the issue of having to justify those fixed values from other sources of experimental data in addition to ensuring that model predictions are not sensitive at the values to which the parameters are fixed.

In the event that the posited model suffers from model discrepancy or identifiability-related issues, alternative paradigms, which we refer to as "non-mechanistic models," exist that can forecast time series data in the absence of mechanistic models. These include empirical dynamical modeling [Ye15], autoregressive models [Bil13], and machine learning [Par00]. Since these methods do not rely on developing a mechanistic model based on real-world knowledge, they do not include parameters that correspond to physically/biologically interpretable quantities. In general, it has been noted that a primary drawback of utilizing non-mechanistic modeling for predicting unseen dynamics is that one forfeits the transferability and theoretical understanding afforded by a validated mechanistic model [Har13].

In the case that no mechanistic model is known, recent efforts have also investigated methods for learning underlying governing systems of equations directly from observed spatiotemporal data. Examples of such methods include the Sparse Identification of Nonlinear Dynamics (SINDy) algorithm [Bru16a] and the Equation Learning (EQL) neural network [Sah18], both of which are used for discovering systems of ODEs, and the PDE Functional Identification of Nonlinear Dynamics (PDE-FIND) algorithm [Rud17a], which is used to identify PDE systems. The goal of these methods is to specify the correct terms of the true equations given a large library of potential candidate terms (such as polynomial combinations of independent variables, function values, partial derivatives, etc.) with the aid of sparse regression or nonlinear optimization. A practical challenge arises, however, because one typically does not have access to noiseless function values or partial derivatives. Instead, one must approximate these values from noisy experimental data. Finite difference and polynomial spline-based techniques have been reported as state-of-the-art denoising methods for equation learning in a recent study which also considered several other numerical methods such as Gaussian kernel smoothing, Tikhonov differentiation, and spectral differentiation with high-frequency term thresholding [Rud17a]. However, both of these methods have been found to be prone to inaccuracies in the presence of noise [Rud17a].

Alternative equation learning approaches have used function-approximating deep neural networks as surrogate models for the solution of the governing dynamical system [Rai19; Yan20]. In these approaches, the assumed mechanistic mechanistic model is pre-specified and then used as a form of regularization in the neural network objective function. The parameters of the specified model are allowed to be "learnable," meaning that the parameters of the governing PDE are calibrated while the neural network is trained to minimize the error between the network outputs and the observed noisy data. This methodology ensures that the neural network solution satisfies the physical laws described by the specified

model while simultaneously fitting the spatiotemporal data. While these methods have been demonstrated with data in the presence of noise, they have so far only been applied to problems where the governing mechanistic PDE is known *a priori*.

## 1.2 Dissertation outline

In Chapter 2, we consider the mathematical challenge of modeling noisy experimental data using equations with practically unidentifiable parameters. We address this challenge by presenting a new method for the data-driven prediction and uncertainty quantification of multivariate systems. Traditionally, either mechanistic or non-mechanistic modeling methodologies have been used for prediction; however, it is uncommon for the two to be incorporated together. We compare the forecast accuracy of mechanistic modeling, using Bayesian inference, a non-mechanistic modeling approach based on state space reconstruction, and a novel hybrid methodology composed of the two. The method is evaluated using real-world age-structured population data from cannibalistic flour beetles, in which it is observed that the adults preying on the eggs and pupae result in chaotic non-equilibrium population dynamics. Additionally, uncertainty quantification methods for the hybrid models are outlined and illustrated for these data. We perform an analysis of the results from Bayesian inference for the mechanistic model and hybrid models to suggest reasons why hybrid modeling methodology may enable more accurate prediction of multivariate systems compared to traditional approaches.

In Chapter 3, we consider the mathematical challenge of discovering the unknown dynamics of a set of observation data when the governing system of equations is unknown. To address this challenge, we develop methods based on artificial neural networks and sparse regression for learning partial differential equation (PDE) models from spatiotemporal data under biologically realistic levels and forms of noise. Recent progress in learning mechanistic models from data have used sparse regression to select candidate terms from a denoised set of data, including approximated partial derivatives. We analyze the performance in utilizing previous methods to denoise data for the task of discovering the governing system of PDEs. We also develop a novel methodology that uses artificial neural networks (ANNs) to denoise data and approximate partial derivatives. We test the methodology on three PDE models for biological transport, i.e., the advection-diffusion, classical Fisher-KPP, and nonlinear Fisher-KPP equations. We show that the ANN methodology outperforms previous denoising methods, including finite differences and both local and global polynomial regression splines, in the ability to accurately approximate partial derivatives and learn the correct PDE model.

In Chapter 4, we consider the mathematical challenge of discovering the unknown dynamics of a set of observation data when the measurements are sparse and only basic conservation laws governing the system are known. To address this challenge, we extend scientific machine learning methods to the class of *equation learning* to be feasible for biological applications with nonlinear dynamics and where data are often sparse and noisy. Physics-informed neural networks [Rai19] have recently been shown to approximate solutions of PDEs from simulated noisy data while simultaneously optimizing the PDE parameters. However, the success of this method requires the correct specification of the governing PDE, which may not be known in practice. Here, we present an extension of the algorithm that allows neural networks to learn the nonlinear terms of the governing system without the need to specify the mechanistic form of the PDE. Our method is demonstrated on real-world biological data from scratch assay experiments and used to discover a previously unconsidered biological mechanism that describes delayed population response to the scratch.

Finally, in Chapter 5, we summarize our findings and contributions and discuss possibilities for future work.

## 1.3   Publication overview

The content of this dissertation corresponds to three articles which are at various stages of review and publication. Of these three articles, [Lag18; Lag20b], corresponding to Chapters 2 and 3, respectively, are published at this time, and [Lag20a], corresponding to Chapter 4, was accepted for publication in October, 2020. In addition, the author has six other publications and two working papers whose content is not included in this thesis.

CHAPTER

2

# DATA-DRIVEN HYBRID MODELING

## 2.1   Introduction

Mechanistic modeling strategies for predicting multivariate biological systems involve the use of dynamical models, e.g., differential equations, to describe the biological mechanisms and interactions that affect the evolution of the system [Ban14a; Smi14; Ban09a]. Applications of this strategy to genetic networks, neuronal networks, and population dynamics have enabled the prediction of complex and emergent behaviors in these systems [Ban09a]. A central challenge when utilizing a mechanistic model for prediction is the ability to accurately parameterize it from available time series data, which can often be sparse and noisy in biological settings [Ban15a; Ado15b]. Commonly encountered challenges that confound the ability to accurately parameterize a model can be attributed to problems related to model discrepancy and parameter identifiability [Vos02; CA09; Smi14]. Thus, the development of methodologies to reduce the effects or presence of these challenges may enable the application of the mechanistic modeling strategy to a wider spectrum of intractable data sets arising from complex biological systems.

Model discrepancy is an inherent issue when developing a mathematical model that approximates a biological system [Vos02; Ban14a]. Ideally, a mathematical model is developed to achieve a balance between model complexity and the ability to parameterize the

model using available data, with the ultimate goal of maximizing predictive accuracy for out-of-sample data. A general principle is to reduce the mathematical model description to the lowest dimension possible, i.e., with the least number of variables and parameters. Whereas "hold-out" validation approaches are often used to evaluate the ability of the model to predict out-of-sample data [Gei93], to the best of our knowledge, no systematic methodologies exist for minimizing model dimensionality while simultaneously maximizing prediction accuracy. Even in cases of full observability, i.e., when every variable in the model is a longitudinal covariate in the available time series data, accurate parameterization may still be a challenge due to identifiability-related issues [III15; Mes14; Cob80; Rau09]. If parameters are not structurally identifiable with respect to an observed set of data, one can attempt to reparameterize the model and then estimate aggregated parameters [Mes14]. One drawback with this technique is that if the primary goal is to accurately estimate parameters, e.g., to infer the kinetic rates of biological interactions, the aggregated parameters may not be biologically meaningful or interpretable. An alternative approach is to use subset selection techniques to find identifiable combinations of parameters and then fix the non-indentifiable parameters to constant values [CA09; Ban15b]. However, one then encounters the issue of having to justify those fixed values from other sources of experimental data and be able to ensure that model predictions are not sensitive at the values to which the parameters are fixed.

Alternative paradigms exist to forecast time series data without a mechanistic model; we refer to these approaches as "non-mechanistic models". These include empirical dynamical modeling [Ye15], autoregressive models, e.g., NARX [Bil13], and machine learning, e.g., multi-step ahead prediction [Che06; Par00]. Since these methods do not rely on developing a mechanistic model based on biological knowledge, they do not include parameters that correspond to biologically interpretable quantities, e.g., kinetic rates. Thus, if the goal is to estimate biologically meaningful parameters from time series data, the mechanistic modeling strategy is a reasonable first approach. In general, it has been noted that a primary drawback of utilizing non-mechanistic modeling in forecasting is that one forfeits the transferability and theoretical understanding afforded by a validated mechanistic model [Har13]. While these concerns have been previously noted, our focus here is to present a hybrid strategy that leverages the advantages of both mechanistic and non-mechanistic modeling to maximize predictive accuracy and minimize forecast uncertainty.

We merge two well-known methods, (i) state space reconstruction and (ii) Bayesian inference, to investigate whether their combination could minimize the drawbacks encountered when utilizing each method separately. The state space reconstruction (SSR) methodology relies on Takens' theorem of delay embedding and uses time series data to

generate a manifold that is one-to-one with the attractor manifold of a dynamical system [Far87; Cas89; Sug90; Smi92; Jim92; Sau94; Sug94; Sch98; Kug98; Yua04; Hsi05; Str06; Reg05; Sch06; Ham16]. In theory, since the reconstructed manifold is one-to-one with the attractor of the real system, one can use it to forecast future dynamics using a nearest neighbor approach as described in Section 2.2.3. A critical limitation of using SSR for prediction is the amount of data needed to accurately reconstruct the attractor manifold. Since no biological knowledge is leveraged, SSR, similar to other non-mechanistic modeling approaches, can require a large amount of data to build a purely data-driven representation of the underlying dynamical system. This attribute can be especially limiting in biological scenarios for which data are collected at sparse time points. On the other hand, Bayesian inference methods have been widely applied in modeling of biological systems with this level of data [Smi14]. However, Bayesian inference relies on fitting parameters for a mechanistic model, and therefore is also subject to the previously discussed modeling related issues.

Here we describe a hybrid implementation of SSR and Bayesian inference methodologies in which we reduce model dimensionality by systematically omitting system variables and replacing them with either data or SSR predictions. To validate our methodology, we consider a real-world biological data set consisting of 21 time series of cannibalistic flour beetle, (*Tribolium casteneum*), population dynamics [Con97]. It was previously found that combining SSR with mechanistic models enabled more accurate predictions of chaotic systems, including the flour beetle data set [Ham17]. Our goal in this work is to provide a hybrid methodology that quantifies the uncertainty, both for the model predictions and estimated parameters, in addition to forecasts of future time series data. We also provide a deeper investigation of the hybrid approach than in previous efforts by analyzing the uncertainty quantification results. We discuss our analysis and suggest possible reasons why the hybrid approach may yield more accurate predictions.

## 2.2   Methods

### 2.2.1   Age-structured population data

We use longitudinal data of total counts for larvae, pupae, and adults in flour beetle populations. The data come from 7 different experimental conditions in which adult mortality rates were altered, resulting in non-equilibrium dynamics; 3 replicates were performed in each condition for a total of 21 data sets [Con97]. Data were sampled every other week over an 82 week period for a total of 41 data points per time series. To test our methodology under noisy observation conditions similar to ecological systems, we added normally distributed

random observation error to each time series using a coefficient of variation (CV) of 0.2, which is consistent with reported noise levels in survey data [Fra03; Per13]. The data from one experiment, shown in Figure 2.1 as black ×'s, exemplify the typical non-equilibrium time series behavior of the beetle system. As denoted by the vertical line in Figure 2.1, each time series is divided into a training set (first 32 time points) and a testing set (last 9 time points). The training set is used for either Bayesian inference or SSR and the testing set is used to evaluate the accuracy of the considered models.



**Figure 2.1** Empirically recorded population levels of adult flour beetles over a period of 82 weeks with normally distributed random observation error added. Data to the left of the vertical bar are used for parameter fitting while the data to the right are left for forecast validation.

### 2.2.2 Mathematical model

We use the previously validated discrete-time age-structured model

$$L(t) = b A(t-1) e^{-c_{el} L(t-1) - c_{ea} A(t-1)}, \tag{2.1}$$

$$P(t) = L(t-1)(1-\mu_1), \tag{2.2}$$

$$A(t) = P(t-1) e^{-c_{pa} A(t-1)} + A(t-1)(1-\mu_a), \tag{2.3}$$

for flour beetle population dynamics. The total number of larvae, pupae, and adults at time

8

$t$, are denoted by $L(t)$, $P(t)$, and $A(t)$, respectively. One unit of time is equal to 2 weeks, which matches the time scale of the data. This model quantifies the stage progression of beetles from the larval to pupae stage, and pupae to adult stage. Adult larvae are reproductive and contribute to the recruitment rate in equation (2.1). The exponential terms in equations (2.1) and (2.3) respectively represent cannibalization of larvae by adults or larvae, and cannibalization of pupae by adults. A more thorough description of the model and parameters can be found in [Con97]. For reference in the following material we note that the parameters $c_{pa}$ and $\mu_a$ are assumed to be experimentally known; see [Con97] for further details.

### 2.2.3 State space reconstruction

For non-mechanistic prediction we choose *direct prediction*, a state space reconstruction technique based on Takens' method of delayed embedding [Far87; Cas89; Sug90; Smi92; Jim92; Sau94; Sug94; Sch98; Kug98; Yua04; Hsi05; Str06; Reg05; Sch06; Ham16]. We summarize direct prediction (denoted throughout as SSR) here, and refer the reader to the supplemental of [Sug12] for a more in depth description of the methodology.

Let the $i$-th state variable of the system be denoted by $Y_i$. Assume we have $N+1$ observations of $Y_i$, namely $\{Y_i(t)\}_{t=0}^N$, which we refer to as the training data of $Y_i$. We start by building $d+1$ dimensional delay-coordinate vector $Y_i^d(N) = [Y_i(N), Y_i(N-\tau), Y_i(N-2\tau), \ldots, Y_i(N-d\tau)]$ where $d$ is the number of delays and $\tau$ is the time delay. The general idea of SSR is to use the $K$ nearest neighbors, found within the training data in delay-coordinate space, to predict $Y_i$ beyond the end of the training set. As such, SSR begins by building a library of delay-coordinates from the training data. The nearest neighbors to the current delay-coordinate vector $Y_i^d(N)$,

$$\begin{aligned}
Y_i^d(t') &= [Y_i(t'), Y_i(t'-\tau), Y_i(t'-2\tau), \ldots, Y_i(t'-d\tau)] \\
Y_i^d(t'') &= [Y_i(t''), Y_i(t''-\tau), Y_i(t''-2\tau), \ldots, Y_i(t''-d\tau)] \\
&\vdots \\
Y_i^d(t^K) &= [Y_i(t^K), Y_i(t^K-\tau), Y_i(t^K-2\tau), \ldots, Y_i(t^K-d\tau)]
\end{aligned}$$

are found as a function of Euclidean distance, and the known $Y_i(t'+P), Y_i(t''+P), \ldots, Y_i(t^K+P)$ are used in a local model to generate prediction $\hat{Y}_i(N+P)$.

Here, our local model is a weighted average of the nearest neighbors

$$\hat{Y}_i(N+P) = w_1 Y_i(t'+P) + w_2 Y_i(t''+P) + \ldots + w_K Y_i(t^K+P),$$

where $w_j$ is the weight of the $j^{th}$ nearest neighbor defined as

$$w_j = \frac{e^{-(d_j/\sigma)^2}}{\sum_{j=1}^{K} e^{-(d_j/\sigma)^2}}. \tag{2.4}$$

Here, $d_j$ is the distance of the $j^{th}$ neighbor to $Y_i^d(N)$ and $\sigma$ is a bandwidth parameter we set as $\sigma = 2$.

Additionally, we adopt the methodology in [Ye15] for computing the variance of the SSR prediction. The weights $w_j$ describe the probability of selecting the $j$th neighbor, where we assume that the variance of the prediction is given by $\text{Var}(\hat{Y}_i(N+P)) = \mathbf{E}[(Y_i(t^j+P) - \hat{Y}_i(N+P))^2]$. For the results presented here, in building our delay-coordinates we set $\tau = 1$, which corresponds to the sample rate of the observed data, and use $d = 1$ delays. 5 nearest neighbors are used to build the local model for non-mechanistic prediction.

### 2.2.4   Bayesian inference

We performed Bayesian inference using a delayed rejection adaptive metropolis (DRAM) algorithm implemented in MATLAB [Haa06]. Parameter values to initialize the parameter chains were generated using a weighted least squares algorithm, see Section 3.2.3 of [Ban14a], for each data set and model, i.e., full model or hybrid model, separately. We used the following lower and upper bounds for each parameter: the initial conditions are given by $L_0 \in [-50, 550]$, $A_0 \in [-50, 550]$, $P_0 \in [-200, 500]$ and the model parameters are $b \in [-5, 30]$, $c_{el}, c_{ea} \in [-0.02, 1]$, and $\mu_1 \in [-1, 1]$. These bounds were set by initially choosing an interval +/-50% around previously estimated parameter values from [Den95] and then increasing the size of the interval until the tails of each parameter's posterior distribution was contained completely within the interval. We used flat prior distributions defined between the upper and lower bounds for each parameter. We set the chain length to 20,000 with a burn-in length of 20,000. We performed uncertainty quantification, i.e., computation of 95% prediction and credible intervals by sampling from posterior distributions as described in [Haa06; Smi14].

### 2.2.5 Hybrid methodology

Our approach for combining SSR and Bayesian techniques is split into a training and testing phase, where the latter phase evaluates prediction accuracy. We describe the procedure for generating the hybrid prediction model, illustrated in Figure 2.2, and note that the procedure for generating other hybrid models is similar.

In the training phase, Bayesian inference uses a subset of the dynamical systems model by dropping one or more of the variables. For example, a hybrid model can be constructed for inference of parameters in the equation for the $A$ variable by using a partial model consisting of only equation (2.3) and substituting training data for the remaining $L$ and $P$ variables as depicted in Figure 2.2. More precisely, since (2.3) contains terms involving the variable $P(t)$ describing the number of pupae, we simply replace this variable by the actual number of pupae observed at time $t$ up to the last training time point at $t = 32$. This enables (2.3) to be used with Bayesian inference to estimate parameters for the $A$ equation. In this case the only parameter to estimate is the initial condition $A_0$, since $c_{pa}$ and $\mu_a$ are experimentally known. We note that SSR is not utilized for Bayesian inference during the training phase. Instead, SSR is used to train an equation-free prediction for the L and P variables.

In the testing phase, time series for the unmodeled variables after the training data are forecasted using SSR predictions. Similar to the training phase, unmodeled variables are substituted by their SSR predictions inside equations for modeled variables. For example, once the parameter $A_0$ is estimated for (2.3), we can generate a prediction with this parameterized equation by continuing to substitute the SSR prediction for $P$ for the missing time series $P(t)$ beyond the training data.

Using this approach, there are a total of 7 models corresponding to subsets of the variables for the full model $\{L, P, A\}$. For example, we may choose to use equations (2.1) and (2.2) to model the $L$ and $P$ variables, and then use either training data or SSR predictions as a substitution for the $A(t)$ time series.

### 2.2.6 Hybrid uncertainty quantification

Our methodology for uncertainty quantification in hybrid models is to make a modification to the input of the forward solution when computing prediction and credible intervals with the DRAM algorithm [Haa06; Lai11]. We will illustrate the method on the hybrid model using an equation for the $A$ variable and data/SSR for the $L$ and $P$ variables, as illustrated in Figure 2.2, and note that computations are similar for other hybrid model choices. We note that, in this scenario, a posterior distribution is obtained for the parameter $A_0$, and

**Figure 2.2** Illustration of the hybrid technique applied to modeling the *A* variable in the flour beetle system given by equations (2.1)-(2.3). The depicted hybrid model uses a partial model for *A*, i.e., only equation (2.3), and training data for *L* and *P* to estimate parameters in the partial model. In this case, the only parameter to estimate for (2.3) is the initial condition *A*(0). SSR predictions are used to continue substituting time series for *L* and *P* after the last time point in the training data, thereby merging nonparametric with parametric predictions.

the equation $A(t) = P(t-1)e^{-c_{pa}A(t-1)} + A(t-1)(1-\mu_a)$ is used to model the *A* variable. The method for computing prediction and credible intervals for times $t \leq 32$ is unchanged from the DRAM algorithm outlined in [Haa06], since we assume that the replacement of $P(t)$ with training data at times $t \leq 32$ are an exact model for the pupae population. At times $t \geq 33$, the trajectory of the forward solution for *A* is affected by the uncertainty in the SSR prediction for *P*. Adopting the notation from Section 2.2.3, the SSR uncertainty is given by the probability distribution generated by the sample space of nearest neighbors $\Omega$, where the trajectory of each neighbor is selected with probability $w_j$, given by equation (2.4). Thus, prediction and credible intervals are obtained for the hybrid model for $t \geq 33$ by sampling from the joint density of the posterior distribution for $A_0$ and the sample space of SSR nearest neighbor trajectories.

## 2.2.7   Evaluation

The standardized root mean square error (SRMSE) is used to quantify prediction accuracy for the testing data, e.g., the last 9 time points in Figure 2.1. The SRMSE quantifies the mean of the squared error over all of the training data, and thus represents an accuracy score at each of the 9 time points predicted, aggregated over the entire 21 time series in the experimental data set. We note that the SRMSE normalizes the prediction score with respect to the standard deviation of the training data. Thus, SRMSE < 1 implies that the prediction is better than using the mean of the training data, referred to as "naive prediction", to forecast future time series. It is expected that, in the long-term, prediction accuracy will converge to

SRMSE $= 1$ since only short-term prediction is possible in chaotic systems.

## 2.3   Results

We performed parameter estimation and uncertainty quantification for the full model in equations (2.1)-(2.3) and every hybrid model option corresponding to each subset of the variables $\{L, P, A\}$. For clarity, we will refer to the full model, which does not use SSR, as "LPA". The hybrid models will be referred to as SSA, SPS, SPA, LSS, LSA, and LPS where L, P, and A denote modeled variables and S denotes unmodeled variables replaced with data and SSR predictions. For example, SSA corresponds to modeling the $A$ variable with Eq. (2.3) and using data or SSR predictions for the $L$ and $P$ variables. Table 2.1 indicates the model used for predicting $L$, $P$, and $A$, either data/SSR or one of equations (2.1)-(2.3), corresponding to each of the seven choices. Parameter estimates, parameter densities, and prediction intervals for each model/dataset combination can be found in Appendix A or downloaded in the Open Science Framework [Osf].

**Table 2.1** Models used to predict the $L$, $P$, and $A$ variables for the various hybrid choices. Either SSR or one of equations (2.1)-(2.3) if used for each variable. LPA corresponds to the full model, while the first six choices represent the possible hybrid models. The corresponding estimated parameters are also indicated for each hybrid model.

| Choice | $L$ | $P$ | $A$ | Parameters |
|---|---|---|---|---|
| SSA | data/SSR | data/SSR | Eq. (2.3) | $A_0$ |
| SPS | data/SSR | Eq. (2.2) | data/SSR | $P_0, \mu_1$ |
| SPA | data/SSR | Eq. (2.2) | Eq. (2.3) | $P_0, A_0, \mu_1$ |
| LSS | Eq. (2.1) | data/SSR | data/SSR | $L_0, b, c_{el}, c_{ea}$ |
| LSA | Eq. (2.1) | data/SSR | Eq. (2.3) | $L_0, A_0, b, c_{el}, c_{ea}$ |
| LPS | Eq. (2.1) | Eq. (2.2) | data/SSR | $L_0, P_0, b, c_{el}, c_{ea}, \mu_1$ |
| LPA | Eq. (2.1) | Eq. (2.2) | Eq. (2.3) | $L_0, P_0, A_0, b, c_{el}, c_{ea}, \mu_1$ |

### 2.3.1   Forecast accuracy

We evaluated the forecast accuracy of the full model (LPA), each hybrid model (SSA, SPS, SPA, LSS, LSA, LPS), and SSR-only model for each of the $L$, $P$, and $A$ variables, which are shown in Figure 2.3. The points in Figure 2.3 show the standardized root mean square errors (SRMSEs) averaged across all 21 data sets. These quantify the prediction accuracy for the testing set as

detailed in Section 2.7. The vertical bars denote the standard errors at each point. A smaller SRMSE is interpreted as having a better predictive ability. Just using the mean of the training data as a constant model to forecast would yield an SMRSE equal to 1. We found that the full model outperformed the SSR method for the $L$ variable for up to 12 weeks of prediction as shown in the left panel of Figure 2.3. Comparisons between the full model and the SSR method were less clear for the $P$ and $A$ variables. For example, in the short term, at the first time point of prediction, SSR outperformed the full model for the $P$ variable, but the full model outperformed SSR at later time points, e.g. forecast horizon at weeks 6-12 as shown in the middle panel of Figure 2.3. We observed that the hybrid method, corresponding to at least one of the hybrid models from the six choices, was able to outperform both the full model and the SSR method for each of the $L$, $P$, and $A$ variables in the first 2-3 time points of prediction and was comparable to the full model in subsequent time points of the forecast horizon. These results are similar to previous work when combining SSR methods with Kalman filtering techniques for parameter estimation [Ham17]. Focusing on the comparison between the hybrid models and the SSR method, we found that hybrid models were able to outperform SSR and stay below a mean SRMSE of 0.8 for up to 10 weeks of prediction in the forecast horizon. These results indicate that the hybrid models are the most accurate choice for predicting future time series as compared to the full model or SSR alone. This is of particular interest since by replacing modeled variables with data and SSR predictions, the hybrid approach, while being a more complicated methodology, is actually a simpler mathematical model in the context of inverse problems. In the following sections we perform an analysis to investigate several reasons why the hybrid method simplifies the parameter estimation task. To the best of our knowledge, such an analysis has not been carried out for a hybrid modeling approach.

### 2.3.2   Uncertainty quantification

Here we tested our proposed methodology for uncertainty quantification with hybrid models on the flour beetle data set and model. The top of Figure 2.4 illustrates the full model fit to experimental data for one time series; only the $A$ variable is shown for comparison to the hybrid model that does not use the $L$ or $P$ variables. This instance is representative of a scenario in which the full model prediction is approximately equivalent to the mean of the training data, and thus performs no better than the naive prediction. In contrast, the bottom of Figure 2.4 shows the hybrid model for the same experimental data. In addition to the error between the testing data interval being lower, the hybrid model model also has narrower 95% credible intervals than the full model, indicating higher confidence in the

**Figure 2.3** Forecast errors (SRMSE) for each forecasting method for the larvae (L), pupae (P), and adult (A) population data. Points are the mean SRMSE over 21 data sets; bars are the standard errors. The full model corresponds to model choice LPA for each population. The hybrid models for the larvae, pupae, and adult populations corresponds to model choices LSA, SPS, and SPA, respectively. These specific hybrid models shown here had the smallest SRMSE values over all 21 data sets.

predicted population size at any given time point.

### 2.3.3 Practical identifiability analysis

The results from Bayesian inference on a model and time series data can be used to produce other outputs besides credible and prediction intervals. Two outputs in particular, pairwise parameter plots and the fisher information matrix (FIM), provide key information about the practical identifiability of parameters in the model, i.e., the ability to estimate parameters with reasonably low uncertainty levels from the available data in the presence of noise. Thereby, analysis of correlation plots and the FIM provides insight into difficulties encountered in the parameter estimation task. We analyzed parameter correlations and the FIM for the full model and hybrid models to test the hypothesis that hybrid models can simplify the parameter estimation task by maximizing the level of data information content with respect to the set of estimated parameters. For example, even though estimation for the full model (LPA) uses 3 time series, corresponding to the $L$, $P$ and $A$ variables, and the estimation for the hybrid SSA only uses a single time series, corresponding to the $A$ variable, it is not immediately obvious that the gain in the amount of data afforded by using the full model justifies the additional number of parameters that need to be estimated.

15

**Figure 2.4** Prediction with uncertainty quantification for the $A$ variable using a hybrid model (bottom) and full model (top). Data (black x's) are from experiment 1 for which $c_{pa}$ was experimentally set to zero. The vertical black line separates the training data used for parameter estimation from the testing data used for evaluating forecast accuracy. The 95% credible and prediction intervals are shown as dark and light grey, respectively. The black line represents the mean of the credible interval (best fit).

Moreover, it is unclear how the presence of observation noise or the structure of the model may independently or synergistically affect the balance between model dimension and the amount of data used for parameter estimation.

In Figures 2.5 and 2.6 we illustrate the parameter correlation information that is output from applying Bayesian inference to the full model, LPA, and one of the hybrid models, SPA, on one of the 21 experiments in the data set. Briefly, parameter correlation information can be seen graphically by plotting the values obtained in the Markov chain Monte Carlo (MCMC) chain used to build the joint posterior distribution of the estimated parameters. For example, if there are $p$ parameters to be estimated, then a $p$-dimensional vector $(\hat{\theta}_{1,k}, \ldots, \hat{\theta}_{p,k})$ of estimated parameters is obtained at step $k$ in the chain. If the number of chain iterations used to construct the posterior distribution is equal to $M$, then a parameter pair correlation plot for the $i$-th vs. $j$-th parameter is made by plotting $M$ pairs of points $(\hat{\theta}_{i,k}, \hat{\theta}_{j,k})$ for $k = 1, \ldots, M$. Figure 2.5 shows that correlations exist when using the

full model between the pairs $(A_0, P_0)$, $(b, P_0)$, $(b, L_0)$, $(c_{el}, L_0)$, $(c_{el}, b)$, $(c_{ea}, P_0)$, $(c_{el}, b)$, $(\mu_1, L_0)$, $(\mu_1, b)$, and $(\mu_1, c_{el})$. In contrast, Figure 2.6 shows that no correlations exist for the hybrid model among the three possible parameter pairs when using SPA, i.e., data/model for the $L$ variable, and (2.2) and (2.3) for the $P$ and $A$ variables, respectively.



**Figure 2.5** Pairwise parameter plots for the full model (LPA). The data from experiment 1 were used for parameter estimation.

To investigate the presence of parameter correlations among all 21 time series and all model choices, we computed the linear correlation coefficients ($\rho$) between all parameter pairs for any given model as plotted in Figure 2.7. We found that, with the exception of one pair ($b$ versus $c_{ea}$), on average the correlation coefficients were lower for any hybrid model compared to the full model (LPA) and that the full model had many coefficients

**Figure 2.6** Pairwise parameter plots for a hybrid model (SPA). The data from experiment 1 were used for parameter estimation

close to 1. We note that correlation, such as that observed between $\mu_1$ and $b$, does not necessarily indicate that parameters are nonidentifiable unless the Pearson coefficients are close to ±1 or pairwise plots are single valued. However, the variability in pairwise plots can be influenced by the level of observation noise, 20% in this case, and the correlation coefficients, shown in Figure 2.7, are close to 1.

This motivates the analysis of the rank deficiency of the FIM. The FIM has been previously used in subset selection algorithms that seek to predict which subsets of parameters are identifiable for a given model and available set of data. Importantly, these FIM based methods take into account the sensitivity of the model output with respect to parameters and combines this information with the effect of parameter correlations. For example, insensitive parameters are more difficult to identify from data since a large change in the parameter doesn't affect the model output, and in particular, doesn't affect how well the model fits the data. If the number of estimated parameters is given by $p$, then the FIM is a $p \times p$ matrix, and the rank of the FIM can be used to estimate the number of parameters that are practically identifiable [Olu13]. The FIM is equal to $\sum_{i=1}^{N} \chi^T(t_i)\chi(t_i)$, where the matrix $\chi(t_i)$ contains sensitivities of the model with respect to parameters at time point $t_i$ in the training data. The $k, j$-th entry of $\chi(t_i)$ is given by $\{\frac{\partial y_k(t_i)}{\partial \theta_j}\}$ where $y_k$ is the $k$-th observable, e.g., $L$, $P$, or $A$, and $\theta_j$ is the $j$-th model parameter. We used the parameters estimated from Bayesian inference to compute the FIM for each of the 21 time series and model choices SSA, SPS, SPA, LSS, LSA, LPS, and LPA.

We determined the rank of the FIM by first computing the singular value decomposition (SVD), i.e., FIM $= USV^T$, where $S$ is a diagonal matrix of the singular values of the FIM

**Figure 2.7** Correlation coefficients among all pairs of estimated parameters. Computations were performed for each model choice where the choices are 1=SSA, 2=SPS, 3=SPA, 4=LSS, 5=LSA, 6=LPS, and 7=LPA. Coefficients were computed only for pairs of parameters that existed among the hybrid models. Each star within each subplot represents a correlation coefficient for a single time series, with 21 total possible time series.

listed in decreasing order, and $U$ and $V$ are orthogonal matrices containing left and right singular vectors. Since $S$ is a diagonal matrix, it can be viewed as a list $\{s_1, ..., s_p\}$. We used the location in the list, if any, where the ratio $\frac{s_m}{s_{m+1}} > 10^{10}$ to indicate that the rank of the FIM was equal to $m$. Thus, the rank deficiency is given by $p - m$, and the number of parameters that are not practically identifiable increases as a function rank deficiency. We found that the full model (LPA) had the largest average rank deficiency over all the 21 time series compared to any of the hybrid models, shown in Figure 2.8. Among the hybrid models, LPS had more rank deficiency than LSS and LSA, and each of choices LSS and LSA had more rank deficiency than choices SSA, SPS, and SPA. Hybrids SSA and SPS had no rank deficiency and SPA had only a rank deficiency of 2 for a single data set out of the 21 total data sets that were analyzed.

**Figure 2.8** Rank deficiency for each model choice. The full model is denoted by LPA and hybrid models are denoted by SSA, SPS, SPA, LSS, LSA, LPS. The y-axis indicates the number of time series (out of 21 total) that a particular choice resulted in a certain level of rank deficiency given by the x-axis.

## 2.4 Summary of contributions

We introduced a hybrid methodology composed of mechanistic Bayesian inference and non-mechanistic state space reconstruction for modeling multivariate real-world experimental data. The results in this work suggest hybrid modeling can maximize predictive accuracy by reducing the dimensionality of the inverse problem (i.e., replacing problematic modeled variables with non-mechanistic SSR predictions). This was further corroborated using correlation plots and rank deficiency of the Fisher information matrix to show that, by eliminating equations with parameters that may not be practically identifiable, the remaining parameters can be more accurately estimated with the available data. Finally, we demonstrated how to quantify the uncertainty of the hybrid predictions in the form of 95% prediction and credible intervals by jointly sampling from the posterior distributions of the modeled parameters and the nearest neighbor sample space for unmodeled variables.

## 2.5 Discussion

We illustrated a hybrid methodology that may be used for minimizing the dimensionality of parameter inference from time series data while also maximizing predictive accuracy of the resulting parameterized model that is broadly applicable to multivariate systems. Since this balance is often the goal when developing predictive models of biological systems for which data are typically sparse by having either low frequency or few time points relative to the number of estimated parameters, we hypothesize that the methodology illustrated here may enable the application of the mechanistic modeling paradigm to a wider range of biological scenarios for which data limitations or high dimensionality of the system inhibits accurate parameterization. In particular, in cases where system variables are replaced by interpolated data during parameter estimation (an example of this is the *forcing function method* which has seen applications in physiological modeling [DiS13; DG00]), the SSR/Bayesian hybrid can serve as a more advanced alternative, giving the opportunity to substitute unmodeled variables with data while simultaneously having the ability to infer their dynamics beyond the recorded data. An important feature of the combination of Bayesian inference with SSR is that uncertainty quantification, i.e., the computation of 95% prediction and credible intervals, is readily obtained by jointly sampling from the posterior distribution of the estimated parameters and the SSR sample space defined by the nearest neighbor prediction. In practice, the ability to ascertain uncertainty in predictions is a necessary feature for applying the hybrid methodology to real-world scenarios. For example, it would be clearer to make an ecological management decision based on population densities forecasted by a hybrid model if one can also associate a level of confidence with those predictions. We note that hybrid modeling using a Kalman filter approach instead of Bayesian inference was able to achieve similar levels of improvement in forecast accuracy and parameter estimates over non-parametric or parametric modeling. The main contrast between using a Kalman filter and the Bayesian approach presented here is that a Bayesian hybrid model inference allowed us to readily achieve uncertainty quantification for hybrid models. Additionally, results from Bayesian inference enabled the analysis of why the hybrid modeling scheme could be more accurate than the traditional parametric approach. Investigating how uncertainty quantification can be ascertained under a Kalman filter hybrid modeling framework is left to future work.

In addition to illustrating uncertainty quantification for hybrid models applied to a real biological data set, our intent in combining Bayesian inference with SSR in this work was to also use the results from Bayesian inference to investigate the mechanisms by which the hybrid technique may alleviate some of the practical identifiability issues that

21

commonly arise in difficult parameter estimation problems. We analyzed the correlation plots and found that for every pair of parameters, with the exception of $b$ vs. $c_{ea}$, the average correlation among all 21 time series was lower for any choice of hybrid model compared to the full model. We note that both $b$ and $c_{ea}$ are parameters located in the equation for the $L$ variable, see (2.1), suggesting that estimating parameters for this variable may be the source of difficulty in general for the flour beetle system. Under this hypothesis, it is expected that hybrid models containing the $L$ variable will have the same or increased level of parameter correlations as the full model, since the hybrid models will use less data for bayesian inference by removing either or both of the equations and corresponding time series for $P$ and $A$.

Our analysis of the rank deficiency of the FIM corroborates the finding that the equation for the $L$ variable is problematic for parameter estimation. For clarity, we show the different model choices and the corresponding parameters contained within each of them in Figure 2.9. The equation for the $L$ variable contains four parameters, while the $P$ variable contains two parameters, and the $A$ variable only contains one parameter. The rank deficiency histograms in Figure 2.8 indicate that the source of parameter non-identifiability stems from the $L$ variable. For example, the only hybrid model choice with one modeled variable (SSA, SPS, or LSS) that has any rank deficiency is LSS, which corresponds in part to using an equation for the $L$ variable and data/SSR for the $P$ and $A$ variables. This deficiency can be mitigated by including the equation for the $A$ variable in the hybrid model, since it only contains one parameter, and we see that LSA ($L$ and $A$ variables modeled with equations) has lower average rank deficiency among the entire experimental data set than LPA. We observed that the same mitigating effect is not present for LPS ($L$ and $P$ variables modeled with equations). These findings are in agreement with the forecast accuracy results in Figure 2.3; the hybrid model choices used in these plots are those with the best forecast accuracy and these choices also have the lowest rank deficiency. Together with our parameter correlation analysis, these results suggest that one likely mechanism by which hybrid modeling increases forecast accuracy is by eliminating variables for which parameters may not be practically identifiable from the available data and replacing their inaccurate estimation with non-mechanistic model-based forecasting.

**Figure 2.9** Venn diagram illustration of the model parameters estimated for each choice. Grey boxes contain the model choice, where LPA corresponds the full model, and SSA, SPS, SPA, LSS, LSA, and LPS represent the hybrid models.

CHAPTER

# 3

# NEURAL NETWORKS AND SPARSE REGRESSION

## 3.1  Motivation

In the previous work, we outlined a general hybrid strategy for circumventing parameter estimation issues in mechanistic modeling by replacing problematic model equations with data-driven techniques. However, the hybrid method (i) assumes the *a priori* knowledge of the set of governing equations that describe a dynamical system and (ii) relies on the use of state space reconstruction (SSR) for non-mechanistic modeling, which remains to be tested in broader contexts. In the event that the mechanisms that drive the governing dynamics are partially (or completely) unknown, one must resort to using an SSR-only model which lacks the physical/biological interpretability that a mechanistic model can provide. Other non-mechanistic approaches exist as alternatives to SSR. For example, machine learning methods, such as neural networks, have successfully been used in forecasting longitudinal data within a multi-step-ahead prediction framework [Che06; Par00]. Autoregressive models are also commonly used in statistical forecasting from longitudinal data and have some similarities to SSR, e.g., predictions of the future state are based on a non-mechanistic model of the recent history [Bil13]. However, similarly to SSR, these methods do not provide the

interpretability and out-of-sample generalizability of well-calibrated mechanistic models. This motivates the use of *equation learning*, a recent field of study that employs data-driven techniques to discover parsimonious mechanistic models directly from observation data. Though constantly evolving, there are two dominant approaches to equation learning in the present moment: (i) sparse regression and (ii) theory-informed neural networks. Below, we consider the extension of (i) using function-approximating neural networks.

## 3.2 Introduction

Recent research has investigated methods for discovering systems of differential equations that describe the underlying dynamics of spatiotemporal data. There are key advantages to learning and then using mathematical models for prediction instead of using a purely machine learning based method, e.g., neural networks. First, if the learned mathematical model is an accurate description of the processes governing the observed data, it has the ability to generalize from the set of training data to data outside of the training domain. Second, the learned mathematical model is interpretable, making it informative for scientists to hypothesize the underlying physical or biological laws governing the observed data. Examples of recent methods for inferring the underlying governing equations include the Sparse Identification of Nonlinear Dynamics (SINDy) algorithm [Bru16a] and the Equation Learner (EQL) neural network [Mar16; Sah18], both of which are used for discovering systems of ordinary differential equations (ODEs), and the PDE Functional Identification of Nonlinear Dynamics (PDE-FIND) algorithm [Rud17a], which is used to identify PDE systems. Boninsegna et al. [Bon18] recently extended the SINDy algorithm to recover stochastic dynamical systems. Model selection criteria (such as Akaike Information Criteria and Bayesian Information Criteria) have been combined with the SINDy algorithm to increase robustness to errors, although incorrect models were still selected at noise levels we consider here [Man17]. The discovery methods mentioned above assume that the measured data arise from a parameterized $n$-dimensional dynamical system of the form

$$u_t(x,t) = F(x,t,u,u_x,u_{xx},\ldots,\theta), \qquad x \in [x_0, x_f], \quad t \in [t_0, t_f] \qquad (3.1a)$$

$$u(x,t_0) = u_0(x), \qquad x \in [x_0, x_f] \qquad (3.1b)$$

that describes the time evolution of some quantity of interest, $u(x,t)$, with parameter vector $\theta \in \mathbb{R}^k$ and appropriate boundary conditions.

Discovering a general function, $F$, from noisy data for $u$ is an active area of research. The PDE-FIND approach from [Rud17a] builds a large library of candidate terms,

$$\Theta = \begin{bmatrix} 1 & u & \cdots & u^p & u_x & \cdots & u^p \odot u_x & u_{xx} & \cdots & u^p \odot u_{xx} & u_x^2 & u_x \odot u_{xx} & u_{xx}^2 \end{bmatrix}, \quad (3.2)$$

whose columns represent potential terms comprising $F$ and $\odot$ represents element-wise multiplication. Identifying $F$ can now be re-cast as a sparse regression problem for the following linear system

$$u_t = \Theta \xi \quad (3.3)$$

where the unknown vector $\xi$ is a sparse vector whose nonzero entries correspond to the terms predicted to comprise $F$. In this work, we will assume that the true dynamical system, $F$, is only comprised of a few simple terms (e.g., for the diffusion-advection equation, $F = D u_{xx} - c u_x, D, c \in \mathbb{R}$), so that a simple library with $p = 2$ and up to second order derivatives is sufficient to discover $F$. Even in this limited setting, a practical challenge arises because one typically does not have access to the noiseless values of $u(x, t)$ or its partial derivatives. Instead, one must approximate these values from noisy experimental data. Here, our goal is to investigate the performance of existing denoising methods (to limit the amount of noise in $u_t$ and $\Theta$) that are used in conjunction with PDE-FIND and to present a novel denoising methodology relying on artificial neural networks (ANNs). We note that in this work we assume that the denoising methods considered sufficiently reduce the noise levels in order to neglect the so-called "error-in-variables" problem, i.e., where the covariates on the right hand side of the sparse regression equation (3.3) are inaccurately observed or contain noise. The error-in-variables problem may lead to biased estimates in high-dimensional sparse regression settings when the covariate noise is large [Ros10]. In future work, we will explore the use of alternative sparse regression methods that are able to construct unbiased estimates when there is additive noise in the observed covariates [Dat17; Loh12; Che13].

Several methods have been used for denoising data to approximate $u(x, t)$ and its partial derivatives ($u_t$, $u_x$, $u_{xx}$, etc.). The most prevalent methods that have been proposed are (i) finite difference approximations or (ii) the use of cubic splines for interpolation followed by partial differentiation of the fitted splines. However, both of these methods have been found to be prone to inaccuracies in the presence of noise [Rud17a]. Recent work has considered recovery of dynamical systems with high amounts of noise added to the time derivative

measurement ($u_t$) by transforming the data into a spectral domain [Sch17]. Zhang et al. [Zha18] recently proposed using sparse Bayesian regression, which allows for error bars for each candidate term in the discovered equation. However, although their method was robust, the noise in this study was also added only to the time derivative term ($u_t$) instead of the observed data ($u(x, t)$). Importantly, it has been noted that introducing noise to the observed data itself ($u(x, t)$) hinders the recovery of the correct PDE, thus, developing a method of denoising data for $u(x, t)$ has been identified as a current challenge for learning PDEs [Sch17]. To the best of our knowledge, the use of finite differences or utilizing splines are the two methods that, in practice, yield the most accurate approximations for the library terms involved in PDE learning. The primary challenge involved with using these methods for numerical differentiation, which we further test in this work, is that they are sensitive to noise levels and can amplify noise as the order of the derivative increases. This challenge inhibits learning PDEs for practical biological applications where data may have large noise levels due to many sources of error, including the data collection process, imprecise measurement tools, and the inherent stochastic nature of biological processes [Ban14a; Cod08]. For example, for ecological measurements of population abundance, typical data sets can have noise levels on the order of a coefficient of variation equal to 0.2 [Fra03; Per13]. Notably, adding this biologically relevant level of noise to the observation $u(x, t)$ has not been considered in previous PDE learning work.

In this work, we are concerned with the performance of various denoising methods in recovering $u$ and its derivatives from noisy data. There is much theory on approximating noiseless continuous functions, e.g., the Stone-Weierstrass Theorem ensures that any compactly-supported and continuous function can be uniformly approximated arbitrarily well with polynomial functions. Similarly, Hornik [Hor91] proved that a single hidden layer ANN can approximate a continuous function and its derivatives arbitrarily well under some reasonable assumptions. How such methods perform in the presence of noise is not as well understood. Evaluating the performance of these methods in the presence of noise is a crucial task for equation learning methods, as poor derivative estimation hinders these methods' ability to uncover the correct underlying equations. We thus perform a systematic investigation in this work on the accuracy of these different denoising strategies in (i) estimating $u$ and its derivatives, and (ii) learning equations from these computations in the presence of varying amounts of noise. We include finite difference and local spline computations in this study because of their successful use in recent equation learning work which used lower amounts of noise than is considered here [Rud17a]. We investigate the use of ANNs as an alternative approach to these methods since they have not been considered previously. Because ANNs must be fit to an entire set of spatiotemporal data

(affording a global context that may help to decrease overfitting), we also consider a global spline method for a comparison between different global methods.

An additional, yet realistic, complication that has not been considered is the presence of non-constant error noise in the spatiotemporal data used for PDE learning. For example, proportional error noise can occur when the variance of the data is proportional to the size of the measurement, e.g., population size or density [Ban11]. Non-constant error noise may also occur when the observed processes occur on different time scales [Joh14]. In the scenario that one has a mathematical model for the biological process generating the data, e.g., $u(x, t)$, the non-constant error noise can be accounted for with a statistical model used in conjunction with the mathematical model [And74]. For example, for a set of observed data at space points $x_i$, $i = 1, \ldots, M$ and time points $t_j$, $j = 1, \ldots, N$, a general statistical model is given by

$$U_{i,j} = u(x_i, t_j) + w_{i,j} \odot \mathcal{E}_{i,j}, \tag{3.4}$$

where the noiseless observations are corrupted by noise modeled by the random variable $w_{i,j} \odot \mathcal{E}_{i,j}$. Finite difference methods assume $w_{i,j} \odot \mathcal{E}_{i,j} = 0$ while regression methods using splines often assume the variance of $w_{i,j} \odot \mathcal{E}_{i,j}$ is constant. More generally, the error term $\mathcal{E}_{i,j}$ may instead be generated by a probability distribution that is weighted by

$$w_{i,j} = \left( \beta_1 u_1^\gamma(x_i, t_j), \ldots, \beta_n u_n^\gamma(x_i, t_j) \right)^{\mathrm{T}} \tag{3.5}$$

for $\gamma \geq 0$ and $\beta_1, \ldots, \beta_n \in \mathbb{R}$. Constant error noise is modeled by assuming $\gamma = 0$ and $\beta_1, \ldots, \beta_n = 1$. Proportional error noise is modeled by assuming $\gamma > 0$, $\beta_1, \ldots, \beta_n \neq 0$ [Ban09b]. We hypothesize that the assumption of constant variance error leads finite difference and spline approximations to yield poor estimates of the noiseless data $u(x, t)$ and its partial derivatives when the data contain proportional error noise. In this work, we investigate this hypothesis and develop a methodology using ANNs as a model for $u(x, t)$ in conjunction with an appropriate statistical error model that accounts for the presence of proportional error when denoising spatiotemporal data.

The denoising methods present in this work focus on spatiotemporal data for learning PDEs, however the methods we describe can be readily applied to learning ODEs. We choose to focus our study on a specific set of diffusive PDE models, which have provided a wealth of insight into many biological transport phenomena, including ecological migra-

tion and invasion [Has05], neuronal transport [Jon84], cancer progression [Bal14; Ste15; Roc15; Rut17], and wound healing [Mai04; Nar16]. We demonstrate that an ANN can be used with a non-constant error statistical model to accurately approximate $u(x,t)$ from noisy proportional error data better than finite difference and spline methods. We further demonstrate that the PDE-FIND algorithm can more accurately infer the governing PDE equations from data when its library of terms is constructed using an ANN-based method than these other methods.



**Figure 3.1** The two components to learning PDEs from data. The first component is to approximate $u$, $u_t$, $u_x$, $u_{xx}$, etc., from noisy data. The second component uses the output from the first component as an input for the PDE-FIND algorithm to learn a PDE. We also employ a pruning algorithm after the PDE-FIND step, not depicted here.

## 3.3   Methods

The process of learning a system of equations from noisy data can be divided into two main components: (1) the data denoising and library construction component, in which the underlying dynamical system $u(x,t)$ and its partial derivatives are approximated from the noisy realizations of $\left\{ U_{i,j} \right\}_{i=1,j=1}^{M,N}$ in (3.4), and (2) the equation learning component, in which, given approximations for $u$, $u_t$, $u_x$, $u_{xx}$, etc., one employs an algorithm that can effectively uncover the mechanistic form of $F$ in (3.1a) (Figure 3.1). Below, we describe the mathematical models used for data generation, the method of constructing a library from noisy data, and equation learning. All of the denoising methods were implemented in Python 2.7 using the Scipy package for polynomial splines and the Keras machine learning library for ANNs. All code and accompanying animations are available at `https://github.com/biomathlab/PDElearning/`.

### 3.3.1 Data generation

We consider three diffusive PDE models for biological transport in this work, each of which has been used previously to interpret biological data [Jin16b; Mai04; She90a; Sib99]. These models include the diffusion-advection equation

$$u_t = -c\,u_x + D\,u_{xx}, \quad D, c \in \mathbb{R} \tag{3.6}$$

the classical Fisher-Kolmogorov-Petrovsky-Piskunov (Fisher-KPP) Equation

$$u_t = D\,u_{xx} + r\,u - r\,u^2, \quad D, r \in \mathbb{R} \tag{3.7}$$

and the nonlinear Fisher-KPP Equation

$$u_t = D\,u\,u_{xx} + D\,u_x^2 + r\,u - r\,u^2, \quad D, r \in \mathbb{R} \tag{3.8}$$

where $D$ is the diffusion coefficient, $r$ is the intrinsic population growth rate, and $c$ is the advection rate.

Assume $u(x,t)$ denotes the solution to one of the above mathematical models. We generate noisy data by using Equation (3.4) with $w_{i,j} = \sigma\,u(x_i, t_j)$ (i.e., $\beta = \sigma$ and $\gamma = 1$) in which all $\mathscr{E}_{i,j}$ terms are simulated as i.i.d. normal random variables with mean zero and variance one. We generate six data sets for each mathematical model, setting $\sigma = 0$, 0.01, 0.05, 0.10, 0.25, and 0.50. For Equation (3.6), we use its analytical solution to compute $u(x,t)$. For Equations (3.7) and (3.8) we use finite difference computations to numerically approximate $u(x,t)$. The numerical step sizes in these computations were chosen small enough to not introduce significant noise into the solution. We use $M = 101$ spatial locations and $N = 300$ time points to generate data for the diffusion-advection equation and $M = 199$ spatial locations and $N = 99$ time points for the Fisher-KPP and nonlinear Fisher-KPP Equations. As a preprocessing step, each data set is scaled to $[0,1]$ in order to consistently measure errors across the various data sets and noise levels.

### 3.3.2 Data denoising and library construction

In data denoising, we are interested in approximating the noisy data set, $\{U_{i,j}\}_{i=1,j=1}^{M,N}$, with some representation, $f(x,t|\theta)$. Computation of $f$ depends on parameters, $\theta$, such as knots

for polynomial splines or weights and biases for an ANN. This representation should match the dominant pattern of the data, so $\theta$ is chosen by finding the values that minimize the generalized least squares (GLS) cost function [1]

$$\mathcal{J}_f(\theta) = \frac{1}{MN} \sum_{i=1, j=1}^{M,N} \left( \frac{f(x_i, t_j|\theta) - u_{i,j}}{|f(x_i, t_j|\theta)|^\gamma} \right)^2. \tag{3.9}$$

Note that this cost function accounts for the statistical error model in (3.4) and also reduces to the ordinary least squares (OLS) estimator when $\gamma = 0$. In the case when the statistical error model is incorrectly specified, the resulting residuals can exhibit non-i.i.d. behavior [Ban14a], violating our assumptions on Equation (3.4). When the appropriate error model and cost function are not known *a priori*, residual computations and difference-based methods can provide insight into how to select these [Ban16]. Some methodology for choosing an appropriate error model is explained in Section B.1 in Appendix B. Thus, $\hat{\theta} = \arg\min_\theta \mathcal{J}_f(\theta)$ so that $f(x, t|\hat{\theta}) \approx u(x, t)$. From this representation, we can then take partial derivatives of $f$ with respect to $x$ and $t$ to estimate the partial derivatives of $u$.

### 3.3.2.1 Finite differences and spline approximations

**Finite differences.** We use central difference formulas on interior points and forward differences at the boundaries to obtain first order derivative approximations. For higher order derivatives (e.g. $u_{xx}$), first order finite difference rules are repeated on the corresponding previous order derivative approximations.

Finite difference approximations can be obtained accurately, efficiently, and directly from noiseless data, however, their accuracy quickly deteriorates in the presence of observation error. Following [Rud17a], we also employ polynomial spline regression. To thoroughly investigate the performance of spline computations for data denoising, we will consider three separate methods of spline computation in this work: local splines with a constant variance (CV) error model (i.e., $\gamma = 0$), local splines with a nonconstant variance error (NCV) model (i.e., $\gamma = 1$), and global splines with a NCV error model.

**Local splines.** For a given data point $u_{i,j}$ in the set of observations, we fit a cubic bi-spline on a small two-dimensional neighborhood of size $11 \times 11$ centered at $u_{i,j}$ by minimizing Equation (3.9). Denoised function and derivative approximations are then obtained using evaluations and the analytic derivatives of the fitted polynomial at the center point $u_{i,j}$. See Figure 3.2 for a visual diagram of the local spline method. Note that since we only

---

[1]Note that function values $f(x_i, t_j|\theta)$ less than 1e-4 in absolute value are set equal to one in the denominator $|f(x_i, t_j|\theta)|^\gamma$ during all training and evaluation for practical implementation.

approximate derivatives up to second order, higher order splines are not considered. For values of $u_{i,j}$ close to the boundary, we evaluate along the spline approximation that is nearest to the boundary. We found that cubic bi-spline approximations were generally more robust than the one-dimensional cubic splines used in [Rud17a], see Section B.2 in Appendix B, and Figures B.2-B.4.



**Figure 3.2** Diagram for using local splines for data denoising and partial derivative approximation. For a given data point $u_{i,j}$, shown on the left in red, in the set of observations, a cubic bi-spline, shown on the right as a red surface, is fit on a small two-dimensional neighborhood of size $11 \times 11$ centered at $u_{i,j}$.

**Global splines.** We perform global spline approximation of $\{U_i, j\}_{i=1,j=1}^{M,N}$ by minimizing Equation (3.9) for $\gamma = 1$ with

$$S(x,t) = \sum_{k=1,l=1}^{K,L} c_{k,l} S_{k,l}(x,t),$$ (3.10)

where $S_{k,l}(x,t)$ are normalized bivariate cubic bi-splines defined on the knot locations $(x,t)_{k,l} = \{\tilde{x}_k, ..., \tilde{x}_{k+4}\} \times \{\tilde{t}_l, ..., \tilde{t}_{l+4}\}$ [Die81]. To implement this approximation, we estimate a smoothness coefficient, the knot locations, and the spline coefficients by splitting the data into training and validation sets (50%/50%). The optimal values were chosen to be those that minimized the error (3.9) on the validation set. Details on implementation of this procedure are provided in Section B.3 in Appendix B. Derivative estimates are obtained analytically from the final form of $S(x,t)$ by evaluating the analytic derivatives of each

$S_{k,l}(x, t)$ term.

### 3.3.2.2 Artificial neural network approximations

An artificial neural network (ANN), denoted $h(\vec{x}|\theta)$, was used to approximate $u(x, t)$, with one hidden layer of the form

$$h(\vec{x}|\theta) = a_2\left( W_2\Big(a_1(W_1\vec{x} + b_1)\Big) + b_2 \right) \tag{3.11}$$

where $\vec{x} = [x \ t]^{\mathrm{T}}$ and $a_i(\cdot)$ represents the continuous nonlinear activation function for the $i^{\mathrm{th}}$ layer. The matrices $W_i$ and vectors $b_i$ (typically called weights and biases) comprise the total set of trainable parameters, $\theta = \{W_1, b_1, W_2, b_2\}$, of the network. We note that the use of one hidden layer in a neural network is a sufficient condition to make it a universal function approximator under the assumption that the activation function is bounded and non-constant [Hor91]. This result extends to ANNs with multiple hidden layers, however, we found that while training multi-layer ANNs resulted in faster convergence, the derivative approximations were worse. The task is therefore to find the optimal parameters $\theta^*$ such that $h(x, t|\theta^*) \approx u(x, t)$. The fitted surface function $h(x, t|\theta^*)$ and the computation of analytic derivatives of this function are used to approximate $u(x, t)$ in (3.4) and its partial derivatives for library construction in the PDE learning task. See Figure 3.3 for a diagram of the ANN method.



**Figure 3.3** Diagram for using ANNs for data denoising and partial derivative approximation [Nar20]. An ANN, shown on the left, inputs $x$ and $t$ pairs and outputs the corresponding approximations for $u(x, t)$. Using automatic differentiation, the analytic partial derivatives of the ANN can be used to construct a library of terms for the PDE learning task, shown on the right.

We found that the choice of activation function, $a_i(\cdot)$, in the ANN plays an important role in the accuracy of the partial derivative approximations. Typical activations like sigmoid and hyperbolic tangent yield oscillations in higher order derivative terms (e.g., see Supplementary Movie S1 available at `https://github.com/biomathlab/PDElearning/`). To mitigate this, we chose to use the "softplus" activation function which takes the form $\log(1 + e^z)$. This function has many desirable properties for approximating $u(x, t)$ (e.g. smoothness and infinitely many derivatives) which help ensure that the ANN can approximate the true function and its partial derivatives sufficiently well [Hor91]. However, the softplus function is unbounded, which violates an assumption of ANNs as universal approximators. While the assumptions on activation functions in [Hor91] are sufficient conditions and not necessary, one can address the unboundedness of the softplus function by including an $\ell_2$-regularization penalty on the activations $a_i$ in (3.11), but we found that no regularization was needed for the 18 data sets considered in this paper.

It becomes necessary to include an additional squared error term in the loss function to penalize function values outside $[0, 1]$. Without this term, the function values can blow up during training since $\mathscr{I}_h(\theta) \to 1$ as $h(\vec{x}|\theta) \to +\infty$ when $\gamma = 1$. Thus, the complete loss function used for training the ANN is

$$\mathscr{L}(\theta) = \frac{1}{MN} \sum_{i=1, j=1}^{M,N} \left( \frac{h(x_i, t_j|\theta) - u_{i,j}}{|h(x_i, t_j|\theta)|^\gamma} \right)^2 + \frac{1}{MN} \sum_{h \notin [0,1]} h^2 \tag{3.12}$$

where the first term corresponds to the generalized least squares cost function (3.9) and the second term corresponds to the additional squared error term to penalize function values outside $[0, 1]$.

We used 1,000 neurons in the hidden layer of the ANN. This choice was large enough to have maximal capacity to fit the data, while still allowing the optimization of $\theta$ to be computationally feasible on a desktop computer (3.4 GHz Intel Core i5 processor, 8gb RAM) without the need for GPU processing. The network parameters $\theta$ are optimized using the first-order gradient-based "Adam" optimizer [Kin14] with default parameters and a batch size of 10. We note that a small batch size paired with the adaptive moment estimation in Adam helps the ANN escape local minima during training which stabilizes convergence across various datasets. In order to prevent overfitting, the data were randomly split into training and validation sets (90%/10%) when training the ANN. The optimal network parameters were chosen to be those that minimized the error (3.9) on the validation set. We did not train the ANNs for some fixed number of epochs since (i) the parameters of each network are randomly initialized and (ii) the networks are trained on different data

sets which can lead to faster or slower convergence rates. Instead, early stopping of 50 (i.e. stopping training once validation error had not decreased for 50 consecutive epochs) was used to ensure convergence regardless of the data set or initial parameter values.

Representative examples of results from the bi-spline and ANN methods are shown in Supplementary Movies S2 and S3, respectively. All movies for all methods and noise levels considered in this work can be found at `https://github.com/biomathlab/PDElearning/animations/`.

### 3.3.3  Equation learning

We use the PDE-FIND algorithm [Rud17a] to discover the form of $F$ in Equation (3.1a) using computations of $u, u_x, u_{xx}$, and $u_t$ from the ANN, spline, and finite difference methods. Prior to implementing PDE-FIND, the numerical approximations are scaled from $[0, 1]$ back into their original scales. We discuss the PDE-FIND implementation in Section 3.3.3.1 and an additional pruning method in Section 3.3.3.2 that is used to remove extra terms from the final learned equation. We further discuss how we analyze our results in Section 3.3.3.3.

#### 3.3.3.1  PDE-FIND implementation

Once $u(x, t)$ and its partial derivatives have been computed, a large library of potential PDE terms is formed column-wise in the matrix, $\Theta$, given by

$$\Theta = \begin{bmatrix} 1 & u & \cdots & u^p & u_x & \cdots & u^p \odot u_x & u_{xx} & \cdots & u^p \odot u_{xx} & u_x^2 & u_x \odot u_{xx} & u_{xx}^2 \end{bmatrix} \quad (3.13)$$

where each column of $\Theta$ is some vectorization of the written term. All spline methods and the ANN have difficulty capturing the early dynamics of the diffusion-advection equation, so we skip the first 20 timepoints from the denoised data when building $\Theta$ for all data sets and denoising strategies. To reduce the computational time, only every fifth remaining timepoint is included in $\Theta$. Hence, while the data sets for the diffusion-advection equation begin with $N = 300$ timepoints, only (300-20)/5=56 timepoints are used in constructing $\Theta$. We set $p = 2$ resulting in $d = 12$ columns in $\Theta$. Each column of $\Theta$ thus represents a candidate term comprising $F$, so we assume

$$u_t \approx \Theta \xi, \quad (3.14)$$

where $\xi$ is a vector whose nonzero entries correspond to the true terms of $F$. The vector

$\xi$ is estimated using methods from sparse regression [Ris15]. Sequential threshold ridge regression was found to be a suitable method for estimating $\xi$ for PDE-FIND in a previous study [Rud17a]. However, we found that the Greedy algorithm performed well for the data and models we considered in this work. The Greedy algorithm computes

$$\hat{\xi} = \arg\min_{\xi \in \mathbb{R}^d} \frac{1}{MN} \| u_t - \Theta\xi \|_2^2, \quad \text{subject to } \|\xi\|_0 \le k \tag{3.15}$$

for some sparsity parameter, $k$ [Zha09].

The tolerance for which we solve Equation (3.15) for a given data set is treated as a hyperparameter that is found by splitting the library data into separate training and validation sets, and then optimizing over the validation set. In this training-validation procedure, we randomly divide our data points for $u_t$ into 5-by-5 tiles of adjacent spatiotemporal points and then randomly assign 50% of these tiles to a training data set, $u_t^{\text{train}}$, and the remaining 50% to a validation set, $u_t^{\text{validate}}$. We split the corresponding rows of $\Theta$ into $\Theta^{\text{train}}$ and $\Theta^{\text{validate}}$. We perform our hyperparameter search over 51 tolerance values, $k$, between 0 and $10^3$. For each value, we estimate $\hat{\xi}$ from the training set. For each estimate, we then compute its mean-squared error (MSE) over the validation set. We choose the hyperparameter corresponding to the $\hat{\xi}$ estimate with the smallest MSE on the validation data. The equation that results from sparse regression with this hyperparameter is our final equation from the PDE-FIND algorithm. We refer to the validation MSE from the final equation "val$_0$" in the remaining text.

### 3.3.3.2 Pruning method

We chose a 50-50 training and validation split for the data to avoid overfitting to the training data with a large validation set. Even so, we will demonstrate in Section 3.4.2 below that PDE-FIND is able to learn small but systematic biases from the ANN's fit to $u$ and its derivatives by incorporating extra terms into the final equations. Pruning methods have previously been developed that remove extra terms that do not significantly increase an algorithm's performance, see for example [And99; Mac95]. Accordingly, we implement the following pruning method after the PDE-FIND implementation described in Section 3.3.3.1 for all methods in order to delete the extra terms from the final equation.

The pruning procedure starts with a reduced library of candidate terms, $\tilde{\Theta}$, for the right hand side of Equation (3.1a) that correspond to the nonzero entries of $\hat{\xi}$ that resulted from our training-validation procedure. We then perform a sensitivity test for the remaining terms as follows. Suppose $\tilde{d}$ terms remain in $\tilde{\Theta}$, and let $\tilde{\Theta}_i, i = 1, \dots, \tilde{d}$ denote the further-

reduced library where the $i$th column of $\tilde{\Theta}$ has been removed. For each value of $i$, we find the least squares solution (without regularization) on the training data to the equation

$$u_t^{\text{train}} = \tilde{\Theta}_i^{\text{train}} \xi_i. \tag{3.16}$$

We then use our $\hat{\xi}_i$ estimate and compute the MSE over the validation data when the $i^{\text{th}}$ term has been removed and call this computation $\text{val}_i$. We then remove any candidate terms for our library that result in $\text{val}_i/\text{val}_0 < 1 + \alpha$ for some $\alpha > 0$. After this pruning step, we perform one final round of training without regularization over the fully reduced library to find the final form of our underlying equation.

It is important to note that choice of the $\alpha$ pruning threshold value warrants careful decision. If this value is chosen too high, then too few terms will be selected and the learned equation will be incomplete. If the chosen value is too small, then the final equation will admit extra terms arising from the systematic errors in derivative estimation. We will demonstrate below that the arbitrary choice of $\alpha = 0.25$ provides promising results for the diffusion-advection and Fisher-KPP Equations, while $\alpha = 0.05$ is suitable for the nonlinear Fisher-KPP Equation.

### 3.3.3.3   Accuracy metrics

To quantitatively assess the accuracy in recovering the correct PDE that generated the data, i.e., using the combined PDE-FIND with pruning methodology described above, we introduce the true positive ratio (TPR) for a given vector $\xi$ as:

$$\text{TPR}(\xi) = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}, \tag{3.17}$$

where TP ("True Positive") denotes the number of correctly-specified nonzero coefficients in $\xi$, FN ("False Negative") denotes the number of coefficients in $\xi$ that are incorrectly specified as zero, and FP ("False Positive") denotes the number of coefficients in $\xi$ that are incorrectly specified as nonzero. Recall that the nonzero entries of $\xi$ correspond to the relevant terms in an equation (i.e., for a library of $\Theta = [1\ u\ u_x\ u_{xx}]$, $\xi = [0\ 1\ 2\ 0]^T$ corresponds to the equation $u_t = u + 2u_x$). For example, when trying to learn Equation (3.6), an equation of the form $u_t = u_{xx} + u u_x$ would have TP $= 1$ (the nonzero coefficient for $u_{xx}$ is correct), FN $= 1$ (the missing $u_x$ term is incorrect), FP $= 1$ (the nonzero $u u_x$ term is incorrect), resulting in a final score of TPR $= 1/3$. Note that the TPR value is similar in nature to the Jaccard index: larger TPR values suggest that the true equation form has been better approximated,

and TPR $= 1$ signifies that the correct equation form has been recovered.

We note that the learned equation from the PDE-FIND method with pruning was often found to be sensitive to the random split of $u_t$ and $\Theta$ into training and validation data. Therefore, we performed PDE learning for 1,000 different random training-validation data splits of $u_t$ and $\Theta$ for each data set and for each computational method (finite differences, splines, and ANN). We then consider the distribution of $\text{TPR}(\hat{\xi})$ scores to assess the overall performance of the methodology. We declare the most commonly-learned equation among the 1,000 data splits as the final learned equation for each data set and computational method.

## 3.4    Results

In this section, we detail our results using the ANN to denoise data for $u(x, t)$ and compute partial derivatives. In addition, we test the accuracy of using the ANN method in conjunction with PDE-FIND to learn PDEs. Analogous results are presented for finite differences and splines. We begin by demonstrating the accuracy of the partial derivative calculations in Section 3.4.1, we explain why PDE-FIND finds small systematic bias terms in Section 3.4.2, and then we detail the accuracy in learning of the diffusion-advection, Fisher-KPP, and nonlinear Fisher-KPP equations in Sections 3.4.3-3.4.5.

### 3.4.1    Derivative calculations

We found that the finite difference method most accurately approximates $u$ and its derivatives for the advection-diffusion equation for $\sigma = 0$ (Table 3.1). This result is not surprising, as finite difference computations assume that there is no error in the data. For all other values of $\sigma$, we observe that the ANN produces the most accurate derivative calculations, although either the local or global NCV splines outperform the ANN at inferring $u(x, t)$ when the data are very noisy ($\sigma > 0.25$). It is important to note that the ANN's derivative calculations are often several  orders of magnitude more accurate than the spline and finite difference approximations, and this disparity between the computations appears to increase with $\sigma$ (Table 3.1). For example, at $\sigma = 0.01$, the ANN's relative mean squared error (RMSE) for $u_t$ is four orders of magnitude smaller than the RMSE for finite differences and at least two orders of magnitude smaller than the RMSE for all spline methods. At $\sigma = 0.50$, the ANN's RMSE for $u_t$ has become six orders of magnitude smaller than the RMSE for finite differences and at least four orders of magnitude smaller than the RMSE for all spline methods. The other derivative computations show similar results.

Similarly, we find that the ANN is most accurate for computing derivatives from noisy data from the Fisher-KPP Equation (Table 3.2) and the nonlinear Fisher-KPP Equation (Table 3.3). Recall that we do not have analytical solutions to these equations, so we used finite difference computations on the noiseless data ($\sigma = 0$) as an estimate for the analytical derivative values for the Fisher-KPP and nonlinear Fisher-KPP Equations. Again, in both cases we observe that the finite difference calculations perform best in computing the RMSE for $\sigma = 0$, but on average, the ANN provides the best calculations for the derivatives for larger values of $\sigma$. One of the local spline methods is consistently the most accurate at inferring $u$ from noisy data. The disparity between the RMSE derivative calculations for the ANN as compared to the splines or finite differences again appears to increase with $\sigma$ for these two equations.

### 3.4.2 PDE-FIND without pruning

We found that, in general, the PDE-FIND method learns the wrong equation, even when no noise is added to the data (Section B.4 in Appendix B, Figures B.5-B.7). Each denoising method resulted in accurate estimates for $u(x,t)$ and its partial derivatives in this case, however. For example, the residuals between the ANN model and the analytical values for $u, u_t, u_x$, and $u_{xx}$ were small when $\sigma = 0$ (Figure 3.4). We observed that, while small, the ANN residuals include systematic biases comprised of regions of over- and under-prediction. For example, all points near $(x,t) = (0.6, 0.4)$ for the ANN's calculation for $u_x$ appear to over-predict the true value for $u_x$ in this region (Figure 3.4). This contradicts the assumption of independence in $\{\epsilon_{i,j}\}_{i=1,j=1}^{M,N}$ for the statistical model in Equation (3.4). As we will now demonstrate, these small, systematic error terms from the ANN cause PDE-FIND to learn the incorrect equation.

We illustrate here that PDE-FIND learns the incorrect equation when training data for $u_t$ is comprised of $u_t$ at all spatial points for the first half of the given time points and the validation data is comprised of all spatial points for the second half of all time points. Recall that in our actual implementation discussed below, we randomly split the training and validation data in $5 \times 5$ bins of adjacent spatiotemporal points. Using denoised values for $u(x,t)$ and its partial derivatives from the ANN in the case where $\sigma = 0$ in the data, our training-validation procedure without pruning learns an equation of the form

$$u_t = a + b\,u_x + c\,u_{xx} + d\,u^2 + e\,u + f\,u^2 u_x + g\,u^2 u_{xx}, \ a,...,g \in \mathbb{R}. \qquad (3.18)$$

Similarly, the learned equations using finite difference and spline computations are

**Table 3.1** The relative mean-squared error (RMSE) between the noiseless data or true derivative values and our denoised data or derivative computations using finite differences, local splines with constant or nonconstant variance, global splines with nonconstant variance, and the ANN for the diffusion-advection equation. "FD" denotes finite differences, "LCVSP" denotes local splines with constant variance, "LNCVSP" denotes local splines with nonconstant variance, "GNCVSP" denotes global splines with nonconstant variance, and "ANN" denotes the ANN.

| $\sigma$ | Method | $U$ RMSE | $U_t$ RMSE | $U_x$ RMSE | $U_{xx}$ RMSE |
|---|---|---|---|---|---|
| 0.00 | FD | **0.00e+00** | **5.39e-05** | **5.77e-04** | **3.69e-02** |
| 0.00 | LCVSP | 1.22e-02 | 1.12e+00 | 3.33e-02 | 4.86e+01 |
| 0.00 | LNCVSP | 8.08e-04 | 1.41e+02 | 2.83e+00 | 5.56e+01 |
| 0.00 | GNCVSP | 7.47e-03 | 1.76e+02 | 1.51e+04 | 2.88e+07 |
| 0.00 | ANN | 2.86e-04 | 3.96e-01 | 8.47e-03 | 3.75e-01 |
| 0.01 | FD | **1.02e-04** | 2.08e+02 | 4.34e-01 | 3.52e+01 |
| 0.01 | LCVSP | 1.11e-02 | 7.34e+00 | 4.93e-02 | 4.87e+01 |
| 0.01 | LNCVSP | 7.83e-04 | 4.39e+02 | 3.29e+00 | 5.72e+01 |
| 0.01 | GNCVSP | 2.63e-03 | 1.05e+01 | 1.36e+02 | 1.19e+05 |
| 0.01 | ANN | 8.40e-04 | **7.71e-02** | **1.15e-02** | **6.93e-01** |
| 0.05 | FD | 2.51e-03 | 2.42e+03 | 9.97e+00 | 1.01e+03 |
| 0.05 | LCVSP | 1.19e-02 | 7.47e+01 | 2.62e-01 | 5.01e+01 |
| 0.05 | LNCVSP | 1.10e-03 | 1.95e+04 | 3.04e+01 | 2.49e+02 |
| 0.05 | GNCVSP | 1.64e-02 | 2.82e+06 | 4.34e+04 | 1.50e+09 |
| 0.05 | ANN | **5.61e-04** | **1.95e-01** | **7.71e-03** | **7.90e-01** |
| 0.10 | FD | 1.00e-02 | 4.04e+03 | 7.59e+01 | 3.78e+03 |
| 0.10 | LCVSP | 1.04e-02 | 2.32e+02 | 1.38e+00 | 5.45e+01 |
| 0.10 | LNCVSP | 1.97e-03 | 1.44e+04 | 2.13e+01 | 2.77e+02 |
| 0.10 | GNCVSP | 1.81e-01 | 1.76e+02 | 1.59e+03 | 7.96e+06 |
| 0.10 | ANN | **9.51e-04** | **1.23e-01** | **1.44e-02** | **7.69e-01** |
| 0.25 | FD | 6.28e-02 | 9.04e+04 | 2.20e+02 | 3.65e+04 |
| 0.25 | LCVSP | 2.51e-02 | 6.21e+03 | 5.83e+00 | 2.42e+02 |
| 0.25 | LNCVSP | **6.76e-03** | 1.64e+05 | 3.98e+02 | 3.42e+03 |
| 0.25 | GNCVSP | 3.73e-01 | 1.31e+03 | 1.42e+02 | 1.93e+06 |
| 0.25 | ANN | 7.29e-03 | **1.53e+00** | **4.49e-02** | **7.21e-01** |
| 0.50 | FD | 2.41e-01 | 2.58e+06 | 1.39e+03 | 1.05e+05 |
| 0.50 | LCVSP | 3.71e-02 | 3.78e+04 | 6.68e+01 | 7.95e+02 |
| 0.50 | LNCVSP | 3.49e-02 | 3.86e+05 | 1.32e+03 | 7.24e+03 |
| 0.50 | GNCVSP | **2.08e-02** | 3.51e+02 | 2.96e+04 | 2.58e+10 |
| 0.50 | ANN | 6.34e-02 | **3.43e+00** | **1.05e-01** | **1.44e+00** |

**Table 3.2** The relative mean-squared error (RMSE) between the noiseless data or true derivative values and our denoised data or derivative computations using finite differences, local splines with constant or nonconstant variance, global splines with nonconstant variance, and the ANN for the Fisher-KPP equation. "FD" denotes finite differences, "LCVSP" denotes local splines with constant variance, "LNCVSP" denotes local splines with nonconstant variance, "GNCVSP" denotes global splines with nonconstant variance, and "ANN" denotes the ANN.

| $\sigma$ | Method | $U$ RMSE | $U_t$ RMSE | $U_x$ RMSE | $U_{xx}$ RMSE |
|---|---|---|---|---|---|
| 0.00 | FD | **0.00e+00** | 3.16e-05 | **4.46e-04** | **4.67e-03** |
| 0.00 | LCVSP | 6.28e-05 | 2.83e-04 | 2.83e-03 | 2.01e-01 |
| 0.00 | LNCVSP | 3.80e-06 | 9.43e-02 | 5.54e-01 | 5.17e-01 |
| 0.00 | GNCVSP | 2.56e-02 | 3.36e+00 | 6.51e+01 | 2.96e+03 |
| 0.00 | ANN | 4.86e-04 | 6.98e-02 | 1.18e-01 | 2.66e+00 |
| 0.01 | FD | 9.82e-05 | 4.84e+00 | 6.47e+01 | 1.19e+03 |
| 0.01 | LCVSP | 6.91e-05 | 1.89e-01 | 1.56e+00 | **2.19e+00** |
| 0.01 | LNCVSP | **1.00e-05** | 1.49e+01 | 6.29e+00 | 9.63e+00 |
| 0.01 | GNCVSP | 1.01e-01 | 7.29e+00 | 2.91e+02 | 6.24e+03 |
| 0.01 | ANN | 3.90e-04 | **3.85e-02** | **1.19e-02** | 2.55e+00 |
| 0.05 | FD | 2.52e-03 | 9.49e+01 | 9.43e+02 | 6.78e+04 |
| 0.05 | LCVSP | 2.20e-04 | 8.85e+00 | 5.83e+01 | 2.39e+02 |
| 0.05 | LNCVSP | **1.71e-04** | 6.06e+02 | 2.20e+03 | 6.22e+02 |
| 0.05 | GNCVSP | 8.33e-03 | 1.90e+01 | 3.96e+03 | 3.39e+03 |
| 0.05 | ANN | 4.67e-04 | **1.03e-02** | **1.63e-02** | **1.57e+00** |
| 0.10 | FD | 9.95e-03 | 4.05e+02 | 2.81e+03 | 2.02e+05 |
| 0.10 | LCVSP | 6.80e-04 | 1.64e+01 | 6.66e+01 | 9.35e+02 |
| 0.10 | LNCVSP | **6.17e-04** | 1.22e+03 | 1.73e+03 | 2.65e+03 |
| 0.10 | GNCVSP | 3.90e-02 | 1.13e+01 | 3.20e+03 | 5.94e+03 |
| 0.10 | ANN | 8.46e-04 | **4.26e-02** | **5.29e-02** | **1.23e+00** |
| 0.25 | FD | 6.30e-02 | 2.92e+03 | 3.29e+04 | 3.92e+05 |
| 0.25 | LCVSP | 4.14e-03 | 2.26e+02 | 5.65e+02 | 1.44e+04 |
| 0.25 | LNCVSP | **4.07e-03** | 5.17e+03 | 9.49e+03 | 2.05e+04 |
| 0.25 | GNCVSP | 1.81e-02 | 4.53e+02 | 1.28e+05 | 3.27e+03 |
| 0.25 | ANN | 6.41e-03 | **6.94e-02** | **1.04e-01** | **5.90e+00** |
| 0.50 | FD | 2.38e-01 | 9.22e+03 | 1.01e+05 | 1.04e+06 |
| 0.50 | LCVSP | **1.52e-02** | 5.44e+02 | 1.32e+03 | 3.22e+04 |
| 0.50 | LNCVSP | 3.29e-02 | 2.36e+04 | 2.19e+04 | 7.33e+04 |
| 0.50 | GNCVSP | 1.87e-01 | 2.82e+01 | 1.14e+03 | 2.33e+04 |
| 0.50 | ANN | 6.60e-02 | **5.48e-01** | **8.98e-01** | **3.52e+01** |

**Table 3.3** The relative mean-squared error (RMSE) between the noiseless data or true derivative values and our denoised data or derivative computations using finite differences, local splines with constant or nonconstant variance, global splines with nonconstant variance, and the ANN for the Nonlinear Fisher-KPP equation. "FD" denotes finite differences, "LCVSP" denotes local splines with constant variance, "LNCVSP" denotes local splines with nonconstant variance, "GNCVSP" denotes global splines with nonconstant variance, and "ANN" denotes the ANN.

| $\sigma$ | Method | $U$ RMSE | $U_t$ RMSE | $U_x$ RMSE | $U_{xx}$ RMSE |
|---|---|---|---|---|---|
| 0.00 | FD | **9.71e-36** | **2.21e-05** | **1.52e-04** | **1.34e-01** |
| 0.00 | LCVSP | 1.53e-05 | 4.19e-05 | 1.09e-03 | 5.71e+00 |
| 0.00 | LNCVSP | 2.88e-06 | 3.14e-02 | 2.12e-02 | 5.59e+00 |
| 0.00 | GNCVSP | 1.61e-02 | 3.91e+01 | 1.87e+02 | 1.42e+02 |
| 0.00 | ANN | 8.73e-04 | 1.41e+01 | 6.16e+00 | 1.40e+02 |
| 0.01 | FD | 1.02e-04 | 2.21e+02 | 9.16e+03 | 8.26e+02 |
| 0.01 | LCVSP | 2.01e-05 | 1.02e+01 | 6.20e+02 | **3.32e+00** |
| 0.01 | LNCVSP | **7.98e-06** | 2.99e+01 | 6.26e+02 | 1.77e+01 |
| 0.01 | GNCVSP | 2.95e-03 | 5.76e+00 | 1.05e+02 | 1.30e+03 |
| 0.01 | ANN | 6.87e-04 | **5.55e+00** | **6.58e+01** | 1.90e+02 |
| 0.05 | FD | 2.43e-03 | 5.45e+03 | 2.65e+05 | 3.35e+04 |
| 0.05 | LCVSP | 1.57e-04 | 2.74e+02 | 6.84e+03 | **6.69e+01** |
| 0.05 | LNCVSP | **1.35e-04** | 8.13e+02 | 6.56e+03 | 5.10e+02 |
| 0.05 | GNCVSP | 2.73e-03 | **6.13e+00** | 1.55e+02 | 1.26e+03 |
| 0.05 | ANN | 1.08e-03 | 8.26e+00 | **1.68e+00** | 1.70e+02 |
| 0.10 | FD | 1.01e-02 | 2.30e+04 | 1.16e+06 | 5.69e+04 |
| 0.10 | LCVSP | 6.20e-04 | 1.32e+03 | 1.73e+04 | 7.23e+01 |
| 0.10 | LNCVSP | **5.37e-04** | 1.12e+04 | 1.77e+04 | 4.97e+03 |
| 0.10 | GNCVSP | 8.57e-04 | 3.94e+03 | 5.23e+03 | **6.72e+01** |
| 0.10 | ANN | 1.84e-03 | **1.93e+01** | **1.10e+01** | 2.04e+02 |
| 0.25 | FD | 6.25e-02 | 1.47e+05 | 4.95e+06 | 9.86e+05 |
| 0.25 | LCVSP | 4.00e-03 | 5.69e+03 | 1.24e+05 | 1.54e+03 |
| 0.25 | LNCVSP | **3.88e-03** | 2.89e+04 | 1.39e+05 | 4.75e+03 |
| 0.25 | GNCVSP | 5.79e-03 | 3.15e+03 | 3.09e+03 | 8.96e+02 |
| 0.25 | ANN | 6.34e-03 | **2.58e+01** | **2.34e+01** | **2.59e+02** |
| 0.50 | FD | 2.43e-01 | 5.73e+05 | 1.92e+07 | 2.46e+06 |
| 0.50 | LCVSP | **1.38e-02** | 2.87e+04 | 4.16e+05 | 1.13e+04 |
| 0.50 | LNCVSP | 1.49e-02 | 3.94e+04 | 4.45e+05 | 3.38e+04 |
| 0.50 | GNCVSP | 3.57e-02 | 1.42e+04 | 1.11e+04 | 7.11e+02 |
| 0.50 | ANN | 6.89e-02 | **6.49e+01** | **1.97e+02** | **4.88e+02** |

**Figure 3.4** Contours of the residual values between the ANN model ($u$, $u_x$,, etc.) and analytical solutions ($u_0$, $u_x 0$, etc.) for the diffusion-advection data with $\sigma = 0$. Top left: Residuals for $u$, top right: residuals for $u_x$, bottom left: residuals for $u_{xx}$, bottom right: residuals for $u_t$.

$$u_t = a\,u_x + b\,u_{xx} + c\,u^2 u_x,\ a, b, c \in \mathbb{R} \tag{3.19}$$

and

$$u_t = a\,u_x + b\,u_{xx} + c\,u^2 u_x + d\,u^2 u_{xx} + e\,u_x^2,\ a, ..., e \in \mathbb{R}, \tag{3.20}$$

respectively.

Each of these equations are incorrect and have extra terms on the right hand side of the learned PDE for the diffusion-advection equation. In Figure 3.5, we depict illustrative portions of the training and validation sets comparing the analytical values of $u_t$ against the computed values of $u_t$ and PDE-FIND's selected equation using ANN approximations. We found that PDE-FIND selects Equation (3.18) in place of the true diffusion-advection

43

equation because it recovers the ANN's incorrect computations of $u_t$ in both the training and validation data. In doing so, PDE-FIND fits the erroneous $u_t$ computations from the ANN approximation by including extra terms in the learned PDE.



**Figure 3.5** Results from the training (left) and validation (right) procedures in PDE-Find. The Red dashed line denotes the analytical value of $u_t$, the blue dots denote the computed $u_t$ values from the ANN, and the black lines denotes the equation for $u_t$ that has been computed with PDE-FIND.

### 3.4.3 Diffusion-advection equation

We tested whether implementing an additional pruning step with PDE-FIND could remove the extra terms resulting from the biases discussed in Section 3.4.2. For the diffusion-advection equation, we found that for all values of $\sigma$ except $\sigma = 0.01$, PDE-FIND with pruning achieves the highest median TPR when using ANN approximations (Figure 3.6). The ANN's median value is TPR $= 1$ (meaning that over half of the simulations yielded the correct equation form) for $\sigma = 0, 0.05, 0.10,$ and $0.25$. The ANN resulted in a median TPR $= 0.667$ at $\sigma = 0.50$. In contrast, the spline methods only achieve a median TPR $= 1$ at the lower noise levels $\sigma = 0, 0.01,$ and $0.05$ for the local methods, and never achieve a median TPR $= 1$ for the global method. For $\sigma \geq 0.10$, the medians for the all spline methods were TPR$\leq 0.5$. The finite difference method resulted in a median TPR $= 1$ at $\sigma = 0.01$, but the median TPR $= 0$ for larger values of $\sigma$.

Table B.1 in Appendix B shows the most commonly learned PDEs for each denoising method at each noise level. We found that the ANN method, used in conjunction with PDE-FIND with pruning, resulted in the correct PDE for $\sigma = 0, 0.05, 0.10, 0.25$. The ANN specifies

**Figure 3.6** TPR values for the diffusion-advection equation. We calculated the TPR, see Equation (3.17), for 1,000 different training-validation splits. These plots demonstrate the range of TPR values for each case. In each plot, the lower line in the colored box portion provides the 25% quartile of the data and the upper line denotes the 75% quartile. The "x" on each box plot denotes the median TPR value for that scenario. The length of the upper and lower whiskers are 1.5 times the interquartile range of the distribution, and diamonds denote outlier points. Any plot depicted as a solid horizontal line (e.g., the neural net computations for $\sigma = 0$) denotes that that this value is the majority of the range of the distribution.

the incorrect equation for $\sigma = 0.01$ and $\sigma = 0.50$. However in both of these cases, the extra terms have small parameter values (e.g. 0.001) that a scientist with an understanding of the system under consideration may manually neglect. On the other hand, PDE-FIND cannot discover the correct equation with finite difference or spline computations for $\sigma \geq 0.10$. These results suggest that the ANN method enables PDE-FIND with pruning to learn the diffusion-advection equation accurately at biologically realistic noise levels, e.g, $\sigma = 0.05, 0.10,$ and $0.25$.

### 3.4.4 Fisher-KPP equation

We tested the PDE-FIND with pruning method in conjuction with several denoising strategies using data from the Fisher-KPP Equation. We found that the ANN method had a median TPR = 1 (meaning that the correct equation is specified for at least half of the

training-validation data splits) for $\sigma = 0, 0.01, 0.05,$ and $0.10$ (Figure 3.7). In contrast, the finite difference calculations only had a median TPR $= 1$ at $\sigma = 0$, and the local spline methods only have median TPR $= 1$ at $\sigma = 0, 0.01$ while the global spline method never achieves a median TPR $= 1$. The accuracy in using PDE-FIND with the spline and finite difference methods quickly deteriorates for high noise levels. The finite difference method resulted in a median TPR $= 0$ for $\sigma \geq 0.05$, and the all spline methods result in a median TPR less than $0.667$ at $\sigma = 0.05, 0.10$. At $\sigma = 0.25, 0.50$, the local spline methods have median TPR $= 0$ while the global spline method maintains a median TPR $= .667$. The ANN had a median TPR $= 0.6$ and $0.5$ for $\sigma = 0.25$ and $0.5$, respectively.



**Figure 3.7** TPR values for the Fisher-KPP Equation. We calculated the TPR, see Equation (3.17), for 1,000 different training-validation splits. These plots demonstrate the range of TPR values for each case. In each plot, the lower line in the colored box portion provides the 25% quartile of the data and the upper line denotes the 75% quartile. The "x" on each box plot denotes the median TPR value for that scenario. The length of the upper and lower whiskers are 1.5 times the interquartile range, and diamonds denote outlier points. Any plot depicted as a solid horizontal line (e.g., the finite difference computations for $\sigma = 0$) denote that this value is the majority of the distribution.

Table B.2 in Appendix B shows the most commonly chosen PDEs resulting from the PDE-FIND with pruning method. We found that PDE-FIND with pruning is able to discover the

correct equation form for $\sigma = 0, 0.01, 0.05,$ and 0.10 when using the ANN approximations. While PDE-FIND is unable to specify the correct equations with ANN data for $\sigma = 0.25$ and 0.50, all of the terms in the Fisher-KPP equation were included in the learned PDEs. In contrast, using the local spline methods for denoising resulted in only learning the correct PDE for $\sigma \leq 0.01$. For $\sigma$, the spline methods resulted in large errors in the derivative approximations and did not recover the $u_{xx}$ terms for $\sigma = 0.05, 0.10$ and did not yield any terms on the right hand side of the learned PDE for $\sigma = 0.25, 0.50$. The global spline method never recovers the $u_{xx}$ term in its final recovered equation, but does recover the $u$ and $u^2$ terms. Similarly, the finite difference method resulted in only learning the true equation form for $\sigma = 0$. These results suggest that only the ANN method enables PDE-FIND with pruning to learn the Fisher-KPP equation for reasonably high noise levels of $\sigma = 0.05, 0.10$, whereas the spline methods cannot for $\sigma > 0.05$.

### 3.4.5   Nonlinear Fisher-KPP equation

We found that the PDE-FIND with pruning method was not able to recover the correct PDE from data that has been generated by the nonlinear Fisher-KPP Equation for all denoising strategies considered. PDE-FIND could not achieve a median TPR = 1 for any of these methods, meaning that the correct equation was never specified for over half of the training-validation data splits (Figure 3.8). All methods have median TPR $\leq 0.8$ at $\sigma = 0$. When using ANN approximations, the PDE-FIND with pruning method has median TPR = 0.8 at $\sigma = 0.01$ and 0.05, TPR = 0.6 at $\sigma = 0.10$ and 0.25, and TPR = 0.5 at $\sigma = 0.50$. When using local spline computations, PDE-FIND with pruning has median TPR = 0.8 at $\sigma = 0.01$, TPR = 0.5 at $\sigma = 0.05, 0.10,$ and 0.25, and TPR = 0 at $\sigma = 0.50$. When using the global spline computation, PDE-FIND with pruning has median TPR = 0.4 at $\sigma = 0, 0.01, 0.05, 0.25,$ and 0.50 and TPR = 0.5 at $\sigma = 0.10$. When using finite difference computations, PDE-FIND with pruning has median = 0.50 at $\sigma = 0.01$ and 0.05 and TPR = 0 for $\sigma \geq 0.10$.

While all of the denoising strategies lead to incorrect equations, the ANN strategy recovers the most relevant terms in its final equations. Table B.3 in Appendix B shows the most-commonly chosen PDEs resulting from the PDE-FIND with pruning method. We found that all methods except the global splines predict the true equation with an extra Fickian diffusion term, $u_{xx}$, for $\sigma = 0$. The global splines do not recover the $uu_x$ and $u_x^2$ terms. When using the ANN approximations, the PDE-FIND with pruning algorithm recovers the true equation with added Fickian diffusion at $\sigma = 0.05$, and it recovers three of the correct terms but excludes $uu_{xx}$ at $\sigma = 0.01$. For larger values of $\sigma$ with the ANN approximations, PDE-FIND with pruning recovers three correct terms, excludes the $uu_{xx}$

**Figure 3.8** TPR values for the nonlinear Fisher-KPP Equation. We calculated the TPR, see Equation (3.17), for 1,000 different training-validation splits. These plots demonstrate the range of TPR values for each case. In each plot, the lower line in the colored box portion provides the 25% quartile of the data and the upper line denotes the 75% quartile. The "x" on each box plot denotes the median TPR value for that scenario. The length of the upper and lower whiskers are 1.5 times the interquartile range, and diamonds denote outlier points. Any plot depicted as a lone solid horizontal line (e.g., the finite difference computations for $\sigma = 0$) denotes that that this value is the whole range of the data.

term, and includes an extra Fickian diffusion term (as well as an additional constant term at $\sigma = 0.50$). When using local spline computations, the PDE-FIND with pruning algorithm recovers three correct terms, excludes the $u u_{xx}$ term, and adds an extra Fickian diffusion term at $\sigma = 0.01$. For $\sigma = 0.05 - 0.25$, the final equation recovers two correct terms but excludes the $u u_{xx}$ and $u_x^2$ terms. At $\sigma = 0.50$, all terms are deleted when using local spline computations. When using global spline computations, PDE-FIND with pruning recovers only the reaction terms when $\sigma = 0.10$. For all other noise levels, it replaces the nonlinear diffusion terms with Fickian diffusion. When using finite difference computations, the PDE-FIND with pruning algorithm correctly recovers two terms but excludes the $u u_{xx}$ and $u_x^2$ terms at $\sigma = 0.01$ and 0.05. For larger values of $\sigma$, no correct terms are included in the final equation form.

We investigated if the recovered terms from the PDE-FIND with pruning algorithm using ANN approximations can be used as the specified mathematical model in an inverse

problem methodology (cf., [Ban09b]) to recover the final parameter estimate values from the nonlinear Fisher-KPP Equation. If we take the union of all terms that are included in the final equations in Table B.3 in Appendix B for the ANN method using noisy data ($\sigma > 0$), then we have an equation of the form

$$u_t = a\,u_{xx} + b\,u\,u_{xx} + c\,u_x^2 + d\,u + e\,u^2 + f; a,\ldots,f \in \mathbb{R}. \tag{3.21}$$

We estimated the parameters $a,\ldots,f$ in Equation (3.21) for each value of $\sigma$ by simulating the solution to this PDE using the method of lines and minimizing Equation (3.9) using the Nelder-Mead algorithm. We input the equations from Table B.3 in Appendix B for the ANN method as the initial guess for each data set. We find that performing this inverse problem leads to accurate parameter estimates for the true terms in the nonlinear Fisher-KPP Equation and small coefficient values for the incorrect terms ($u_{xx}$ and 1) for $\sigma = 0, 0.01, 0.05,$ and 0.25 (Table 3.4). At $\sigma = 0.10$ and 0.50, this inverse problem methodology leads to small coefficient estimates for $u\,u_{xx}$ in addition to $u_{xx}$ and 1. Note that this same process would not lead to ultimately recovering the true equation and parameter estimates from the spline or finite difference approximations because their final equations never included the correct $u\,u_{xx}$ term in the final equation for noisy data ($\sigma > 0$).

**Table 3.4** Inferred parameters for the nonlinear Fisher-KPP Equation data when performing an inverse problem on Equation (3.21).

| | True Equation |
|---|---|
| | $u_t = .02\,u\,u_{xx} + .02\,u_x^2 + 10\,u - 10\,u^2$ |
| $\sigma$ | **Revised Equation** |
| 0.00 | $u_t = 2.3 \times 10^{-12}\,u_{xx} + .017\,u\,u_{xx} + .021\,u_x^2 + 10.0\,u - 10.0\,u^2 - 1.60 \times 10^{-8}$ |
| 0.01 | $u_t = 1.4 \times 10^{-5}\,u_{xx} + .021\,u\,u_{xx} + .019\,u_x^2 + 10.0\,u - 10.0\,u^2 - 1.38 \times 10^{-8}$ |
| 0.05 | $u_t = 2.5 \times 10^{-4}\,u_{xx} + .0034\,u\,u_{xx} + .029\,u_x^2 + 9.9\,u - 10.0\,u^2 - 2.86 \times 10^{-8}$ |
| 0.10 | $u_t = 6.1 \times 10^{-4}\,u_{xx} + 8.42 \times 10^{-4}\,u\,u_{xx} + .023\,u_x^2 + 9.38\,u - 9.52\,u^2 - 2.44 \times 10^{-4}$ |
| 0.25 | $u_t = 7.2 \times 10^{-4}\,u_{xx} + .015\,u\,u_{xx} + .024\,u_x^2 + 9.8\,u - 9.14\,u^2 + 3.53 \times 10^{-7}$ |
| 0.50 | $u_t = 2.1 \times 10^{-3}\,u_{xx} - 3.72 \times 10^{-3}\,u\,u_{xx} + .032\,u_x^2 + 9.7\,u - 7.5\,u^2 + 1.38 \times 10^{-6}$ |

## 3.5   Summary of contributions

We introduced a novel ANN-based method for robust data denoising and partial derivative approximation in the presence of biologically realistic forms and levels of noise. Based on the numerical experiments in this work, the neural network derivative approximations were up to four orders of magnitude more accurate compared to polynomial splines, a state-of-the-art denoising method. Further, using automatic differentiation to construct a library of candidate terms thought to comprise the unknown governing system, the superior denoising capabilities of the ANN also translated to more accurate discovery of the unknown PDE with noise levels up to 25%. Finally, this work also showcased significant extensions to the PDE-FIND sparse regression algorithm involving training/validation partitioning and pruning to identify sensitive PDE terms that constitute a strong signal in the data.

## 3.6   Discussion

The novel use of the ANN method presented here is a significant step toward making PDE learning more achievable in realistic scenarios with noisy biological data. Because an ANN is a fully differentiable function, it can be used to approximate derivative computations to build the library of terms needed for PDE learning. The current practice to build a library of terms for learning PDEs when noise is present in the observed data $u(x, t)$ is to use finite difference or local spline approximations for small amounts of noise [Rud17a; Zha18]. Finite difference and spline-based techniques have been reported as state-of-the-art denoising methods for equation learning in a recent study by Rudy et al. [Rud17a]. We note that Rudy et al. also considered several other numerical methods such as Gaussian kernel smoothing, Tikhonov differentiation, spectral differentiation with high-frequency term thresholding and found that the most reliable and robust method was polynomial interpolation. We expanded on their use of local uni-variate polynomial splines by implementing local bi-variate splines which we found to be more robust for function and derivative approximation, and also implemented global spline approximations as a closer comparison to the performance of the ANN as a global approximation method. Our findings suggest that finite difference and spline methods are highly sensitive to the amount of noise in the data in the range of less than 5% noise. We showed that the ANN method outperforms spline and finite difference approximations of the partial derivatives of $u(x, t)$, even when some spline methods better approximated $u(x, t)$ than the ANN, in the presence of significant levels of non-constant error noise . It is important to note that this level of noise and nonconstant

variance are typical phenomena encountered in biological data [Ban11].

We compared polynomial spline approximations and ANNs as global methods to denoise data because there exist previous theoretical results regarding the ability of polynomials and ANNs to approximate continuous functions up to arbitrary precision on compact domains [Hor91; Sto48]. This theory is only relevant in the case of noiseless data and in the limit of a large number of neurons or right hand side terms. In practice, neither of these assumptions is realistic when dealing with biological data, which tends to have a large noise-to-signal ratio, and where one is required to optimize hyperparameters (e.g., number of neurons, smoothness, knots) for a given data set. To the best of our knowledge, there has not been a direct comparison between the practical performance of polynomial splines and ANNs in approximating a continuous function and its partial derivatives from noisy non-constant variance data. We also note that derivative estimation is often ignored in the optimization for various methods, e.g., splines and ANNS, used to approximate a function $f(x, t|\theta)$ for $u(x, t)$. For example, Equation (3.9) does not consider the ability of $f_x(x, t|\theta)$ to approximate $u_x(x, t)$ or $f_t(x, t|\theta)$ to approximate $u_t(x, t)$. The numerical results presented here are the first such comparison between various local and global approximation methods in their ability to estimate both $u(x, t)$ and its derivatives. We found that the ANN is comparable to spline methods in approximating $u(x, t)$ and more accurate in approximating its spatial and temporal derivatives. The disparity between the accuracy of ANNs and other methods in estimating derivatives sharply increased with higher noise in the data. These results suggest the need for further theoretical investigation into the fidelity of derivative estimation from the function approximation methods considered here in the presence of realistic forms and levels of noise. Since a number of heuristics were used in the development of our proposed method, this suggests that future modifications to the network architecture may increase denoising and derivative approximation accuracy. We postulate that one reason for the superior performance of our proposed ANN for derivative approximation is that the chosen activation function (i.e. softplus) is sufficiently smooth. This satisfies conditions in [Hor91] that ensure the capability of ANNs to converge sufficiently close to the true partial derivatives in addition to the underlying dynamical system. Our results suggest that ANNs outperform other methods considered here in this respect for the biological transport models we considered (Tables 3.1 - 3.3).

The importance of our denoising methodology results is underscored by the need for accurate derivative estimates in equation learning techniques. We note that while methods such as Bayesian inference and Kalman Filtering have been commonly used in the mathematical modeling literature [Smi13] for denoising data, we did not compare these approaches in this work as they require specifying a mathematical model for the data *a priori*.

In contrast, the two step procedure we performed is used to learn a mathematical model from the data, which requires a model-free method to first denoise the data and numerically approximate derivatives. Indeed, we found that the PDE-FIND algorithm can successfully recover the true equations underlying $u(x, t)$ when $u(x, t)$ and its derivatives have been accurately recovered. For example, the ANN outperforms all splines methods in computing the derivatives of $u$ for the advection-diffusion equation when $\sigma > 0$ (See Table 3.1), and in turn the ANN computations allow PDE-FIND to learn the correct underlying equation form for $\sigma = 0.05, 0.10$, and $0.25$, whereas all finite difference and spline computations fail in identifying the underlying equation for $\sigma > 0.05$ (See Figure 3.6). Using the ANN estimates, we found that PDE-FIND could be used to learn the Fisher-KPP Equation for up to 10% noise levels ($\sigma = 0.1$), whereas the other methods fail at 5% ($\sigma = 0.05$). None of the methods we considered were able to learn the nonlinear Fisher-KPP Equation, but the ANN computations lead to an equation that can then be used with an inverse problem methodology [Ban09b] to infer which terms are meaningful. We found in this study that identifying the correct underlying equation under the high levels of observation noise we considered was a challenging problem for current state-of-the-art equation learning methods, even under the restriction of setting $p = 2$ in Equation (3.2). Future work will investigate correctly inferring models equations from much larger libraries for the current equation learning method.

We focused on three common transport models in this work. These models are broadly relevant across several fields in biology and have been used previously to describe the movement and growth of a varying array of spatiotemporal processes, e.g., wound healing [Jin16b; Joh14; Mai04; She90a], cancer progression [Bal14; Rut17; Ste15], and animal development and herd migration [Fra03; Has05; Sib99]. In future work, we aim to extend the results of this study to more complex PDE models in biology, such as systems of PDEs used to model several subpopulations, processes in higher spatial dimensions, and more complex non-linearities used to describe migration.

There are several reasonable explanations for the difficulties in learning the nonlinear Fisher-KPP Equation. Three of the terms in Equation (3.8) are the product of two terms including $u$ or its derivatives ($u^2, u_x^2$, and $u u_{xx}$). In practice, these terms may be inaccurate approximations from noisy data (as demonstrated in Table 3.3). Multiplying two inaccurate terms may lead to an even larger amount of uncertainty associated with these estimates. We postulate that the high level of uncertainty in these types of terms resulting from the product of inaccurate estimates likely increases the difficulty of learning to include them in the process of PDE learning. Furthermore, it must be noted that the data for this equation was generated numerically by the finite difference method. Though some analytical solutions

to the Fisher-KPP equations are known, they come either in series form, which is beyond the scope of this article as it would require learning an infinite number of polynomial terms ($p = \infty$), or in traveling wave form, which would be indistinguishable for PDE-FIND from the advection equation [Mur02]. The finite difference methods used to approximate the spatial derivatives in the nonlinear Fisher-KPP equation introduce second and fourth-order error terms which lead to numerical dispersion and diffusion effects that may account for the recovery of some unexpected terms.

We found that the use of pruning following our implementation of the PDE-FIND algorithm increased our ability to recover the correct equation in terms of the TPR. It may be argued that these additional terms learned from the PDE-FIND implementation without pruning would be removed from the final equation if more regularization (i.e., a larger value of $k$ in implementation of the Greedy algorithm) were used. However, when faced with the issue of learning the governing equation from actual data in practice, one will not have the ability to know when the specified equation is correct or not. We thus need to identify the correct hyperparameters without any *a priori* knowledge. We observed that the systematic biases in our ANN (depicted in Figure 3.4) make it difficult to choose a hyperparameter value that leads to the correct equation because these biases are present in the training and validation data. The pruning algorithm is a way to correct for when the incorrect hyperparameter has been chosen by ensuring that all terms in the learned PDE are sufficiently sensitive to constitute a strong signal in the data, instead of resulting from a bias in our approximation methods.

Our use of pruning to remove terms from learned PDEs could be improved in future work. While effective, our implementation is somewhat crude, in which we pre-specify a threshold level to prune parameters based on out-of-sample MSE values on the validation data set. Previous studies have discussed F-statistics as one way to infer the increase in variance that pruning a variable will lead to, but there are still many different interpretations of these results which makes a definitive statistical pruning method challenging to ascertain [Bur95].

CHAPTER

4

# BIOLOGICALLY-INFORMED NEURAL NETWORKS

## 4.1 Motivation

In the previous work, we outlined a general strategy for utilizing artificial neural networks (ANNs) for the denoising and partial derivative approximation of data in the presence of biologically realistic forms and levels of noise. Automatic differentiation is used on the network outputs to construct a library of candidate terms thought to comprise the unknown governing PDE that describes the time evolution of the dynamical system. By (i) using the library of candidate terms from the ANN method and (ii) extending PDE-FIND, a sparse regression algorithm, with training/validation data partitioning and a pruning algorithm that automatically removes unnecessary candidate terms, we developed a methodology that significantly exceeds state-of-the-art methods in their ability to both denoise and learn equations from data. However, our method (and the sparse regression approach in general) assumes that the governing system of equations can be written as a linear combination of nonlinear candidate terms. In the event that the true system contains terms that can not be represented this way, e.g., if a true term contains exponents that also must be estimated, then this method will ultimately fail to capture the true dynamics. Further, ANNs

are severely prone to overfitting, and can require complex regularization strategies based on domain expertise (if applicable) or heuristics. Other approaches exist as alternatives, e.g., physics-informed neural networks [Rai19], however, these methods require the *a priori* specification of the governing system of ODEs/PDEs, which significantly limits their use for equation learning. Below, we take insights from this work together with physics-informed neural networks to produce a robust end-to-end neural-network-based equation learning methodology, nicknamed biologically-informed neural networks, that can be applied on a wide spectrum of physical and biological problems (for both ODE and PDE systems) in which the available data are noisy and sparse and the governing dynamics are unknown and highly nonlinear.

## 4.2   Introduction

*Collective migration* refers to the coordinated migration of a group of individuals [Fri09; Vic95]. This process arises in a variety of biological and social contexts, including pedestrian dynamics [Hel95], tumor progression [Gal13], and animal development [McL15]. In the presence of many individuals, differential equation models provide a flexible framework to investigate collective behavior as a continuum [Arc11; Dys15; Joh12; Nar16; Top12]. A challenge for mathematicians and scientists is to use mathematical models together with spatiotemporal data of collective migration to validate assumptions about the underlying physical and biological laws that govern the observed dynamics. Several factors contribute to the difficulty of this task, even for simple systems/data, some of which include biological forms and levels of noise in the observation process, poor understanding of the underlying dynamics, a large number of candidate mathematical models, implementation of computationally expensive numerical solvers, etc. This work provides a data-driven tool which can alleviate many of these problems by enabling the rapid development and validation of mathematical models from sparse noisy data. The methodology is demonstrated using a case study of scratch assay experiments.

Scratch assays are a widely adopted experiment in cellular biology used to study collective cell migration *in vitro* as cell populations re-colonize empty spatial regions. These experiments have been used previously to observe population-wide behavior in many different contexts, including wound healing [Aok17; Cha14; Mat04; Nik06] and cancer progression [Har17]. Mathematical modeling of scratch assays plays an important role in the quantification and analysis of population dynamics. This is because (i) the equations and parameters comprising mathematical models are *interpretable*, providing information about the underlying physical and biological mechanics that drive the observed system,

and (ii) when properly calibrated, they are *generalizable*, affording the ability to make accurate predictions beyond the data set used for calibration.

Reaction-diffusion partial differential equations (PDEs) are frequently used to model scratch assay experiments [Jin16a; Arc11; Nar16; Joh15; Mai04]. The general one-dimensional reaction-diffusion equation that describes the rate of change of a quantity of interest $u(x, t)$ (e.g. cell density) is

$$u_t = (\mathscr{D} u_x)_x + \mathscr{G} u, \quad x \in [x_0, x_f], \quad t \in [t_0, t_f], \tag{4.1}$$

in which the rate of change of $u$ (i.e. $u_t$) is a function of diffusion, modeled by the function $\mathscr{D}$, and reaction or growth, modeled by the function $\mathscr{G}$. Note that $\mathscr{D}$ and $\mathscr{G}$ depend on the application, and choosing the correct/optimal mechanistic models for these terms is the focus of many current research efforts and remains an open question. The classical Fisher–Kolmogorov–Petrovsky–Piskunov (FKPP) equation is a reaction-diffusion equation that has been used to model a wide spectrum of growth and transport of biological processes. In particular, the FKPP model assumes a scalar diffusivity function $\mathscr{D} = D$ and logistic growth function $\mathscr{G} = r(1 - u/K)$ with intrinsic growth rate $r$ and carrying capacity $K$ [Bal14; Mai04]. Variants of the reaction-diffusion equation have also been used to account for different types of cell interactions during scratch assay experiments. For example, the nonlinear diffusivity function $\mathscr{D} = 1 - \alpha^4/3 + 3\alpha(u/K - 2/3)^2$ with cell-to-cell adhesion coefficient $\alpha$ was used to model dynamics in which neighboring cells prevent other cells from migrating [Ang09]. Alternatively, a diffusivity function of the form $\mathscr{D} = D(1 + \alpha(u/K)^2)$ can be used to model dynamics in which cells promote the migration of others [Nar16]. Additional variants of reaction-diffusion equation models have captured cell migration in the presence of growth factors [Dal95], during melanoma progression [Har17], and in response to different drug treatments [Joh15].

A recent study quantitatively investigated the role of initial cell density by conducting a suite of scratch assay experiments on PC-3 prostate cancer cells with systematically varying initial cell densities [Jin16a]. The experimental data was used to calibrate the FKPP equation as well as a variant model known as the Generalized Porous-FKPP equation, which assumes that diffusivity increases with cell density $u$ by using a diffusivity function $\mathscr{D} = D(u/K)^m$ with diffusion coefficient $D$, carrying capacity $K$, and exponent $m$. Like the FKPP equation, the growth term is also described by the logistic growth function $\mathscr{G} = r(1 - u/K)$. While the calibrated models approximated the experimental data well in many cases in [Jin16a], the presence of systematic biases between the model solutions and experimental data indicate the existence of additional governing mechanisms that may not be accounted

for in these mathematical models. However, the existence of a large number of possible biophysical mechanisms that could play a role in scratch assay dynamics makes the testing of mathematical models against these experimental data computationally challenging. Thereby, this scenario motivates the use of equation learning methods to discover the diffusion and reaction terms directly from the experimental data.

Enabled by advances in computing power, algorithms, and the amount of available data, the field of equation learning has recently emerged as a powerful tool for the automated identification of underlying physical laws governing a set of observation data. The basic assumption in this field is that measured data arise from some unknown $n$-dimensional dynamical system of the form

$$u_t = \mathscr{F}(x, t, u, u_x, u_{xx}, \ldots; \theta), \quad x \in [x_0, x_f], \quad t \in [t_0, t_f], \tag{4.2}$$

with quantity of interest $u = u(x, t)$, parameter vector $\theta \in \mathbb{R}^k$, and appropriate initial and boundary conditions. An example quantity of interest for modeling cell migration dynamics is the cell density (cells/mm²) at location $x$ and time $t$. The measured data $\{u_{i,j}\}_{i,j=1}^{M,N}$ for a set of spatial points $x_i$, $i = 1, \ldots, M$, and set of time points $t_j$, $j = 1, \ldots, N$, are assumed to be corrupted by some form of observation error that may be known or unknown in practice. The goal of equation learning methods is to identify the closed form of $\mathscr{F}$ in Equation (4.2) directly from the noisy measurements $u_{i,j}$. Note that, in order to simulate the learned equation, either the noisy or a denoised version of the initial condition can be used along with an assumed boundary condition (e.g. no-flux) that describes the biological process generating the data.

Two primary sets of methodology have been used in field of equation learning to date: sparse regression [Bru16b; Rud17b] and theory-informed neural networks [Rai19; Yan20]. In the sparse regression framework, numerical methods (e.g. finite differences or polynomial splines) are used to denoise $u$ and approximate the partial derivatives $u_t$, $u_x$, $u_{xx}$, etc. from a set of data. The approximations are then used to construct a library of nonlinear candidate terms (e.g. 1, $u$, $u^2$, $u_x$, $\ldots$, $u_x^2 u_{xx}^2$, etc.) thought to comprise the governing system of ordinary differential equations (ODEs) or PDEs. The data relating $u_t$ to all possible model terms inside the library are formulated as a linear regression problem in which sparsity promoting techniques are used to select a small subset of library terms that produce the most parsimonious model. While the sparse regression framework has been successfully demonstrated to circumvent searching through a combinatorially large space of possible candidate models, it can require large amounts of training data and the numerical methods used for denoising and differentiation are not robust to biologically realistic forms and

levels of noise, leading to inaccuracies in both the constructed library and learned equations [Lag20d]. Further, the method assumes the unknown function $\mathscr{F}$ in Equation (4.2) can be written as a linear combination of nonlinear candidate terms, which may not be true in practice.

An alternative approach uses function-approximating deep neural networks, i.e., multilayer perceptrons (MLPs), as surrogate models, $u_{\mathrm{MLP}}(x, t)$, for the solution of the governing dynamical system [Rai19; Yan20]. In this approach, the assumed mechanistic form of $\mathscr{F}$ in Equation (4.2) is pre-specified and then used as a form of regularization in the neural network objective function. The parameters of $\mathscr{F}$ are allowed to be "learnable," meaning that the parameters of the governing PDE are calibrated while the neural network is trained to minimize the error between $u_{\mathrm{MLP}}(x_i, t_j)$ and the data, $u_{i,j}$. This methodology ensures that the neural network solution satisfies the physical laws described by $\mathscr{F}$ while simultaneously fitting the spatiotemporal data. Theory-informed neural networks have been demonstrated with smaller amounts of data in the presence of noise, however, they have so far only been applied to problems where the governing mechanistic PDE is known *a priori*.

Hybrid approaches that combine neural networks and sparse regression have also been suggested to address some of the issues surrounding the above methods [Lag20d; Bot19]. In these approaches, neural networks are used as surrogate models for $u(x, t)$ and then used to construct the library of candidate terms for sparse regression using automatic differentiation. These methods have been shown to accurately learn the governing system of equations for a variety of reaction-diffusion models from spatiotemporal data with biologically realistic levels of noise [Lag20d].

All three approaches (i.e. sparse regression, theory-informed neural networks, and hybrids) however, suffer from the *model specification problem*, in which the governing ODE/PDE model must be specified *a priori* either explicitly or as a library of candidate terms. Thus, (i) if the true dynamical system contains terms that are not included in the regularization term for theory-informed neural networks, or (ii) if the true terms cannot be represented as a linear combination of nonlinear candidate terms for sparse regression, then these methods will ultimately fail to recover the true system. Further, detecting this issue when determining what the "true" system is in real-use cases is an open question. Where systems with scalar or linear dynamics may be suitable for these approaches, biological systems pose a particular challenge in this respect, since many of the underlying mechanics driving these systems are nonlinear. For example, the Generalized Porous-FKPP model contains a nonlinear diffusivity function $\mathscr{D} = D(u/K)^m)$ with unknown exponent $m$. These issues help explain why, to the best of our knowledge, equation learning methods have not yet been successfully applied to real-world biological population-level data.

In this work, biologically-informed neural networks (BINNs), an extension of physics-informed neural networks (PINNs) [Rai19], are presented as a solution to the library specification problem for systems with biological/physical constraints. In this framework, the right-hand-side function $\mathscr{F}$ of the PDE in Equation (4.2) is assumed to be a combination of biologically relevant terms. For example, the general form of reaction-diffusion models can be described by the two right-hand-side terms in Equation (4.1) meaning that the equation learning problem is transformed from learning $\mathscr{F}$ to learning the diffusivity and growth functions $\mathscr{D}$ and $\mathscr{G}$. Rather than assigning mechanistic forms to each function as in previous equation learning studies, each function is replaced with a separate neural network. This approach leverages the ability of deep neural networks to approximate continuous functions arbitrarily well [Hor91]. Importantly, the form of each learned neural network function can be visualized, thereby enabling a data-driven tool for user-guided conjecture of new mathematical equations that describe each separate term in $\mathscr{F}$. Moreover, formulating the equation learning task within the BINNs framework enables the modeler to use domain expertise to include qualitative constraints on the parameter networks (e.g. specifying nonlinear functions that are non-negative, monotone increasing/decreasing, etc.) by selecting appropriate activation functions and loss terms for the optimization.

While BINNs can be used to discover a wide range of governing equations across the biological and physical sciences, including systems of ODEs and PDEs, in this work they are demonstrated using reaction-diffusion PDEs. The BINNs methodology is first tested using synthetic data and then demonstrated on experimental data from scratch assay experiments with variable initial cell densities [Jin16a]. Notably, each data set is noisy and sparse, containing only five time measurements across 38 spatial locations. BINNs are used to discover the nonlinear forms of the diffusivity function and growth term of the governing reaction-diffusion equation. Persistent model discrepancy is used to motivate the incorporation of a novel delay term which may have important implications for the reproducibility and modeling of scratch assays. The learned nonlinear forms of the diffusion, growth, and delay terms are used to guide the selection of a mechanistic model with biologically interpretable parameters that remove virtually all of the model discrepancy.

### 4.2.1 Scratch assay data

Biologically-informed neural networks (BINNs) are evaluated on experimental scratch assay data from [Jin16a]. A typical scratch assay involves (i) growing a cell monolayer up to some desired initial cell density, (ii) creating a "scratch" in the interior of the monolayer to produce an empty region, and (iii) recording longitudinal measurements of the cell density

59 of 150

during re-colonization of the area. One-dimensional cell density profiles are obtained by manually counting the cells within vertical columns of the two-dimensional image data. See Figure C.1 in Appendix C for a visualization of the experiment. For these data, the cell density profiles were reported for six varying initial cell density levels (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). To make the cell density profiles compatible with neural network training, the data are pre-processed by rescaling the $x$ and $t$ variables to the scales of millimeters (mm) and days, respectively (see Methods Section for more details). Further, the cell density profile at the left boundary is removed from the data because it was identified as an outlier across each of the six data sets. The resulting pre-processed cell densities at 37 spatial points and five time points are shown in Figure 4.1.



**Figure 4.1** Experimental scratch assay data. Pre-processed cell density profiles from scratch assay experiments with varying initial cell densities [Jin16a]. Each subplot corresponds to an experiment with a different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). The cell densities are reported at 37 equally-spaced positions and five equally-spaced time points.

### 4.2.2 Biologically-informed neural networks

BINNs are centered around a function-approximating deep neural network, or MLP, denoted by $u_{\mathrm{MLP}}(x, t)$ which acts as a surrogate model that approximates the solution to the

governing equation described by Equation (4.2) (Figure 4.2A). In this work, the governing PDE is assumed to contain two terms, $\mathscr{D}$ and $\mathscr{G}$, that describe the general reaction-diffusion model in Equation (4.1). Since the true forms of the diffusivity and growth functions are unknown, they are approximated by neural networks $\mathscr{D} = D_{\mathrm{MLP}}(u)$ and $\mathscr{G} = G_{\mathrm{MLP}}(u)$ (Figure 4.2B). Both $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ are continuously differentiable functions that input the predicted cell density $u_{\mathrm{MLP}}(x, t)$ and output the corresponding diffusivity or growth value. The advantage of using MLPs for the terms of the governing PDE is that the nonlinear forms of these terms can be learned without specifying them explicitly (or as a library of candidate terms), thus circumventing the model specification problem. Automatic differentiation (Figure 4.2C) is used to numerically differentiate compositions of $u_{\mathrm{MLP}}$, $D_{\mathrm{MLP}}$, and $G_{\mathrm{MLP}}$ in order to construct the general reaction-diffusion model in Equation (4.1). The resulting PDE (Figure 4.2D) is used to regularize $u_{\mathrm{MLP}}$ during training so that $u_{\mathrm{MLP}}$ not only fits the data $u_{i,j}$ but also satisfies the governing reaction-diffusion system.



**Figure 4.2** Biologically-informed neural networks for reaction-diffusion models. (A) BINNs are deep neural networks that approximate the solution of a governing dynamical system. (B) By allowing the terms of the dynamical system (e.g. diffusivity function $\mathscr{D}$ and growth function $\mathscr{G}$) to be function-approximating deep neural networks, the nonlinear forms of these terms can be learned without the need to specify a mechanistic model or library of candidate terms. (C) Automatic differentiation is used on compositions of the different neural network models (e.g. $u$, $D$, and $G$) to construct the PDE that describes governing dynamical system. (D) The governing system is used in the neural network objective function to jointly learn and satisfy the governing PDE while minimizing the error between the network outputs and noisy observations.

To ensure that the fit to the data and the fidelity to the governing PDE are simultaneously optimized, the BINNs are trained with gradient-based methods using the following multi-part objective function:

$$\mathscr{L}_{\mathrm{Total}} = \mathscr{L}_{\mathrm{GLS}} + \mathscr{L}_{\mathrm{PDE}} + \mathscr{L}_{\mathrm{Constr}}. \tag{4.3}$$

The first term, $\mathscr{L}_{\text{GLS}}$, concerns the generalized least squares (GLS) distance between $u_{\text{MLP}}(x_i, t_j)$ and the corresponding observed data, $u_{i,j}$. The observation process is assumed to be described by a statistical error model of the form

$$u_{i,j} = u(x_i, t_j) + w_{i,j} \odot \varepsilon_{i,j}, \tag{4.4}$$

in which the measured data $u_{i,j}$ are a combination of the underlying dynamical system $u(x_i, t_j)$ and some random variable $w_{i,j} \odot \varepsilon_{i,j}$ where $\odot$ represents element-wise multiplication [Ban14b]. In general, the independent and identically distributed (i.i.d.) random variable $\varepsilon_{i,j}$ is modeled by an $n$-dimensional normal distribution with mean zero and variance one that is weighted by

$$w_{i,j} = \begin{bmatrix} \omega_1 u_1^\gamma(x_i, t_j) & \dots & \omega_n u_n^\gamma(x_i, t_j) \end{bmatrix}^{\text{T}}, \tag{4.5}$$

for $\gamma \geq 0$ and $\omega_1, \dots, \omega_n \in \mathbb{R}$ where $n$ is the dimensionality of the system. Note that (i) noiseless data are modeled by letting $\omega_1, \dots, \omega_n = 0$, (ii) constant-variance error used in ordinary least squares is modeled by letting $\gamma = 0$, $\omega_1, \dots, \omega_n = 1$, and (iii) non-constant-variance error (e.g. proportional error) used in generalized least squares is modeled by letting $\gamma > 0$, $\omega_1, \dots, \omega_n \neq 0$. Therefore, to account for the statistical error model in Equation (4.4), the GLS objective function

$$\mathscr{L}_{\text{GLS}} = \frac{1}{MN} \sum_{i=1, j=1}^{M,N} \left[ \frac{u_{\text{MLP}}(x_i, t_j) - u_{i,j}}{\left| u_{\text{MLP}}(x_i, t_j) \right|^\gamma} \right]^2, \tag{4.6}$$

is used with proportionality constant $\gamma = 0.2$. Note that $\gamma$ was tuned numerically following the methodology suggested in [Lag20d] (see Methods Section for more details).

The next term $\mathscr{L}_{\text{PDE}}$ ensures $u_{\text{MLP}}$ satisfies the solution of the governing PDE. For ease of notation, let $\hat{u}_{i,j} \equiv u_{\text{MLP}}(x_i, t_j)$, $\hat{D}_{i,j} \equiv D_{\text{MLP}}(u_{\text{MLP}}(x_i, t_j))$, and $\hat{G}_{i,j} \equiv G_{\text{MLP}}(u_{\text{MLP}}(x_i, t_j))$. Then for the reaction-diffusion equation, the error term takes the following form:

$$\mathscr{L}_{\text{PDE}} = \frac{1}{MN} \sum_{i=1, j=1}^{M,N} \left[ \underbrace{\frac{\partial \hat{u}_{i,j}}{\partial t}}_{\text{LHS}} - \underbrace{\left( \frac{\partial}{\partial x} \left( \hat{D}_{i,j} \frac{\partial \hat{u}_{i,j}}{\partial x} \right) + \hat{G}_{i,j} \hat{u}_{i,j} \right)}_{\text{RHS}} \right]^2, \tag{4.7}$$

where LHS and RHS denote the left-hand- and right-hand-sides of the governing PDE, respectively. Thus, by driving $\mathscr{L}_{\text{PDE}}$ to zero, the RHS is trained to match the LHS. Through this process, the nonlinear forms of $D_{\text{MLP}}$ and $G_{\text{MLP}}$ are learned despite not being directly observed. See the Methods Section for additional implementation details, including a

random sampling procedure that enforces this PDE constraint everywhere in the input domain during training.

Biological information and domain expertise are incorporated into the BINNs framework by adding penalties in the loss term $\mathscr{L}_{\text{Constr}}$. For the reaction-diffusion equation, the diffusivity and growth rates are assumed to be within biologically feasible ranges $[D_{\min}, D_{\max}]$ and $[G_{\min}, G_{\max}]$, respectively. Further, diffusion is also assumed to be non-decreasing and growth to be non-increasing with respect to cell density. The corresponding constraints take the form:

$$\mathscr{L}_{\text{Constr}} = \frac{1}{MN}\Bigg[ \sum_{\substack{i=1, j=1 \\ \hat{D} < D_{\min} \\ \hat{D} > D_{\max}}}^{M,N} \left(\hat{D}_{i,j}\right)^2 + \sum_{\substack{i=1, j=1 \\ \partial \hat{D}/\partial \hat{u} < 0}}^{M,N} \left(\frac{\partial \hat{D}_{i,j}}{\partial \hat{u}_{i,j}}\right)^2 \tag{4.8}$$
$$+ \sum_{\substack{i=1, j=1 \\ \hat{G} < G_{\min} \\ \hat{G} > G_{\max}}}^{M,N} \left(\hat{G}_{i,j}\right)^2 + \sum_{\substack{i=1, j=1 \\ \partial \hat{G}/\partial \hat{u} > 0}}^{M,N} \left(\frac{\partial \hat{G}_{i,j}}{\partial \hat{u}_{i,j}}\right)^2 \Bigg].$$

The maximum and minimum diffusivity and growth rates considered in [Jin16a] were used to force $D_{\text{MLP}}$ and $G_{\text{MLP}}$ to stay within biologically realistic ranges. The constraints on $D_{\text{MLP}}$ and $G_{\text{MLP}}$ shown in Equation (4.8) were used for all computational results in this work. See the Methods Section for biological motivations and numerical implementation details of these constraints.

### 4.2.3   Evaluation procedure

Because the model prediction, $u_{MLP}(x,t)$, is only a surrogate model for the dynamical system, $u(x,t)$, it is possible that this approximation may contain errors, particularly in areas where the PDE constraint given by Equation (4.7) is not satisfied. To ensure that the inferred diffusion and growth terms lead to biologically realistic dynamics, the reaction-diffusion equation given by Equation (4.1) is solved numerically with a method-of-lines approach using $\mathscr{D} = D_{\text{MLP}}$ and $\mathscr{G} = G_{\text{MLP}}$. Note that this model is well-defined because $D_{\text{MLP}}$ and $G_{\text{MLP}}$ are continuously differentiable functions of the cell density, $u$. Further, BINNs are retrained multiple times for each data set in which the forward simulation using the learned PDE terms that yields the smallest GLS error (Equation (4.6)) is saved. All fits to the data shown in the Results Section are numerical solutions to the PDE in Equation (4.1) using the learned diffusivity and growth functions. See the Methods Section for numerical implementation details of the PDE forward solver.

## 4.3  Methods

All methods herein were implemented in Python 3.6.8 using the PyTorch 1.2.0 deep learning library. All data and code are made publicly available at `https://github.com/jlager/binns`. The following section is intended to make BINNs feasible for a wide range of biological applications. In particular, this section covers (i) the importance of data pre-processing, (ii) strategies for using real-world knowledge to design effective neural network models, (iii) the complete training protocol ranging from selecting appropriate statistical error models and hyperparameters to balancing the multi-objective error function, and (iv) numerical implementation details for forward solving BINNs-guided PDEs.

### 4.3.1  Data pre-processing

Input and output standardization are common practice to stabilize neural network training [The09]. Since the scratch assay data in [Jin16a] reported cell densities on the order of $u = \mathcal{O}(10^{-3})$ cells/$\mu$m$^2$ at spatial locations on the order of $x = \mathcal{O}(10^3)$ $\mu$m for time points on the order of $t = \mathcal{O}(10)$ hours, these variables needed to be standardized. Without standardization, the neural network models failed to converge for these data because (i) the network inputs ($x$ and $t$) differed by several orders of magnitude from each other and (ii) the network inputs ($x$ and $t$) and outputs ($u$) also differed by several orders of magnitude. By rescaling $x$ and $t$ to millimeters (mm) and days, respectively, the adjusted variables ranged from $x = \mathcal{O}(1)$ mm, $t = \mathcal{O}(1)$ days, and cell density $u = \mathcal{O}(10^3)$ cells/mm$^2$. Standardizing $x$ and $t$ addressed (i) while (ii) is addressed by using scaling factors discussed in the following section. The cell density profile at the left boundary was removed since it was consistently larger than the remaining cell densities across all six data sets.

### 4.3.2  Network design

BINNs are centered around $u_{\text{MLP}}$, a function-approximating multilayer perceptron (MLP) (also known as an artificial neural network). MLPs, like polynomials [Sto48], are in the class of *universal function approximators*, meaning that they can approximate any continuous bounded functions on a closed interval arbitrarily well under some reasonable assumptions [Hor91]. For the scratch assay data in the present work, $u_{\text{MLP}}$ inputs spatiotemporal vectors $x = \begin{bmatrix} x, t \end{bmatrix}$ and outputs the corresponding approximations to the cell density $u$. To give $u_{\text{MLP}}$ sufficient capacity to approximate the solution to the governing PDE, the network is chosen to have three hidden layers with 128 neurons in each layer, resulting in a model

with approximately 30,000 total parameters. Concretely, $u_{\mathrm{MLP}}$ takes the form

$$u_{\mathrm{MLP}}(\boldsymbol{x}) = \alpha \cdot \phi\Big( \sigma\big( \sigma\big( \sigma(\boldsymbol{x} W_1 + b_1) W_2 + b_2 \big) W_3 + b_3 \big) W_4 + b_4 \Big), \qquad (4.9)$$

where the trainable parameters $W_i$ and $b_i$ denote weight matrices and bias vectors for the $i^{\mathrm{th}}$ layer, $\sigma(\cdot)$ and $\phi(\cdot)$ denote nonlinear activation functions, and $\alpha$ denotes a scaling factor. Each hidden layer uses a "sigmoid" activation function (i.e. $\sigma(x) = 1/(1 + e^{-x})$) while the output layer uses a "softplus" activation function (i.e. $\phi(x) = \ln(1 + e^x)$). The softplus activation function is a particular design choice since it is a continuously differentiable function that forces the predicted cell densities to be non-negative, and has been previously shown to be well-suited for biological transport models [Lag20d]. Finally, to account for the difference in scale between the inputs ($x, t = \mathcal{O}(1)$) and outputs ($u = \mathcal{O}(10^3)$), the MLP outputs are post-multiplied by the experimentally validated carrying capacity (i.e. $\alpha = 1.7 \times 10^3$) from [Jin16a]. Note that in practice, if values like this are unknown, one can simply let $\alpha$ be the maximum observed cell density or some other similar quantity. The key here is to ensure the orders of magnitude between the network inputs and outputs are similar so that the parameters of the MLP do not have to account for the change of scale [The09].

The terms of the governing PDE are modeled with neural networks $D_{\mathrm{MLP}}(u_{\mathrm{MLP}})$, $G_{\mathrm{MLP}}(u_{\mathrm{MLP}})$, etc. All of these MLPs share the same number of layers as $u_{\mathrm{MLP}}$ but use 32 neurons per layer. These networks are chosen to be smaller for both computational efficiency and because the parameter dynamics are assumed to be simpler than the cell density dynamics $u$. The hidden layers use sigmoid activation functions. The output layer for $D_{\mathrm{MLP}}$ uses a softplus activation because diffusion is assumed to be non-negative for all cell densities. Since the growth term can be negative (e.g. logistic growth when the cell density exceeds the carrying capacity), a linear output (i.e. no activation function) is used in the final layer for $G_{\mathrm{MLP}}$. A discovered third term (see the Results Section) of the governing PDE, denoted $T_{\mathrm{MLP}}$, uses the sigmoid activation function for the output layer to constrain its outputs to $(0, 1)$. Finally, as with $u_{\mathrm{MLP}}$, the inputs and outputs of $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ are also standardized. In particular, the inputs of both networks (i.e. $u_{\mathrm{MLP}}$) are divided by the carrying capacity $K = 1.7 \times 10^3$ while the outputs of $D_{\mathrm{MLP}}$ are multiplied by 0.096 $^{\mathrm{mm}^2}/_{\mathrm{day}}$ and the outputs of $G_{\mathrm{MLP}}$ are multiplied by 2.4 $^1/_{\mathrm{day}}$. These values were the maximum diffusion and growth values considered in [Jin16a]. Similar to $u_{\mathrm{MLP}}$, the input and output scaling factors ensure the MLP parameters do not have to account for changes in scale. No standardization was used for $T_{\mathrm{MLP}}$ since it is a function of $t$, meaning its inputs and outputs of the same order (i.e. $\mathcal{O}(1)$).

### 4.3.3 Training procedure

The BINN parameters (i.e. weights and biases of $u_{\mathrm{MLP}}$, $D_{\mathrm{MLP}}$, $G_{\mathrm{MLP}}$, and $T_{\mathrm{MLP}}$) are optimized using the first-order gradient-based Adam optimizer [Kin17] with default hyper-parameters and minibatch-optimization. To prevent over-fitting, the scratch assay data were randomly partitioned into 80%/20% training and validation sets. The network parameters were updated iteratively to minimize $\mathcal{L}_{\mathrm{Total}}$ in Equation (4.3) on the training set and saved on *relative* improvement in validation error. In other words, the model parameters were saved if the relative difference between (i) the validation error in the current iteration and (ii) the smallest recorded validation error exceeded 5%. Finally, since the parameters of each BINN are randomly initialized and applied to different data sets, early stopping of 5,000 (i.e. training was stopped if the relative validation error did improve for 5,000 consecutive epochs) was used to guarantee the convergence of each BINN independently. The implementation details of each term in $\mathcal{L}_{\mathrm{Total}}$ (i.e. $\mathcal{L}_{\mathrm{GLS}}$, $\mathcal{L}_{\mathrm{PDE}}$, and $\mathcal{L}_{\mathrm{Constr}}$) are discussed in more detail below.

The first term, $\mathcal{L}_{\mathrm{GLS}}$, in Equation (4.6) corresponds to the generalized least squares (GLS) distance between $u_{\mathrm{MLP}}$ and the observation data $u_{i,j}$. Since the error process is assumed to be i.i.d., the parameters of the statistical model in Equation (4.4) (i.e. $\gamma$) must first be calibrated. Following [Lag20d], $u_{\mathrm{MLP}}$ is trained using $\mathcal{L}_{\mathrm{GLS}}$ as an objective function for $\gamma = 0.0, 0.2, 0.4, 0.6$ (recall that $\gamma = 0.0$ represents the ordinary least squares case) for each data set. After qualitative assessment of the modified residual errors (see Figure C.9 in Appendix C), $\gamma = 0.2$ was identified as the value that produced the most i.i.d. residuals across each of the six data sets. Using the calibrated statistical error model, $\mathcal{L}_{\mathrm{GLS}}$ is evaluated at each training iteration using mini-batches (i.e. randomly selected subsets) of input/output data. In general, using a small batch size acts as an additional form of regularization that helps neural networks escape local minima during training and allows for better generalization [Kes17]. However, this significantly increases the computational cost of training due to the increased number of training iterations needed to converge. Therefore, BINNs were trained using mini-batches of size 37 (i.e. 1/4 the number of points in the training set) which was found to balance the accuracy and computational cost.

To ensure $u_{\mathrm{MLP}}$ satisfies the solution of the governing PDE, the terms $\mathcal{L}_{\mathrm{PDE}}$ in Equation (4.7) and $\mathcal{L}_{\mathrm{Constr}}$ in Equation (4.8) are included in $\mathcal{L}_{\mathrm{Total}}$ as a form of regularization. However, since the scratch assay data are sparse, simply training $u_{\mathrm{MLP}}$ using $\mathcal{L}_{\mathrm{Total}}$ at the observed data locations can result in unrealistic dynamics in between data points. Therefore, to ensure $u_{\mathrm{MLP}}$ satisfies the solution of a governing PDE *everywhere* in the input domain, $\mathcal{L}_{\mathrm{PDE}}$ and $\mathcal{L}_{\mathrm{Constr}}$ are evaluated at 10,000 uniformly randomly sampled points $x_i \in [x_{\min}, x_{\max}]$ and $t_j \in [t_{\min}, t_{\max}]$ at *each* training iteration. Without the random sampling

procedure, $u_\text{MLP}$ can severely overfit to the data. To illustrate the importance of the random sampling procedure, the model fits, GLS errors, and PDE errors are shown in Figure C.10 in Appendix C for three cases in which (i) no PDE regularization is used, (ii) PDE regularization is used at the data locations, and (iii) PDE regularization is used at 10,000 randomly sampled points. In particular, Figure C.10 shows that in option (i) $u_\text{MLP}$ overfits the data practically everywhere in the input domain, (ii) $u_\text{MLP}$ overfits everywhere except at the data locations (see vertical lines in third subplot of row b), and (iii) the random sampling procedure results in the smallest amount of PDE error and the largest amount of GLS error. The desired behavior is shown in option (iii) since $u_\text{MLP}$ fits the data as accurately as allowed by the governing PDE.

The third error term, $\mathscr{L}_\text{Constr}$, constrains $D_\text{MLP}$, $G_\text{MLP}$, and $T_\text{MLP}$ to exhibit biologically re-alistic values and dynamics. Choosing appropriate constraints can be ambiguous when the relevant literature gives conflicting suggestions. For example, when designing a derivative constraint for the diffusivity network $D_\text{MLP}$, [Ang09] suggest that diffusion should decrease with cell density due to cell-to-cell adhesion whereas [Nar16] suggest the opposite in which cells promote the migration of others. To mitigate this, BINNs were trained without any constraints on $D_\text{MLP}$ and $G_\text{MLP}$ in order to visualize the collective behavior of the parameter networks (see Figure C.8 in Appendix C). Note that $T_\text{MLP}$ was still forced to be non-decreasing (see the Results Section). The network evaluations in Figure C.8 showed unrealistic parame-ter dynamics for some data sets, but their collective behavior was used to design derivative constraints that forced $D_\text{MLP}$ to increase as a function of cell density and $G_\text{MLP}$ to decrease with cell density for the set of scratch assay data considered in this work. Concretely, the diffusion term $D_\text{MLP}$ was constrained to values between 0.0 and 0.096 $^\text{mm}^2$/day and the growth term $G_\text{MLP}$ to values between $-0.48$ and 2.4 $^1$/day. The maximum and minimum diffusion values and maximum growth value were chosen based on values used in [Jin16a]. The minimum growth value was chosen to be negative 20% of the maximum growth value to allow $G_\text{MLP}$ to output negative values for cell densities near the carrying capacity if needed. The sigmoid output activation function for the delay term $T_\text{MLP}$ constrained its outputs to between 0 and 1. Derivative terms were used in $\mathscr{L}_\text{Constr}$ to constrain $D_\text{MLP}$ and $T_\text{MLP}$ to be non-decreasing and $G_\text{MLP}$ to be non-increasing. For ease of notation, let $\hat{u}_{i,j} \equiv u_\text{MLP}(x_i, t_j)$, $\hat{D}_{i,j} \equiv D_\text{MLP}(u_\text{MLP}(x_i, t_j))$, $\hat{G}_{i,j} \equiv G_\text{MLP}(u_\text{MLP}(x_i, t_j))$, and $\hat{T}_{i,j} \equiv T_\text{MLP}(t_j)$, then the constraint

term can be written concretely as

$$\mathscr{L}_{\text{Constr}} = \frac{1}{MN}\left[\alpha_1 \sum_{\substack{i=1,j=1 \\ \hat{D}<0.0 \\ \hat{D}>0.096}}^{M,N} \left(\hat{D}_{i,j}\right)^2 + \alpha_2 \sum_{\substack{i=1,j=1 \\ \partial\hat{D}/\partial\hat{u}<0}}^{M,N} \left(\frac{\partial\hat{D}_{i,j}}{\partial\hat{u}_{i,j}}\right)^2 + \right. \qquad (4.10)$$

$$\left. \alpha_3 \sum_{\substack{i=1,j=1 \\ \hat{G}<-0.48 \\ \hat{G}>2.4}}^{M,N} \left(\hat{G}_{i,j}\right)^2 + \alpha_4 \sum_{\substack{i=1,j=1 \\ \partial\hat{G}/\partial\hat{u}<0}}^{M,N} \left(\frac{\partial\hat{G}_{i,j}}{\partial\hat{u}_{i,j}}\right)^2 + \alpha_5 \sum_{\substack{i=1,j=1 \\ \partial\hat{T}/\partial\hat{t}<0}}^{M,N} \left(\frac{\partial\hat{T}_{i,j}}{\partial\hat{u}_{i,j}}\right)^2 \right].$$

Since the parameter networks and their derivatives occur at different scales with respect to each other and with respect to the error terms $\mathscr{L}_{\text{GLS}}$ and $\mathscr{L}_{\text{PDE}}$, each term of Equation (4.10) is weighted by a factor $\alpha_i$. In particular, each constraint is weighted based on the input/output scaling factors of the corresponding neural network (see Network Design subsection). Concretely, the terms in Equation (4.10) are weighted by $\alpha_1 = 1/0.096 \times 10^{10}$, $\alpha_2 = K/0.096 \times 10^{10}$, $\alpha_3 = 1/2.4 \times 10^{10}$, $\alpha_4 = K/2.4 \times 10^{10}$, and $\alpha_5 = 10^{10}$. Note that the weight factors for the derivative constraints on $D_{\text{MLP}}$ and $G_{\text{MLP}}$ (i.e. $\alpha_2$ and $\alpha_4$) include the carrying capacity $K = 1.7 \times 10^3$ since $K$ was used as an input scaling factor for these networks. The factor $10^{10}$ was chosen large enough to guarantee that $D_{\text{MLP}}$, $G_{\text{MLP}}$, and $T_{\text{MLP}}$ exhibited the desired behavior. Boundary conditions can also be included in the $\mathscr{L}_{\text{Constr}}$ term, however, since they were unknown for the scratch assay data considered in this work, no boundary conditions were used to train $u_{\text{MLP}}$.

Finally, the GLS errors at the initial condition (i.e. data locations where $t = 0$) were weighted by a factor of 10 during training. This was found to improve the generalization accuracy of $D_{\text{MLP}}$, $G_{\text{MLP}}$, and $T_{\text{MLP}}$ when evaluated using a numerical PDE solver. The reason for this is because the cell density at $t = 0$ may not satisfy a governing dynamical system since the measurement is taken directly after the scratch assay protocol is performed [Jin16a]. However, the initial condition "sets the stage" for the governing dynamics to drive the temporal evolution of the system. Therefore, by weighting the initial condition more heavily in $\mathscr{L}_{\text{GLS}}$, the PDE error term $\mathscr{L}_{\text{PDE}}$ must conform $u_{\text{MLP}}$ to satisfy the governing system for $t > 0$ as dictated by $u_{\text{MLP}}$ at $t = 0$. This step forced $D_{\text{MLP}}$, $G_{\text{MLP}}$, and $T_{\text{MLP}}$ to learn more generalizable representations of the diffusivity, growth, and delay functions, respectively. The weighting factor was numerically validated using the mean GLS error across each scratch assay experiment for weighting factors 1, 10, and $10^2$. Note that this weighting factor makes BINNs sensitive to the random choice of training/validation split, since some data points in the initial condition may be more informative than others for equation learning and ultimate model generalizability. This observation was also noted in

a recent equation learning study in which the random split of training and validation sets was found to influence the structure of the learned equation [Lag20d]. Adopting a strategy similar to this previous study, BINNs were trained 20 times for each data set (using different random training/validation splits). The BINN for which the numerical simulations resulted in the smallest GLS error was saved as the best model.

### 4.3.4 PDE forward solver

The numerical implementation details are provided for systems describing quantity of interest $u(x, t)$ that are governed by the following equation:

$$u_t = \left(Q(u, u_x, t)\right)_x + F(u),$$ (4.11)
$$u(x, t_0) = \phi(x),$$
$$u_x(x_0, t) = u_x(x_f, t) = 0,$$

for $x \in [x_0, x_f]$, and $t \in [t_0, t_f]$. Note that the reaction-diffusion model in Equation (4.1) is an example of Equation (4.11) where $Q(u, u_x, t) = \mathscr{D}(u, t)u_x$ and $F(u) = \mathscr{G}(u, t)u$. In Equation (4.11), the initial condition is denoted by $\phi(x)$ and the boundary conditions are assumed to be no-flux boundary conditions. Note that the no-flux condition represents a zero net flux boundary condition which does not preclude cells moving across the boundary, but instead reflects the situation in which the flux in the positive and negative $x$-directions are equal, giving rise to zero total flux. The spatial and temporal domains are discretized into equispaced grids as:

$$x_i = i\Delta x, \quad t_j = j\Delta t,$$ (4.12)

for $i = 0, \ldots, 200$ and $j = 0, \ldots, 1,000$. For notational convenience, let $u_i(t) = u(x_i, t)$. Then, the method-of-lines approach is used to solve Equation (4.11) with the numerical discretization from [Kur00] that is given by

$$\left(Q(u, u_x, t)\right)_x \approx \frac{P_{i+1/2}(t) - P_{i-1/2}(t)}{\Delta x},$$ (4.13)

where $P_{i+1/2}(t)$ is an estimate for the rightwards diffusive flux at location $x_i$ that is given by

$$P_{i+1/2}(t) = \frac{1}{2}\left[Q\left(u_i(t), \frac{u_{i+1}(t) - u_i(t)}{\Delta x}, t\right) + Q\left(u_{i+1}(t), \frac{u_{i+1}(t) - u_i(t)}{\Delta x}, t\right)\right].$$ (4.14)

The no-flux boundary conditions at $x_0$ and $x_{200}$ are implemented by incorporating the ghost points $x_{-1}$ and $x_{201}$ satisfying $u_{-1}(t) = u_1(t)$ and $u_{201}(t) = u_{199}(t)$. The Scipy integration subpackage (version 1.4.1) is used to integrate Equation (4.11) over time using an explicit fourth order Runge-Kutta Method.

### 4.3.5 Parameter estimation

The parameters of each mechanistic model were optimized using the Limited-memory BFGS algorithm with bound constraints (L-BFGS-B) in Python's Scipy package with default tolerance values to minimize the generalized least squares error function in Equation (4.6) with the adjusted statistical error model in Equation (4.4) with $\gamma = 0.2$. The parameters for Equations (4.19) and (4.20) were initialized using the values from [Jin16a]. The parameters for Equation (4.17) were initialized by fitting each PDE term in Equations (4.18a)-(4.18c) to the corresponding parameter network solutions in Figure 4.7 using ordinary least squares. Finally, the diffusivity and growth function parameters were bounded using $D_{\min} = 0$ mm²/day, $D_{\max} = 0.096$ mm²/day, $m_{\min} = 0$, $m_{\max} = 4$, $r_{\min} = 0$ 1/day, and $r_{\max} = 2.4$ 1/day (all of which come from [Jin16a]), while the delay function parameters $\beta_0$ and $\beta_1$ were bounded by $[-10, 10]$.

## 4.4 Results

### 4.4.1 Simulation case study

Since the diffusivity and growth terms are inferred by BINNs through learning $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$, respectively, the ability of BINNs to learn biologically accurate representations of these terms must first be tested. To investigate this, data were simulated using the classical FKPP and Generalized Porous-FKPP equations with parameter values from [Jin16a] for the scratch assay data with initial cell density 20,000 cells per well. Additionally, the simulated data were obscured with artificial observation error using the statistical model in Equation (4.4) with $\gamma = 0.2$. Each simulation used the initial condition from the scratch assay data with initial cell density 20,000 cells per well. Using the same level of sparsity (i.e. 37 spatial points and five time points), the BINNs framework was shown to (i) approximate the dynamical system accurately and (ii) approximate the general forms of the diffusivity and growth terms. See Figures C.2 and C.3 in Appendix C for the model and parameter fits, respectively. This case study demonstrates that BINNs are able to learn accurate representations of the diffusivity and growth functions from biologically realistic noisy sparse data, however, further analysis, like model selection and comparison, is omitted here and instead explored

70

using experimental data.

## 4.4.2 Reaction-diffusion BINNs for experimental data

As described in the previous sections, the diffusivity and growth functions are approximated by deep neural networks, $\mathscr{D} = D_{\mathrm{MLP}}(u)$ and $\mathscr{G} = G_{\mathrm{MLP}}(u)$, resulting in a governing PDE of the form

$$u_t = (D_{\mathrm{MLP}}(u)u_x)_x + G_{\mathrm{MLP}}(u)u, \tag{4.15}$$

where $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ are functions of the cell density $u$. A BINN was trained for each data set with varying initial cell density. The resulting numerical PDE solutions using the trained $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ are shown in Figure 4.3.



**Figure 4.3** Reaction-diffusion BINN solutions. Predicted cell density profiles using BINNs with the governing reaction-diffusion PDE in Equation (4.15). Each subplot corresponds to an experiment with a different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Solid lines represent the numerical solution to Equation (4.15) using $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$. The markers represent the experimental scratch assay data.

While the model fits shown in Figure 4.3 are excellent for lower initial cell densities, there still remains a significant amount of model discrepancy at higher initial cell densities. GLS residual errors were computed to provide an additional way of visualizing the model

discrepancy (see Figure C.4 in Appendix C) in which non-i.i.d. residuals are clearly present at higher initial cell densities. To investigate the specific form of the model discrepancy, Figure 4.4 shows the learned diffusivity and growth functions with the corresponding model fit for the data set with an initial cell density of 20,000 cells per well.



**Figure 4.4** Reaction-diffusion BINN terms and discrepancy. Left: learned diffusivity and growth functions, $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$, evaluated over cell density, $u$. Right: Predicted cell density profiles using BINNs with the governing reaction-diffusion PDE in Equation (4.15) for data with initial cell density 20,000 cells per well. Solid lines represent the numerical solution to Equation (4.15) using $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$. The markers represent the experimental scratch assay data.

Figure 4.4 reveals clear model discrepancy in two main areas: (i) at high cell densities (i.e. $x \in [0, 0.25]$ mm and $x \in [1.75, 2.0]$ mm for $t \in [0, 1]$ days) where diffusion is negligible and the dynamics are governed primarily by growth; and (ii) at low cell densities (i.e. $x \in [0.5, 1.5]$ mm for $t \in [0, 1]$ days) where growth is negligible and the dynamics are primarily governed by diffusion. In particular, the discrepancy is largest for early time points where the diffusion and growth dynamics appear too rapid. The solutions of $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ are also qualitatively similar to the classical FKPP equation in which the learned diffusivity function is relatively constant while the learned growth function is approximately linearly decreasing with cell density, $u$. However, despite $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ learning biologically realistic functions for the diffusivity and growth, the persistent model discrepancy observed across multiple data sets with high initial cell densities (see Figure 4.3) suggests that the reaction-diffusion equation described in Equation (4.15) may be insufficient to fully capture the underlying dynamics of cell migration for these data. From a mathematical modeling perspective, the model discrepancy at early time points suggests the existence of a time delay that scales the

magnitude of the density-dependent diffusion and growth rates. Biological reasons behind this phenomenon may include cell damage from the scratch assay protocol or changes in cell functions where more cells become immobile/non-proliferative as the cell density approaches carrying capacity [Dyd20; Pou95; Neu20]. See the Discussion Section for more details.

### 4.4.3 Delay-reaction-diffusion BINNs for experimental data

Motivated by the model discrepancy for data sets with high initial cell density, the reaction-diffusion equation in Equation (4.15) was modified by including a time delay described by an additional neural network function $T_{\mathrm{MLP}}(t)$. The new term $T_{\mathrm{MLP}}(t)$ is a continuously differentiable function of time that is constrained to be non-decreasing and output values between 0 and 1. In this way, $T_{\mathrm{MLP}}$ can scale the strength of the density-dependent diffusivity and growth terms in time. Letting the diffusivity, $\mathscr{D}$, and growth, $\mathscr{G}$, terms of the governing PDE be functions of $u$ and $t$, they are replaced with $\mathscr{D} = T_{\mathrm{MLP}}(t)D_{\mathrm{MLP}}(u)$ and $\mathscr{G} = T_{\mathrm{MLP}}(t)G_{\mathrm{MLP}}(u)$. This results in a governing PDE of the form

$$u_t = \big(T_{\mathrm{MLP}}(t)D_{\mathrm{MLP}}(u)u_x\big)_x + T_{\mathrm{MLP}}(t)G_{\mathrm{MLP}}(u)u,$$

which simplifies to

$$u_t = T_{\mathrm{MLP}}(t)\Big((D_{\mathrm{MLP}}(u)u_x)_x + G_{\mathrm{MLP}}(u)u\Big), \tag{4.16}$$

where $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ are functions of the cell density $u$ and $T_{\mathrm{MLP}}$ is a function of time $t$. Note that $T_{\mathrm{MLP}}$ was chosen to be separable from $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ since the density-dependent dynamics of diffusion and growth are assumed to be consistent throughout time. Further, it was assumed that both $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ are scaled by the same time delay; see the Discussion Section for more details. BINNs governed by the PDE in Equation (4.16) were trained for each data set with varying initial cell density. The resulting forward simulations using the trained $T_{\mathrm{MLP}}$, $D_{\mathrm{MLP}}$, and $G_{\mathrm{MLP}}$ networks are shown in Figure 4.5.

The model fits shown in Figure 4.5 demonstrate that virtually all of the model discrepancy across each initial cell density was removed by including a time delay. This is confirmed further using GLS residual errors (see Figure C.5 in Appendix C) where the residuals are approximately i.i.d. even at higher initial cell densities. Similar to the reaction-diffusion case, Figure 4.6 shows the learned diffusivity, growth, and delay functions with the corresponding model fit for the data set with initial cell density of 20,000 cells per well.

Figure 4.6 shows that the model discrepancy in areas with high and low cell densities

**Figure 4.5** Delay-reaction-diffusion BINN solutions. Predicted cell density profiles using BINNs with the governing delay-reaction-diffusion PDE in Equation (4.16). Each subplot corresponds to an experiment with a different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Solid lines represent the numerical solution to Equation (4.16) using $T_{\text{MLP}}$, $D_{\text{MLP}}$, and $G_{\text{MLP}}$. The markers represent the experimental scratch assay data.

at early time points has been practically eliminated. This is most clearly seen in the delay-reaction-diffusion model solution at the second time point (i.e. $t = 0.5$ days), which matches the data more accurately than the reaction-diffusion model in Equation (4.15) at the same time point (see Figure 4.4). Moreover, $D_{\text{MLP}}$ and $G_{\text{MLP}}$ for the delay-reaction-diffusion BINN learned similar forms of the diffusivity and growth compared to the reaction-diffusion case. However, the delay term $T_{\text{MLP}}$ reveals that the diffusion and growth dynamics described by $D_{\text{MLP}}$ and $G_{\text{MLP}}$ are scaled down for early time points (i.e. $t < 1$) before $T_{\text{MLP}}$ converges to 1, allowing $D_{\text{MLP}}$ and $G_{\text{MLP}}$ to come into full effect. This observation is of particular importance since the majority of scratch assay data are reported within this time delay region (i.e. 4, 6, 12, or 24 hrs) [Aok17; Arc11; Bin07; Mat04]. Importantly, not accounting for a time delay within this region may potentially explain why scratch assay experiments are notoriously difficult to reproduce [Jin16a].

### 4.4.4 Guided mechanistic model selection

The diffusion, growth, and delay networks, $D_{\text{MLP}}$, $G_{\text{MLP}}$, and $T_{\text{MLP}}$, were used to guide the selection of biologically realistic mechanistic models for downstream use in a traditional

**Figure 4.6** Delay-reaction-diffusion BINN terms and discrepancy. Left: learned diffusivity and growth functions, $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$, evaluated over cell density, $u$, and delay function, $T_{\mathrm{MLP}}$, evaluated over time, $t$. Right: Predicted cell density profiles using BINNs with the governing delay-reaction-diffusion PDE in Equation (4.16) for data with initial cell density 20,000 cells per well. Solid lines represent the numerical solution to Equation (4.16) using $D_{\mathrm{MLP}}$, $G_{\mathrm{MLP}}$, and $T_{\mathrm{MLP}}$. The markers represent the experimental scratch assay data.

mathematical modeling framework. Each network solution corresponding to the six scratch assay data sets is shown in Figure 4.7.



**Figure 4.7** Delay-reaction-diffusion BINN terms. The learned diffusivity, $D_{\mathrm{MLP}}$, growth, $G_{\mathrm{MLP}}$, and delay, $T_{\mathrm{MLP}}$, functions extracted from the corresponding BINNs with governing delay-reaction-diffusion PDE in Equation (4.16). Each line corresponds to an experiment with a different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Note that $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ have different lengths since they are evaluated between the minimum and maximum observed cell densities corresponding to each data set.

From Figure 4.7, the learned diffusivities for each experiment with different initial

cell density are non-zero when $u = 0$ and appear increasing and concave up (like power laws) with respect to the cell density $u$. On the other hand, the learned growth terms are approximately linear, which is consistent with logistic models, and the learned delay terms all exhibit sigmoidal dynamics. Note that the outlying $G_{\text{MLP}}$ solution for the scratch assay data set with 10,000 initial cells per well is likely an artifact of the observed cell densities in that experiment not approaching the carrying capacity, and therefore leading to unrealistic learned dynamics. Based on qualitative analysis of these plots, the following mechanistic delay-reaction-diffusion equation is proposed to satisfy each scratch assay data set:

$$u_t = \mathcal{T}(t)\Big((\mathcal{D}(u)u_x)_x + \mathcal{G}(u)u\Big), \tag{4.17}$$

with diffusivity, growth, and delay functions

$$\mathcal{D} = D_0 + D\left(\frac{u}{K}\right)^m, \tag{4.18a}$$

$$\mathcal{G} = r\,u\left(1 - \frac{u}{K}\right), \tag{4.18b}$$

$$\mathcal{T} = \frac{1}{1 + e^{-(\beta_1 t + \beta_0)}}, \tag{4.18c}$$

respectively. The diffusivity function $\mathcal{D}$ in Equation (4.18a) is a combination of the classical FKPP and Generalized Porous-FKPP diffusivity function, with baseline cell diffusivity $D_0$, diffusion coefficient $D$, and exponent $m$. The growth function $\mathcal{G}$ in Equation (4.18b) is chosen to be the logistic growth function with intrinsic growth rate $r$ and carrying capacity $K$. The delay function $\mathcal{T}$ in Equation (4.18c) is represented by the logistic regression function with parameters $\beta_0$ and $\beta_1$. One advantage of using a mathematical model with specified functional forms and parameters described by Equation (4.18a)-(4.18c) is that standard parameter estimation techniques can now be used. This enables a comparison of the BINN-guided model in Equation (4.17) to other mechanistic models, namely, the classical FKPP and Generalized Porous-FKPP equations.

### 4.4.5 Model comparison

The BINN-guided delay-reaction-diffusion model in Equation (4.17) was compared to the classical FKPP equation

$$u_t = (D\,u_x)_x + r\,u\left(1 - \tfrac{u}{K}\right), \tag{4.19}$$

with diffusion coefficient $D$, intrinsic growth rate $r$, and carrying capacity $K$ and Generalized Porous-FKPP equation

$$u_t \quad = (D(\tfrac{u}{K})^m u_x)_x + r\, u \left(1 - \tfrac{u}{K}\right), \tag{4.20}$$

with additional exponent $m$. These models were used as a baseline for comparison since they have been identified as the current state-of-the-art in modeling these data [Jin16a; War19]. The parameters of each model were optimized numerically using the generalized least squares error function in Equation (4.6) with the adjusted statistical error model in Equation (4.4) with $\gamma = 0.2$. Note that the carrying capacity was fixed at $K = 1.7 \times 10^3$ cells/mm² and not optimized because it was empirically validated in [Jin16a]. The resulting model fits and parameter values for the classical FKPP and Generalized Porous-FKPP models are shown in Figures C.6 and C.7 and Tables C.1, and C.2 in Appendix C. The solutions of the BINN-guided delay-reaction-diffusion model in Equation (4.17) to each data set are shown in Figure 4.8.



**Figure 4.8** BINN-guided delay-reaction-diffusion model solutions. Predicted cell density profiles using the delay-reaction-diffusion model in Equation (4.17). Each subplot corresponds to an experiment with a different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Solid lines represent the numerical solution to Equation (4.17) using the parameters that minimize $\mathcal{L}_{\mathrm{GLS}}$ in Equation (4.6). The markers represent the experimental scratch assay data.

The predicted cell density profiles in Figure 4.8 closely matched the scratch assay data which suggests that the proposed model in Equation (4.17) with Equations (4.18a)-(4.18c) successfully captured the learned dynamics from $T_{\mathrm{MLP}}$, $D_{\mathrm{MLP}}$, and $G_{\mathrm{MLP}}$. The optimized parameter values across each data set are shown in Table 4.1. Note that the parameters were rescaled to $\mu$m and hours (hr) for comparison with [Jin16a] and [War19].

**Table 4.1** BINN-guided delay-reaction-diffusion model parameters.

| Parameter | Initial cell density | | | | | |
|---|---|---|---|---|---|---|
| | 10,000 | 12,000 | 14,000 | 16,000 | 18,000 | 20,000 |
| $D_0$ ($\mu$m$^2$/hr) | 95.7 | 353.3 | 482.1 | 604.3 | 804.0 | 675.8 |
| $D$ ($\mu$m$^2$/hr) | 3987.1 | 3166.4 | 3775.0 | 3773.8 | 2201.8 | 1954.9 |
| $m$ (unitless) | 1.5976 | 3.4708 | 1.9060 | 3.5173 | 3.2204 | 0.9876 |
| $r$ (1/hr) | 0.0525 | 0.0714 | 0.0742 | 0.0798 | 0.0772 | 0.0951 |
| $\beta_0$ (unitless) | -1.0292 | -3.3013 | -3.1953 | -2.9660 | -1.2695 | -4.0651 |
| $\beta_1$ (1/hr) | 0.2110 | 0.2293 | 0.2761 | 0.2180 | 0.1509 | 0.4166 |

Table of model parameters for Equation (4.17) calibrated for each scratch assay data set. Each column corresponds to an experiment with different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well).

Table 4.1 reveals that many of the parameters relating to density-dependent diffusion and growth show trends (e.g. $D_0$ and $r$ increasing) with initial cell density similar to [Jin16a]. The implications of this observation are considered in the Discussion Section. To compare the three models quantitatively, the generalized least squares (GLS) errors were computed for each model and data set and reported in Table 4.2.

**Table 4.2** Generalized least squares (GLS) errors.

| Model | Initial cell density | | | | | |
|---|---|---|---|---|---|---|
| | 10,000 | 12,000 | 14,000 | 16,000 | 18,000 | 20,000 |
| classical FKPP | 786.80 | 557.28 | 616.76 | 619.12 | 685.17 | 964.19 |
| Porous-FKPP | 681.18 | 540.29 | 418.57 | 566.89 | 744.44 | 928.38 |
| BINN-guided model | **557.01** | **317.18** | **410.79** | **393.15** | **307.74** | **386.52** |

Table of GLS errors between the model solutions and scratch assay data. Each column corresponds to an experiment with different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Bold numbers represent the minimum GLS error across the three models.

The results in Table 4.2 showed that Equation (4.17) with Equations (4.18a)-(4.18c) fit each data set more accurately than the classical FKPP or Generalized Porous-FKPP models. This behavior is not surprising given that the BINN-guided model is more complex. Therefore, model selection methods, which balance model accuracy with model complexity, were also used to compare the *quality* of each model relative to the others. In particular, the modified Akaike Information Criterion (AIC) from [Ban17] was used to account for the statistical error model in Equation (4.4). See Table 4.3 for the AIC scores across each model and data set.

**Table 4.3** Akaike Information Criterion (AIC) scores.

| Model | Initial cell density | | | | | |
|---|---|---|---|---|---|---|
| | 10,000 | 12,000 | 14,000 | 16,000 | 18,000 | 20,000 |
| classical FKPP | 1239.6 | 1175.8 | 1194.5 | 1195.2 | 1214.0 | 1277.2 |
| Porous-FKPP | 1214.9 | 1172.0 | **1124.8** | 1180.9 | 1231.3 | 1272.2 |
| BINN-guided model | **1183.7** | **1079.5** | 1127.3 | **1119.2** | **1073.9** | **1116.1** |

Table of AIC scores for each model and scratch assay data set. Each column corresponds to an experiment with different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Bold numbers represent the minimum AIC score across the three models.

The results in Table 4.3 showed that the BINN-guided delay-reaction-diffusion model outperforms the classical FKPP and Generalized Porous-FKPP models across all data sets except with initial cell density 14,000 cells per well. This discrepancy follows from Table (4.6) where the additional parameters in Equation (4.17) only slightly decreased the GLS error for the data set with initial density of 14,000. Finally, to quantify the "value" of adding the novel delay term in Equation (4.18c) the differences between AIC scores for each model and the minimum AIC score, denoted by $\Delta$AIC, are shown in Table 4.4.

Table 4.4 suggests that the delay term is most impactful for data sets with large initial density (i.e. 18,000 and 20,000 cells per well) since the $\Delta$AIC scores are significantly larger for these data sets. Biological analysis and explanations for these results are considered in the following Discussion Section.

**Table 4.4** Difference Akaike Information Criterion (ΔAIC) scores.

| Model | Initial cell density | | | | | |
|---|---|---|---|---|---|---|
| | 10,000 | 12,000 | 14,000 | 16,000 | 18,000 | 20,000 |
| classical FKPP | 55.90 | 96.26 | 69.71 | 76.01 | 140.08 | 161.11 |
| Porous-FKPP | 31.23 | 92.54 | 0.00 | 61.70 | 157.42 | 156.11 |
| BINN-guided model | 0.00 | 0.00 | 2.53 | 0.00 | 0.00 | 0.00 |

Table of AIC differences (ΔAIC) between each model and scratch assay data set. Each column corresponds to an experiment with different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Each ΔAIC score represents the difference between a model's AIC score and the minimum recorded AIC score for that data set.

## 4.5 Summary of contributions

We introduced biologically-informed neural networks (BINNs), an extension of physics-informed neural networks, to be feasible for applications with poorly understood dynamics (e.g., only knowledge of basic conservation laws) and sparse noisy observations. In this framework, neural networks are used to (i) approximate the solution of an unknown PDE and (ii) learn the nonlinear terms comprising the governing system without the need to specify the mechanistic form of each term. Applying our methodology to experimental cell migration data of PC-3 prostate cancer cells, to the best of our knowledge, we are the first to successfully use equation learning methods to discover a previously unconsidered real-world biological mechanism describing delayed population response. Using the discovered dynamics, we posited a new mechanistic equation that fit the experimental scratch assay data more accurately and had lower AIC scores compared to state-of-the-art models. Based on the success of this work, we believe BINNs can enable the rapid development and validation of new mathematical models for a broad range of real-world applications throughout the natural sciences.

## 4.6 Discussion

In this work, biologically-informed neural networks (BINNs) were introduced as a flexible and robust equation learning method for real-world biological applications. The BINNs framework was demonstrated using experimental biological data from scratch assays [Jin16a] and used to discover a delay term that had not yet been considered in the modeling of these data. The trained diffusivity, growth, and delay networks were used to guide the selection of the mechanistic model in Equation (4.17) with Equations (4.18a)-

(4.18c), which was shown to model the data more accurately than the current state-of-the-art models (i.e. classical FKPP and Generalized Porous-FKPP equations). The results shown in this work suggest that the BINNs framework can be successfully applied to a wide range of biological and physical problems where the data are sparse and the governing dynamics are unknown. The biological motivations for various aspects of the BINNs framework and significance of the results are discussed in the following paragraphs.

The model solutions in Figure 4.3 and Figure 4.4 indicated that using only density-dependent diffusivity and growth functions $\mathscr{D}(u)$ and $\mathscr{G}(u)$ was not sufficient to fully capture the scratch assay dynamics. Figure 4.4 highlighted this discrepancy at the second time measurement ($t = 0.5$ days) in which the model failed to capture the areas of both high and low cell density, despite letting $\mathscr{D}$ and $\mathscr{G}$ be universal function-approximating neural networks. In particular, the model solutions in the areas of high cell density (i.e. $x \in [0.0, 0.5]$ and $x \in [1.5, 2.0]$) showed exponential convergence to the carrying capacity, which successfully captured the data for later time points ($t \geq 1$ days) but over-predicted for early time points ($t < 1$ days). Similarly, diffusion in areas of low cell density (i.e. $x \in [0.5, 1.5]$) over-predicted the cell density profile for early time points but then matched the data accurately for later time points. From a mathematical perspective, this motivates the existence of a time delay that scales the density-dependent dynamics to be reduced for early time points and larger for later time points. There are also several biological motivations for considering a time delay. For example, [Dyd20] showed how cells are damaged at the borders of the scratch as a result of the experimental scratch assay protocol. Cell damage can potentially inhibit the communication between cells and physically block healthy cells from diffusing into uncolonized spatial regions. Another source of delay may stem from changes in density-dependent cell functions (e.g. differentiation, division, and senescence). Studies have shown that cells are more likely to terminally differentiate when cell populations approach carrying capacity [Pou95; Neu20]. Therefore, scratch assay experiments that are performed for high density populations may contain fewer mobile/proliferative cells at the borders of the scratch, thus causing a time delay in the cell migration dynamics.

A general framework for incorporating the delay term may be to consider diffusivity and growth functions $\mathscr{D} = D_{\mathrm{MLP}}(u, t)$ and $\mathscr{G} = G_{\mathrm{MLP}}(u, t)$, respectively. However, since the dynamics of diffusion and growth are assumed to be consistent throughout time, the diffusion and growth terms were chosen to be separable functions composed of diffusivity $\mathscr{D}(u)$, growth $\mathscr{G}(u)$, and delay $\mathscr{T}(t)$. Additionally, it was assumed that both diffusion and growth were scaled by the same time delay $\mathscr{T}(t)$ as opposed to a diffusion delay $\mathscr{T}_{\mathrm{D}}(t)$ and growth delay $\mathscr{T}_{\mathrm{G}}(t)$. This assumption may not be accurate if the time delay is a result of density-dependent changes in cell function where cells become mobile and proliferative at

different rates. In particular, since migration and proliferation have very different timescales, it might be natural to expect that the delays would also have different timescales. However, since the numerical solutions using $\mathcal{T}(t)$ matched the data sufficiently accurately, this question is left for future work. Finally, $\mathcal{T}(t)$ was constrained to output values between 0 and 1 and forced to be increasing with time. These constraints were chosen to ensure that the delay term modeled the time-dependent changes in cell dynamics for early time points but converged to unity by later time points.

In this work, BINNs revealed that the reaction-diffusion system in Equation (4.1) with cell density-dependent diffusivity and growth functions was insufficient to capture the data dynamics. However, the model discrepancy for data sets with large initial cell density motivated the development of a time delay which significantly improved the model accuracy and resolved the observed discrepancy. The diffusivity, growth, and delay networks were used to posit a mechanistic model (i.e. Equation (4.17) with Equations (4.18a)-(4.18c)). Using the logistic growth model (Equation (4.18b)) for the growth function and logistic regression (Equation (4.18c)) for the delay function followed straightforwardly from the parameter network solutions in Figure 4.7, however, the diffusivity function in Equation (4.18a) warrants further discussion.

Opinions vary between the biological validity of (i) the classical FKPP and (ii) the Generalized Porous-FKPP diffusivity functions. For example, one study compared (i) and (ii) using experimental wound size data and found that (ii) with $m = 4$ provided the best fit to the data [She90b]. Another study fit (i) and (ii) to experimental cell migration data with different cell populations and found that one population was best described by constant diffusivity in (i) and the other by nonlinear diffusivity with $m = 1$ in (ii) [Sen07]. These studies do not reveal which approach is best, but they demonstrate that care is warranted. The diffusivity network ($D_{\mathrm{MLP}}$) solutions showed significant variability across the scratch assay data sets, so the posited mechanistic model was chosen to respect the observed variability while also being as simple as possible. Therefore, Equation (4.18a) was chosen to be a combination of the diffusivity functions in (i) and (ii). This way both (i) and (ii) can be seen as nested models of the posited diffusivity function by setting either $D = 0$ or $D_0 = 0$. Yet the posited diffusivity is still simple, as it only increases the number of parameters with respect to (ii) by one. It may be the case that the true diffusivity function is even more complex, such as a linear combination of powers:

$$\mathcal{D} = D_0 + D_1 \left(\frac{u}{K}\right)^{m_1} + D_2 \left(\frac{u}{K}\right)^{m_2},$$

with baseline diffusivity $D_0$, diffusion rates $D_1$ and $D_2$, carrying capacity $K$, and exponents

$m_1$ and $m_2$. However, these considerations are beyond the scope of the present work and left for future work.

The parameters of (i) the classical FKPP in Equation (4.19), (ii) the Generalized Porous-FKPP in Equation (4.20), and (iii) the BINN-guided model in Equation (4.17) with Equations (4.18a)-(4.18c) were optimized numerically for each scratch assay data set. The optimized parameters for (i) in Table C.1 all fall within the ranges reported in [Jin16a]. However, this is not the case for any set of parameter values for (ii) as shown in Table C.2. This is likely due to the parameter optimization being conducted using the adjusted statistical error model in Equation (4.4) with $\gamma = 0.2$ and since the exponent $m$ in the Porous-FKPP diffusivity function was not fixed at $m = 1$ as in [Jin16a]. However, in both (i) and (ii), the diffusion coefficient, $D$, and intrinsic growth rate, $r$, showed variability with initial cell density, similar to the conclusions drawn in [Jin16a]. Therefore, in theory, if the delay term in Equation (4.18c) accounts for the time it takes for density-dependent growth and diffusion to become active in the system, which may be a function of initial cell density, then the variability among diffusion coefficients and intrinsic growth rates for the BINN-guided delay-reaction-diffusion model should be reduced across the scratch assay experiments. However, from the optimized parameter values in Table 4.1, the baseline diffusion rate $D_0$ and intrinsic growth rate $r$ generally increase with initial cell density and the diffusion coefficient $D$ generally decreases with initial cell density. This observation may indicate (i) practical identifiability issues between the diffusion, growth, and delay terms or (ii) the existence additional mechanisms that are not accounted for in the model. To confirm this, a Bayesian parameter estimation framework can be used to examine practical identifiability of parameters [Lag18; Ado15a]. Then, a possible strategy to mitigate this issue would be to optimize the parameters of Equation (4.17) with Equations (4.18a) and (4.18b) jointly across each scratch assay data set while allowing the delay parameters in Equation (4.18c) to be tuned separately for each set. This exploration is left for future work.

The BINN-guided delay-reaction-diffusion model was compared to the baseline classical FKPP and Generalized Porous-FKPP models using both GLS errors and modified AIC scores. The GLS errors in Table 4.2 showed that the BINN-guided model fits the data more accurately than the baseline models across each scratch assay data set. However, this improvement in accuracy is due to the increased model complexity (i.e. number of parameters and PDE terms) in the BINN-guided model. Therefore, to rank the quality of each model, AIC scores were also computed since they balance model accuracy with model complexity. The AIC scores reported in Table 4.3 indicate that the BINN-guided model also exceeds the baseline models in terms of relative quality across each scratch assay data set except with initial cell density 14,000 cells per well, in which the Generalized Porous-FKPP

model has a slightly smaller AIC score. In other words, Tables 4.2 and 4.3 indicate that the BINN-guided model performs as well or better than the state-of-the-art in modeling the suite of scratch assay experiments from [Jin16a]. In particular, this advantage is afforded by including the delay term in Equation (4.18c). To quantify the relative value of adding the delay term, the AIC scores from Table 4.3 are used to compute difference AIC ($\Delta$AIC) scores in Table 4.4 in which the $\Delta$AIC score for a fixed model and data set is given by the difference between the corresponding AIC score and the minimum AIC score across all models for the given data set. The $\Delta$AIC scores in Table 4.4 indicate that the relative value of the delay term is largest for data sets with initial cell density 18,000 and 20,000 cells per well. This observation is supported by the relevant biology discussed at the beginning of this section, in which large initial cell densities either (i) result in more damaged cells near the borders of the scratch, (ii) cause more cells in the population to have terminally differentiated away from mobile/proliferative cell functions, or (iii) some combination of (i) and (ii) and other potentially unconsidered biological sources, all of which increase the potential time delay before the density-dependent diffusion and growth dynamics become the primary drivers of the temporal evolution of the system.

### 4.6.1 Conclusions

BINNs, a robust and flexible framework for equation learning with sparse and noisy data, was demonstrated and used to posit a mechanistic equation that outperforms the state-of-the-art in modeling experimental scratch assay data. The development, training, and evaluation of BINNs and the resulting model selection and analysis were reported to justify these claims. The discovered time delay term may have important implications for the reproducibility and modeling of scratch assays, since the majority of the reported data fall within the time delay region. Some of the drawbacks of the BINNs method and opportunities for future work and development are discussed below.

Since BINNs rely on multilayer perceptrons (MLPs), the learned dynamics may not generalize well outside the training domain. For example, in the present work, if the observed cell densities for a particular experiment do not approach the carrying capacity (e.g. the scratch assay data set with 10,000 initial cells per well) then the learned dynamics given by $D_{\mathrm{MLP}}$ and $G_{\mathrm{MLP}}$ may lead to biologically unrealistic behavior (see $G_{\mathrm{MLP}}$ solutions in Figure 4.7). Further, since none of the scratch assay data reported values that significantly exceeded the empirically set carrying capacity, $G_{\mathrm{MLP}}$ would likely not generalize well to a scenario with exceedingly large observed cell densities. Options for mitigating this issue include (i) replacing unrealistic MLP terms with mechanistic models (e.g. logistic growth

instead of $G_{\mathrm{MLP}}$) if the particular dynamics are known *a priori,* or (ii) adding additional constraints which force the MLP terms to satisfy specific values (e.g. $G_{\mathrm{MLP}}(u = K) = 0$).

An opportunity for future development is quantifying the uncertainty of both the approximate solution, $u_{\mathrm{MLP}}$, and the parameter networks, $D_{\mathrm{MLP}}$, $G_{\mathrm{MLP}}$, and $T_{\mathrm{MLP}}$. From the frequentist perspective, so called "subagging" (i.e. subsample aggregating) can be used to build posterior distributions of the model solutions and parameter networks [Buh12]. In this framework, one simply samples $N$ training/validation splits and trains a BINN for each split. Then kernel density estimation or some other equivalent methodology can be used to build distributions from the $N$ number of trained BINNs. Alternatively, from the Bayesian perspective, physics-informed neural networks were recently extended to Bayesian physics-informed neural networks (B-PINNs) [Yan20]. In this framework, Bayesian neural networks are substituted for $u_{\mathrm{MLP}}$ and regularized using a pre-specified governing PDE. In the BINNs framework, Bayesian neural networks could also be substituted for $D_{\mathrm{MLP}}$, $G_{\mathrm{MLP}}$, and $T_{\mathrm{MLP}}$ to quantify the uncertainty of the PDE terms in addition to the model solution.

While BINNs were demonstrated using one-dimensional reaction-diffusion PDEs for scratch assay data in this work, they can be applied on a wide spectrum of physical and biological problems (for both ODE and PDE systems) in which the governing dynamics are unknown and highly nonlinear. A straightforward next step for this work would be to evaluate BINNs on the two-dimensional scratch assay image data that were used to construct the one-dimensional cell density profiles in [Jin16a]. Further, more complicated cell dynamics could be incorporated into the governing system in the present work by including PDE terms that describe cell population heterogeneity or additional biological mechanisms for damaged (but not dead) cells at the borders of the scratch.

BINNs were used to address a canonical problem in the field of collective cell migration by analyzing how the combination of density-dependent cell motility and proliferation drive the temporal dynamics of cell invasion during an experimental scratch assay. This novel framework revealed new mechanistic and biological insights into this process by guiding the derivation of a mathematical model that has not been considered previously using traditional mathematical modeling approaches. The classical FKPP and Generalized Porous-FKPP models are ubiquitous in modeling cell migration and proliferation, yet the BINNs methodology presented here revealed that these models may fail to incorporate all of the relevant mechanisms underlying this process. These results suggest that new models incorporating a time delay may be necessary to accurately capture the dynamics within the first day of a scratch assay, i.e., just after the scratch is introduced. Based on the success of this work, BINNs establish a new paradigm for data-driven equation learning from sparse and noisy data that could enable the rapid development and validation of mathematical

models for a broad range of real-world applications throughout biology including ecology, epidemiology, and cell biology.

CHAPTER

5

# CONCLUSIONS

Throughout this work we presented a number of data-driven methods for augmenting and discovering mathematical models using observation data. There are various amounts of data and domain knowledge spanning the field of mathematical modeling, we considered the following three scenarios: (i) full knowledge of the governing system (i.e., a specified mathematical model) with sparse observation data, (ii) no knowledge of the governing system (i.e., unspecified mathematical model) with large amounts of data, and (iii) basic knowledge of the governing system (e.g., conservation laws) with sparse observation data. In the preceding chapters we considered ways of incorporating techniques from data science (e.g., state space reconstruction) and machine learning (e.g., neural networks) to address the mathematical challenges associated with each case. Our proposed methods are widely applicable to real-world scientific problems across the natural sciences, however, biological applications were used specifically as a test bed to evaluate these methods because of their notorious difficulty. The mathematical challenges and our proposed solutions in each of the cases above are discussed in the following section.

## 5.1 Contributions

A key challenge in case (i), and in parameter estimation generally, is that a set of available observation data may be insufficient to accurately parameterize a set of governing equations. This is true particularly in cases where one of more of the modeled variables contain correlated parameters and issues related to practical identifiability. We addressed this challenge by proposing a hybrid method that utilizes both mechanistic and non-mechanistic modeling strategies to systematically eliminate problematic system variables. The equations associated with unmodeled variables are replaced with the raw data for those variables. The remaining equations for the modeled variables are then calibrated to fit the data, in which the unmodeled variables are supplemented with the corresponding data. After estimating the parameters of the subset of equations, in order to make a prediction, state space reconstruction (SSR) is used to obtain model-free predictions of the unmodeled variables. Then, using the SSR predictions to supplement the unmodeled variables together with the calibrated equations for the modeled variables, joint predictions can be made to forecast future time series dynamics.

This work resulted in multiple mathematical and scientific contributions. First, by simply replacing subsets of model equations with data, the dimensionality of the inverse problem is decreased since fewer parameters need to be estimated. This translates to increased prediction accuracy of the hybrid method compared to the full model and model-free methods independently. Further, when utilizing Bayesian inference for modeled variables and SSR for unmodeled ones, we proposed a novel method for quantifying the uncertainty of the model predictions. We obtain 95% credible and prediction intervals by jointly sampling from the posterior parameter distributions of the modeled variables and the nearest neighbor trajectories in delay-coordinate space of the unmodeled variables. Using correlation plots and rank deficiency of the Fisher information matrix (FIM), we confirmed that hybrid models generally reduce the practical unidentifiability of the parameters for the set of equations describing flour beetle population dynamics. In particular, our analysis concluded that the model for the Larva population contained correlated parameters which hindered the ability to obtain accurate parameter estimates. By replacing the corresponding equation with data/SSR predictions, the parameter correlations in the hybrid model were eliminated and the resulting time series predictions were more accurate. In summary, We proposed a hybrid methodology that can minimize the dimensionality of parameter inference while maximizing predictive accuracy of the resulting model. This methodology is broadly applicable to multivariate dynamical systems, particularly when developing predictive models for biological applications where data are typically sparse. We hypothe-

size that this methodology can enable mechanistic modeling to a wider range of problems across the natural sciences in which limited amounts of data and high dimensionality of the system inhibit accurate parameterization with traditional techniques.

Where case (i) dealt with full system knowledge and sparse data, case (ii) considers the opposite, in which no system dynamics are known but lots of data are available. The key mathematical challenge in this scenario is how to utilize the informative set of observation data to discover the dynamics that drive the temporal evolution of the system. A popular equation learning method involves the use of data denoising and partial derivative approximation to formulate the equation learning task as a sparse regression problem. Traditional approaches relied on the use of finite differences and polynomial splines for data denoising, which are prone to inaccuracy in the presence of large amounts of noise. We addressed this challenge by developing an equation learning pipeline based on a combination of neural networks and sparse regression. First, we denoise the observation data using an artificial neural network (ANN) as a surrogate model for the solution of the unknown governing differential equation. Then, using automatic differentiation of the ANN outputs to obtain partial derivative estimates, we construct a library of candidate terms thought to comprise the unknown system. Discovering the governing PDE is formulated as a sparse regression problem for the linear system given by the candidate library, in which the nonzero elements of the solution vector correspond to the small number terms that most accurately describe the data. A final parameter estimation procedure on the discovered equation is then used to fine-tune the model parameters.

The use of ANNs for data denoising and partial derivative approximation resulted in multiple scientific contributions to the field of data-driven discovery and equation learning. First, we showed that ANNs can be used to identify appropriate statistical error models (e.g., the proportionality constant in a generalized least squares problem) describing the error process in a set of observation data. Next, using ANNs together with the calibrated error model to denoise and approximate the partial derivatives of the data, we showed that under biologically realistic forms and levels of noise, ANNs can approximate derivatives up to four orders of magnitude more accurately compared to polynomial splines, and up to six orders of magnitude compared to finite differences for the advection-diffusion and Fisher-KPP models considered in this work. We note that polynomial interpolation was identified as the most robust denoising method by [Rud17b] after comparing against alternative methods such as Gaussian kernel smoothing, Tikhonov differentiation, and spectral differentiation with high-frequency term thresholding. Using the denoised data values and partial derivative estimates from the ANN to construct the library of candidate terms, we demonstrated that accurate PDE models could be discovered with noise levels

up to 25%. Further, we introduced modifications to the PDE-FIND algorithm, including training/validation data partitioning and a pruning algorithm, that significantly increased the robustness of the algorithm. The novel use of ANNs method presented in this work is a major step toward enabling equation learning in realistic scenarios with noisy observation data.

Finally, in case (iii), we consider the intersection of cases (i) and (ii), in which only basic information about the dynamics is known and the available data are sparse and noisy. In this scenario, due to poor understanding of the governing system, one does not have a set of equations with which to proceed using the hybrid modeling strategy. Similarly, without a large data set, unconstrained function-approximating ANNs will easily overfit the data, leading to inaccuracies in the denoised data values and partial derivative approximations. Therefore, to address the mathematical challenge of equation learning with sparse data, we draw from cases (i) and (ii) in order to maximize the benefits of each strategy in the current setting. Namely, we develop biologically-informed neural networks (BINNs), a neural network-based equation learning methodology that learns to approximate a set of observation data while simultaneously learning the governing dynamics. To prevent the network from overfitting, we introduce a regularization term that uses automatic differentiation to force the network outputs to satisfy a governing PDE. However, since only basic information about the dynamics is known (e.g., diffusion, reaction, etc.), the mechanistic models comprising these terms are replaced with additional neural networks that learn the appropriate dynamics. Thus, by training BINNs to minimize the distance to the data while satisfying the general form of a governing PDE, the nonlinear terms of the PDE are also learned. In this work we use simulated scratch assay data from the classical Fisher-KPP and Generalized Porous Fisher-KPP models as a proof of concept, and then apply BINNs to real-world experimental scratch assay data describing the dynamics of cell invasion for a population of PC-3 prostate cancer cells. The discovered dynamics from the trained BINNs are used to posit a new mechanistic model where traditional mathematical modeling techniques are used to conduct a model comparison study.

Our work on BINNs has resulted in multiple contributions to mathematics, machine learning, and biology. First and foremost, BINNs are a novel methodology that extend standard theory-informed neural networks to the class of equation learning algorithms, particularly in scenarios with poorly understood dynamics and sparse noisy data. For the first time, to the best of our knowledge, we used neural networks to identify an appropriate statistical error model for real-world experimental scratch assay data. Further, using our methodology, we are the first to use neural network-based equation learning methods to discover a previously unconsidered biological mechanism describing delayed population

response. This discovery may have important implications in future scratch assay studies, since most experiments record and model data in the early-time scratch assay dynamics, which are most affected by our discovered delay term. In this work we also used BINNs to posit a new state-of-the-art mathematical model for modeling scratch assay experiments. By conducting a model comparison study using model error and Akaike Information Criteria (AIC), we demonstrated that our posited model was more accurate and produced better AIC scores compared to previous state-of-the-art models for the scratch assay experiments considered in this work. Based on the success of this work, BINNs may enable the rapid development and validation of new mathematical models for a broad range of real-world applications throughout the natural sciences.

## 5.2   Future work

Overall, this work is aimed at the larger goal of enabling and accelerating the validation of new mathematical models from real-world data. Since BINNs can be seen as an extension/combination of our previous works, i.e., ANN-based equation learning (i.e., using neural networks for discovering new dynamics) and hybrid modeling (i.e., replacing model terms with data-driven methods), the future work of this dissertation is discussed in the context of BINNs. In particular, below, we discuss quantifying the uncertainty of BINNs, how BINNs can be extended to address practical unidentifiability and overfitting to individual data sets, and finally some plans for developing a fully machine learning pipeline connecting experimental data in the form of vector, image, or unstructured point cloud/graph data to an interpretable and generalizable mechanistic equation that describes and quantifies the uncertainty of the dynamics that drive the data.

In our work on BINNs, we used neural networks to posit a new mechanistic equation for describing scratch assay dynamics of cell populations, however, we did not quantify the uncertainty of the parameters or model solutions. A reasonable and straightforward step to quantifying the uncertainty is to repeat the parameter estimation procedure of the posited PDE model in a Bayesian context. Similar to the hybrid model, Bayesian inference using a delayed rejection adaptive metropolis (DRAM) algorithm can be performed to obtain posterior distributions for the model parameters and to compute 95% credible and prediction intervals for the solution of the posited PDE. However, recent work has also been done to extend physics-informed neural networks (PINNs) to a Bayesian framework using Bayesian neural networks (BNNs) [Yan20]. In this approach, a neural network surrogate for the solution of the dynamical system is replaced with a BNN that acts as a prior for a Hamiltonian Monte Carlo or variational inference framework for estimating posterior

distributions. However, unlike PINNs, which require the specification of a mechanistic PDE, the BINNs framework, which uses neural networks for the PDE terms, could also be replaced with BNNs. Therefore, BINNs would utilize BNNs to quantify the uncertainty of both the model solution as well as the nonlinear terms comprising the unknown PDE. Further, this framework may be even more robust to large noise levels due to the capability of BNNs to avoid overfitting.

Another opportunity for future work involves extending the BINNs framework to incorporate multiple data sets (e.g., experimental replicates) in parallel to address practical unidentifiability of parameters or to learn global physical/biological dynamics that capture broad system behavior. For example, it is well known that *how* data is collected can significantly impact the ability accurately parameterize mathematical models [Ban14c]. In [Nar20], we demonstrated that accurately parameterizing various characteristics of cell populations depends crucially on the time scale at which the data are collected. Therefore, by combining multiple data sets that each describe similar quantities of interest (e.g., cell density) with varying time scales in a parallel equation learning framework, one can conceivably better parameterize mathematical models at the population level. Further, combining multiple data sets in a BINNs-inspired equation learning framework can lead to the discovery of broad scale governing dynamics that encompass multiple data sets rather than capturing more narrow, experiment-specific properties. An example application of this framework is in modeling tumor growth in humans. It is uncommon for patients to receive more than 1-2 tumor scans, however, there is a large number of patients in various stages of cancer. By developing an equation learning framework that can utilize tumor data from a population of patients, one can discover transferable dynamics at the human population level, and can be further augmented and fine-tuned with patient-specific characteristics.

More formally, for a set of $M$ data sets, each describing quantity of interest $u_i$, $i = 1, \ldots, M$, that are functions of some common spatial domain, $\Omega$ (e.g. Cartesian coordinates $x, y, z$, etc.), and time, $t$, a neural network surrogate model can be trained to predict all quantities of interest from the shared input domain. Then, replacing the terms comprising the conservation laws of a governing PDE, the network outputs can be used together with data-set-specific properties (e.g. initial degree of confluence in scratch assays, patient characteristics in cancer modeling, etc.) as inputs to a set of parameter networks, $T_1, \ldots, T_L$, that learn the governing dynamics. By forcing the predicted solutions to each data set to satisfy a common set of governing PDEs, the global dynamics across all data sets can be learned. An example illustration of this framework is shown in Figure 5.1.

Finally, to allow BINNs to be used more effectively in the natural sciences, a set of methods must be developed to accurately process and quantify incoming experimental data. Ex-

**Figure 5.1** Multi-dataset extension to BINNs. A neural network surrogate for the solution of the dynamical system (left) inputs spatiotemporal vectors, $\begin{bmatrix} x_1 & \cdots & x_N & t \end{bmatrix}$, and predicts the corresponding outputs for $M$ different data sets. Then, parameter networks, $T_1, \ldots, T_L$ (center), input the quantities of interest, $u_1, \ldots, u_M$, with the option to include additional conditions, $c_1, \ldots, c_M$, pertaining to individual data sets (e.g. initial degree of confluence in scratch assays, patient characteristics in cancer modeling, etc.) to learn the associated dynamics. Finally, the quantities of interest, $u_1, \ldots, u_M$, and the predicted terms of the governing system, $T_1, \ldots, T_L$, can be combined into a single governing PDE (right) that describes the global dynamics that are common across all data sets.

perimental data typically comes in the form of images or unstructured point clouds/graphs. For example, in the case of scratch assay and cancer patient data, the experimental data typically comes in the form of images. However, quantifying image data is a laborious and time consuming task, often requiring expert supervision and verification. Machine learning has recently emerged as a powerful tool for the fast and accurate processing of structured and unstructured data, however, it can require large amounts of annotated training data, which can be difficult to obtain in practice. Further, traditional machine learning-based methods do not necessarily preserve physically realistic object properties (e.g., contiguity) which can hinder the analysis and equation discovery of the quantified data. Thus, a key challenge in this area is the development of data-efficient learning algorithms that alleviate the burden of manual image processing while maintaining physically meaningful object properties. Though not discussed in this dissertation, we have developed methods based on convolutional neural networks for semantic segmentation that mimics the human task of object tracing [Rut18; Rut19], and an extension of the tracing algorithm into higher dimensions [Lag20c] which were shown to be more accurate than state-of-the-art models, orders of magnitude faster than manual segmentation, and ensured that physically realistic regions are generated for each segmented object. Further, we have developed novel deep

learning methods for the unsupervised processing of unstructured point cloud data [Pan20]. Thus, an opportunity exists for combining the above methods that automatically extract quantifiable data from structured and unstructured sources together with the BINNs framework for a fully end-to-end equation learning methodology. This framework would be a significant step toward the overall goals of this dissertation by enabling the wide-spread use of novel mathematical models in conjunction with real-world data.

# BIBLIOGRAPHY

[Ado15a]   Adoteye, K. et al. "Correlation of parameter estimators for models admitting multiple parameterizations". *International Journal of Pure and Applied Mathematics* **105**.3 (2015). Publisher: Academic Publications, Ltd., pp. 497–522.

[Ado15b]   Adoteye, K. et al. "Optimal design of non-equilibrium experiments for genetic network interrogation". *Applied Mathematics Letters* **40** (2015), pp. 84–89.

[And99]    Anders, U. & Korn, O. "Model selection in neural networks". *Neural Networks* **12**.2 (1999), pp. 309–323.

[And74]    Anderssen, R. S. & Bloomfield, P. "Numerical differentiation procedures for non-exact data". en. *Numerische Mathematik* **22**.3 (1974), pp. 157–182.

[Ang09]    Anguige, K. & Schmeiser, C. "A one-dimensional model of cell diffusion and aggregation, incorporating volume filling and cell-to-cell adhesion". en. *J. Math. Biol.* **58**.3 (2009), p. 395.

[Aok17]    Aoki, K. et al. "Propagating Wave of ERK Activation Orients Collective Cell Migration". *Developmental Cell* **43**.3 (2017), 305–317.e5.

[Arc11]    Arciero, J. C. et al. "Continuum Model of Collective Cell Migration in Wound Healing and Colony Expansion". en. *Biophysical Journal* **100**.3 (2011), pp. 535–543.

[Bal14]    Baldock, A. L. et al. "Patient-Specific Metrics of Invasiveness Reveal Significant Prognostic Benefit of Resection in a Predictable Subset of Gliomas". en. *PLOS ONE* **9**.10 (2014), e99057.

[Ban16]    Banks, H. et al. "Use of difference-based methods to explore statistical and mathematical model discrepancy in inverse problems". *Journal of Inverse and Ill-posed Problems* **24** (2016).

[Ban09a]   Banks, H. T. & Tran, H. T. *Mathematical and Experimental Modeling of Physical and Biological Processes*. 2009.

[Ban09b]   Banks, H. T. & Tran, H. T. *Mathematical and Experimental Modeling of Physical and Biological Processes*. en. Google-Books-ID: SSRapIe8p3QC. CRC Press, 2009.

[Ban11]    Banks, H. T. et al. "Estimation of Cell Proliferation Dynamics Using CFSE Data". en. *Bulletin of Mathematical Biology* **73**.1 (2011), pp. 116–150.

[Ban14a]   Banks, H. T. et al. *Modeling and Inverse Problems in the Presence of Uncertainty*. English. Boca Raton: Chapman and Hall/CRC, 2014.

[Ban15a]    Banks, H. T. et al. "Model comparison tests to determine data information content". *Applied Mathematics Letters* **43** (2015), pp. 10–18.

[Ban15b]    Banks, H. T. et al. "Uncertainty quantification in modeling HIV viral mechanics". eng. *Mathematical biosciences and engineering: MBE* **12**.5 (2015), pp. 937–964.

[Ban14b]    Banks, H. T. et al. *Modeling and inverse problems in the presence of uncertainty*. Chapman and Hall/CRC, 2014.

[Ban14c]    Banks, H. T. et al. *Modeling and inverse problems in the presence of uncertainty*. Chapman and Hall/CRC, 2014.

[Ban17]     Banks, H. & Joyner, M. L. "AIC under the framework of least squares estimation". *Applied Mathematics Letters* **74** (2017), pp. 33 –45.

[Bil13]     Billings, S. A. *Nonlinear System Identification: NARMAX Methods in the Time, Frequency, and Spatio-Temporal Domains*. en. John Wiley & Sons, 2013.

[Bin07]     Bindschadler, M. & McGrath, J. L. "Sheet migration by wounded monolayers as an emergent property of single-cell dynamics". en. *Journal of Cell Science* **120**.5 (2007), pp. 876–884.

[Bon18]     Boninsegna, L. et al. "Sparse learning of stochastic dynamical equations". *The Journal of Chemical Physics* **148**.24 (2018), p. 241723.

[Bot19]     Both, G.-J. et al. *DeepMoD: Deep learning for Model Discovery in noisy data*. 2019. arXiv: 1904.09406 [physics.comp-ph].

[Bru16a]    Brunton, S. L. et al. "Discovering governing equations from data by sparse identification of nonlinear dynamical systems". *Proceedings of the National Academy of Sciences* (2016), p. 201517384.

[Bru16b]    Brunton, S. L. et al. "Discovering governing equations from data by sparse identification of nonlinear dynamical systems". *Proceedings of the National Academy of Sciences* **113**.15 (2016), pp. 3932–3937. eprint: https://www.pnas.org/content/113/15/3932.full.pdf.

[Buh12]     Buhlmann, P. "Bagging, Boosting and Ensemble Methods". *Handbook of Computational Statistics* (2012).

[Bur95]     Burgess, A. N. "Non-linear model identification and statistical significance tests and their application to financial modelling". en. *Artificial Neural Networks, 1995., Fourth International Conference on*. 1995, pp. 312–317.

[Cas89]     Casdagli, M. "Nonlinear prediction of chaotic time series". *Physica D: Nonlinear Phenomena* **35**.3 (1989), pp. 335–356.

[Cha14]    Chapnick, D. A. & Liu, X. "Leader cell positioning drives wound-directed collective migration in TGF beta-stimulated epithelial sheets". en. *Mol. Biol. Cell* **25**.10 (2014), pp. 1586–1593.

[Che13]    Chen, Y. et al. "Robust Sparse Regression under Adversarial Corruption". en. *International Conference on Machine Learning*. 2013, pp. 774–782.

[Che06]    Cheng, H. et al. "Multistep-ahead time series prediction". *Lecture Notes in Computer Science: Advances in Knowledge Discovery and Data Mining* **3918**.765-774 (2006).

[CA09]     Cintrón-Arias, A. et al. "A sensitivity matrix based methodology for inverse problem formulation". *Journal of Inverse and Ill-posed Problems* **17**.6 (2009), pp. 545–564.

[Cob80]    Cobelli, C. & DiStefano, J. J. "Parameter and structural identifiability concepts and ambiguities: a critical review and analysis". en. *American Journal of Physiology - Regulatory, Integrative and Comparative Physiology* **239**.1 (1980), R7–R24.

[Cod08]    Codling, E. A. et al. "Random walk models in biology". *Journal of The Royal Society Interface* **5**.25 (2008), pp. 813–834.

[Con97]    Constantino, R. F. et al. "Chaotic dynamics in an insect population". *Science* **276** (1997), pp. 1881–1882.

[Dal95]    Dale, P. D. et al. "Travelling waves in wound healing". en. *FORMA* **10**.3 (1995), pp. 205–222.

[Dat17]    Datta, A. & Zou, H. "Cocolasso for high-dimensional error-in-variables regression". English (US). *Ann. Statist.* **45**.6 (2017), pp. 2400–2426.

[DG00]     De Gaetano, A. & Arino, O. "Mathematical modelling of the intravenous glucose tolerance test". *Journal of Mathematical Biology* **40**.2 (2000), pp. 136–168.

[Den95]    Dennis, B. et al. "Nonlinear Demographic Dynamics: Mathematical Models, Statistical Methods, and Biological Experiments". en. *Ecological Monographs* **65**.3 (1995), pp. 261–282.

[Die81]    Dierckx, P. "An Algorithm for Surface-Fitting with Spline Functions". en. *IMA J Numer Anal* **1**.3 (1981), pp. 267–283.

[DiS13]    DiStefano, J. *Dynamic Systems Biology Modeling and Simulation*. en. Academic Press, 2013.

[Dyd20]    Dydowiczova, A. et al. "Improved multiparametric scrape loading-dye transfer assay for a simultaneous high-throughput analysis of gap junctional intercellular communication, cell density and viability". *Sci Rep.* **10** (2020).

[Dys15]    Dyson, L. & Baker, R. E. "The importance of volume exclusion in modelling cellular migration". en. *J. Math. Biol.* **71**.3 (2015), pp. 691–711.

[Far87]    Farmer, J & Sidorowich, J. "Predicting chaotic time series". *Phys. Rev. Lett.* **59** (1987), pp. 845–848.

[Fra03]    Francis, C. R. I. C. et al. "Quantifying annual variation in catchability for commercial and research fishing". en. *Fishery Bulletin* **101**.2 (2003), pp. 293–304.

[Fri09]    Friedl, P. & Gilmour, D. "Collective cell migration in morphogenesis, regeneration and cancer". en. *Nat Rev Mol Cell Biol* **10**.7 (2009), pp. 445–457.

[Gal13]    Gallaher, J. & Anderson, A. R. A. "Evolution of intratumoral phenotypic heterogeneity: the role of trait inheritance". en. *Interface Focus* **3**.4 (2013), p. 20130016.

[Gei93]    Geisser, S. *Predictive Inference*. en. Google-Books-ID: wfdlBZ_iwZoC. CRC Press, 1993.

[Haa06]    Haario, H. et al. "DRAM: Efficient adaptive MCMC". *Statistics and Computing* **16**.4 (2006), pp. 339–354.

[Ham16]    Hamilton, F. et al. "Ensemble Kalman filtering without a model". *Physcial Review X* **6** (2016), p. 011021.

[Ham17]    Hamilton, F. et al. "Hybrid modeling and prediction of dynamical systems". *arXiv:1701.08141 [math]* (2017). arXiv: 1701.08141.

[Har17]    Haridas, P. et al. "Quantifying rates of cell migration and cell proliferation in co-culture barrier assays reveals how skin and melanoma cells interact during melanoma spreading and invasion". eng. *J. Theor. Biol.* **423** (2017), pp. 13–25.

[Har13]    Hartig, F. & Dormann, C. F. "Does model-free forecasting really outperform the true model?" en. *Proceedings of the National Academy of Sciences* **110**.42 (2013), E3975–E3975.

[Has05]    Hastings, A. et al. "The spatial spread of invasions: new developments in theory and evidence". en. *Ecology Letters* **8**.1 (2005), pp. 91–101.

[Hel95]    Helbing, D. & Molnár, P. "Social force model for pedestrian dynamics". *Phys. Rev. E* **51**.5 (1995). Publisher: American Physical Society, pp. 4282–4286.

[Hor91]    Hornik, K. "Approximation capabilities of multilayer feedforward networks". *Neural Networks* **4**.2 (1991), pp. 251 –257.

[Hsi05]    Hsieh, C.-H. et al. "Distinguishing random environmental fluctuations from ecological catastrophes for the North Pacific Ocean". *Nature* **435**.7040 (2005), pp. 336–340.

[Osf]       "https://osf.io/hx7bj/" ().

[III15]     III, J. D. *Dynamic Systems Biology Modeling and Simulation*. en. Google-Books-ID: nWoYAgAAQBAJ. Academic Press, 2015.

[Jim92]     Jimenez, J et al. "Forecasting on chaotic time series: A local optimal linear-reconstruction method". *Phys. Rev. A* **45**.6 (1992), p. 3553.

[Jin16a]    Jin, W. et al. "Reproducibility of scratch assays is affected by the initial degree of confluence: Experiments, modelling and model selection". *Journal of Theoretical Biology* **390** (2016), pp. 136 –145.

[Jin16b]    Jin, W. et al. "Reproducibility of scratch assays is affected by the initial degree of confluence: Experiments, modelling and model selection". en. *Journal of Theoretical Biology* **390** (2016), pp. 136–145.

[Joh12]     Johnston, S. T. et al. "Mean-field descriptions of collective migration with strong adhesion". *Phys. Rev. E* **85**.5 (2012), p. 051922.

[Joh14]     Johnston, S. T. et al. "How much information can be obtained from tracking the position of the leading edge in a scratch assay?" en. *Journal of The Royal Society Interface* **11**.97 (2014), p. 20140325.

[Joh15]     Johnston, S. T. et al. "Estimating cell diffusivity and cell proliferation rate by interpreting IncuCyte ZOOM assay data using the Fisher-Kolmogorov model". *BMC Systems Biology* **9**.38 (2015).

[Jon84]     Jones, C. K. R. T. "Stability of the travelling wave solution of the FitzHugh-Nagumo system". en. *Transactions of the American Mathematical Society* **286**.2 (1984), pp. 431–431.

[Kes17]     Keskar, N. S. et al. "On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima". *arXiv:1609.04836 [cs, math]* (2017). arXiv: 1609.04836.

[Kin17]     Kingma, D. P. & Ba, J. "Adam: A Method for Stochastic Optimization". *arXiv: 1412.6980* (2017). arXiv: 1412.6980.

[Kin14]     Kingma, D. P. & Ba, J. L. "Adam: A method for stochastic optimization". *Proc. 3rd Int. Conf. Learn. Representations*. 2014.

[Kug98]     Kugiumtzis, D et al. "Regularized local linear prediction of chaotic time series". *Physica D: Nonlinear Phenomena* **112**.3 (1998), pp. 344–360.

[Kur00]     Kurganov, A. & Tadmor, E. "New High-Resolution Central Schemes for Nonlinear Conservation Laws and Convection–Diffusion Equations". en. *Journal of Computational Physics* **160**.1 (2000), pp. 241–282.

[Lag18]     Lagergren, J. et al. "Forecasting and uncertainty quantification using a hybrid of mechanistic and non-mechanistic models for an age-structured population model." *Bulletin of Mathematical Biology* **80**.6 (2018).

[Lag20a]   Lagergren, J. et al. "Biologically-informed neural networks guide mechanistic modeling from sparse experimental data." *arXiv:2005.13073* (2020).

[Lag20b]   Lagergren, J. et al. "Learning partial differential equations for biological transport models from noisy spatiotemporal data." *Proceedings of the Royal Society A* **476**.2234 (2020).

[Lag20c]   Lagergren, J. et al. "Region growing with convolutional neural networks for biomedical image segmentation". *arXiv:2009.11717* (2020).

[Lag20d]   Lagergren, J. H. et al. "Learning partial differential equations for biological transport models from noisy spatio-temporal data". *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* **476**.2234 (2020), p. 20190800. eprint: `https://royalsocietypublishing.org/doi/pdf/10.1098/rspa.2019.0800`.

[Lai11]      Laine, M. "DRAM - Delayed Rejection Adaptive Metropolis." (2011).

[Loh12]     Loh, P.-L. & Wainwright, M. J. "High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity". EN. *Ann. Statist.* **40**.3 (2012), pp. 1637–1664.

[Mac95]   Mackay, D. J. C. "Probable networks and plausible predictions — a review of practical Bayesian methods for supervised neural networks". *Network: Computation in Neural Systems* **6**.3 (1995), pp. 469–505.

[Mai04]     Maini, P. K. et al. "Travelling Waves in a Wound Healing Assay". *Applied Mathematics Letters* **17** (2004), pp. 575–580.

[Man17]    Mangan, N. M. et al. "Model selection for dynamical systems via sparse regression and information criteria". *Proc. R. Soc. A* **473**.2204 (2017), p. 20170009.

[Mar16]    Martius, G. & Lampert, C. H. "Extrapolation and learning equations". *arXiv preprint arXiv:1610.02995* (2016).

[Mat04]    Matsubayashi, Y. et al. "ERK Activation Propagates in Epithelial Cell Sheets and Regulates Their Migration during Wound Healing". *Current Biology* **14**.8 (2004), pp. 731–735.

[McL15]    McLennan, R. et al. "Neural crest migration is driven by a few trailblazer cells with a unique molecular signature narrowly confined to the invasive front". en. *Development* **142**.11 (2015), pp. 2014–2025.

[Mes14]    Meshkat, N. et al. "On Finding and Using Identifiable Parameter Combinations in Nonlinear Dynamic Systems Biology Models and COMBOS: A Novel Web Implementation". *PLOS ONE* **9**.10 (2014), e110261.

[Mur02]    Murray, J. D. *Mathematical Biology I. An Introduction*. en. Ed. by Antman, S. S. et al. 3rd. Vol. 17. Interdisciplinary Applied Mathematics. New York, NY: Springer New York, 2002.

[Nar20]    Nardini, J. et al. "Learning equations from biological data with limited time samples." *Bulletin of Mathematical Biology* **82**.119 (2020).

[Nar19]    Nardini, J. T. & Bortz, D. M. "The influence of numerical error on parameter estimation and uncertainty quantification for advective PDE models". *Inverse Problems* **35**.6 (2019), p. 065003.

[Nar16]    Nardini, J. T. et al. "Modeling keratinocyte wound healing: cell-cell adhesions promote sustained migration". *Journal of Theoretical Biology* **400** (2016), pp. 103–117.

[Neu20]    Neurohr, G. E. & Amon, A. "Relevance and Regulation of Cell Density". *Trends in Cell Biology* **30**.3 (2020), pp. 213 –225.

[Nik06]    Nikolic, D. L. et al. "Role of boundary conditions in an experimental model of epithelial wound healing". en. *American Journal of Physiology - Cell Physiology* **291**.1 (2006), pp. C68–C75.

[Olu13]    Olufsen, M. S. & Ottesen, J. T. "A practical approach to parameter estimation applied to model predicting heart rate regulation". en. *Journal of Mathematical Biology* **67**.1 (2013), pp. 39–68.

[Pan20]    Paniagua, T. et al. "A simple deconvolutional mechanism for point clouds and sparse unordered data." *Proceedings of the AAAI Conference on Artificial Intelligence* **34**.10 (2020).

[Par00]    Parlos, A. G. et al. "Multi-step-ahead prediction using dynamic recurrent neural networks". *Neural Networks* **13**.7 (2000), pp. 765–786.

[Per13]    Perretti, C. et al. "Model-free forecasting outperforms the correct mechanistic model for simulated and experimental data". *Proceedings of the National Academy of Sciences* **110** (2013), pp. 5253–5257.

[Pou95]    Poumay, Y. & Pittelkow, M. R. "Cell Density and Culture Factors Regulate Keratinocyte Commitment to Differentiation and Expression of Suprabasal K1/K10 Keratins". *Journal of Investigative Dermatology* **104**.2 (1995), pp. 271 –276.

[Rai19]     Raissi, M. et al. "Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations". *Journal of Computational Physics* **378** (2019), pp. 686 –707.

[Rau09]     Raue, A. et al. "Structural and practical identifiability analysis of partially observed dynamical models by exploiting the profile likelihood". en. *Bioinformatics* **25**.15 (2009), pp. 1923–1929.

[Reg05]     Regonda, S et al. "Local polynomial method for ensemble forecast of time series". *Nonlin. Proc. in Geophys.* **12** (2005), pp. 397–406.

[Ris15]     Rish, I. & Grabarnik, G. Y. *Sparse Modeling: Theory, Algorithms, and Applications*. Boca Raton, FL: CRC Press, 2015.

[Roc15]     Rockne, R. C. et al. "A patient-specific computational model of hypoxia-modulated radiation resistance in glioblastoma using 18F-FMISO-PET". *Journal of The Royal Society Interface* **12**.103 (2015).

[Ros10]     Rosenbaum, M. & Tsybakov, A. B. "Sparse recovery under matrix uncertainty". EN. *Ann. Statist.* **38**.5 (2010), pp. 2620–2651.

[Rud17a]    Rudy, S. H. et al. "Data-driven discovery of partial differential equations". *Science Advances* **3**.4 (2017), e1602614.

[Rud17b]    Rudy, S. H. et al. "Data-driven discovery of partial differential equations". *Science Advances* **3**.4 (2017). eprint: `https://advances.sciencemag.org/content/3/4/e1602614.full.pdf`.

[Rut18]     Rutter, E. et al. "Automated object tracing for biomedical image segmentation using a deep convolutional neural network." *Lecture Notes in Computer Science* **11073** (2018).

[Rut19]     Rutter, E. et al. "A convolutional neural network method for boundary optimization enables few-shot learning for biomedical image segmentation." *Lecture Notes in Computer Science* **11795** (2019).

[Rut17]     Rutter, E. M. et al. "Mathematical Analysis of Glioma Growth in a Murine Model". *Scientific Reports* **7** (2017).

[Sah18]     Sahoo, S. et al. "Learning Equations for Extrapolation and Control". *Proceedings of the 35th International Conference on Machine Learning*. Ed. by Dy, J. & Krause, A. Vol. 80. Proceedings of Machine Learning Research. Stockholmsmässan, Stockholm Sweden: PMLR, 2018, pp. 4442–4450.

[Sau94]     Sauer, T. "Time series prediction by using delay coordinate embedding". *Time Series Prediction: Forecasting the Future and Understanding the Past*. Addison Wesley, 1994, pp. 175–193.

[Sch17]  Schaeffer, H. "Learning partial differential equations via data discovery and sparse optimization". *Proc. R. Soc. A* **473**.2197 (2017), p. 20160446.

[Sch06]  Schelter, B. et al. *Handbook of time series analysis: recent theoretical developments and applications*. John Wiley and Sons, 2006.

[Sch98]  Schroer, C. G. et al. "Predicting chaos most of the time from embeddings with self-intersections". *Phys. Rev. Lett.* **80**.7 (1998), p. 1410.

[Sen07]  Sengers, B. G. et al. "Experimental characterization and computational modelling of two-dimensional cell spreading for skeletal regeneration." *Journal of the Royal Society, Interface* **4**.17 (2007), pp. 1107–1117.

[She90a]  Sherratt, J. A. & Murray, J. D. "Models of epidermal wound healing". en. *Proc. R. Soc. Lond. B* **241**.1300 (1990), pp. 29–36.

[She90b]  Sherratt, J. A. et al. "Models of epidermal wound healing". *Proceedings of the Royal Society of London. Series B: Biological Sciences* **241**.1300 (1990), pp. 29–36. eprint: `https://royalsocietypublishing.org/doi/pdf/10.1098/rspb.1990.0061`.

[Sib99]  Sibert, J. R. et al. "An advection–diffusion–reaction model for the estimation of fish movement parameters from tagging data, with application to skipjack tuna (Katsuwonus pelamis)". *Canadian Journal of Fisheries and Aquatic Sciences* **56**.6 (1999), pp. 925–938.

[Smi92]  Smith, L. A. "Identification and prediction of low dimensional dynamics". *Physica D: Nonlinear Phenomena* **58**.1 (1992), pp. 50–76.

[Smi13]  Smith, R. C. *Uncertainty Quantification: Theory, Implementation, and Applications*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2013.

[Smi14]  Smith, R. C. *Uncertainty quantification: theory, implementation, and applications*. English. OCLC: 875327904. Philadelphia: SIAM, 2014.

[Ste15]  Stepien, T. L. et al. "A data-motivated density-dependent diffusion model of in vitro glioblastoma growth". en. *Mathematical Biosciences and Engineering* **12**.6 (2015), pp. 1157–1172.

[Sto48]  Stone, M. H. "The Generalized Weierstrass Approximation Theorem". en. *Mathematics Magazine* **21**.5 (1948), p. 237.

[Str06]  Strelioff, C. C. & Hübler, A. W. "Medium-term prediction of chaos". *Phys. Rev. Lett.* **96**.4 (2006), p. 044101.

[Sug94]     Sugihara, G. "Nonlinear forecasting for the classification of natural time series". *Philosophical Transactions of the Royal Society of London. Series A: Physical and Engineering Sciences* **348**.1688 (1994), pp. 477–495.

[Sug90]     Sugihara, G. & May, R. M. "Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series". *Nature* **344**.6268 (1990), pp. 734–741.

[Sug12]     Sugihara, G. et al. "Detecting Causality in Complex Ecosystems". en. *Science* **338**.6106 (2012), pp. 496–500.

[The09]     Theodoridis, S. & Koutroumbas, K. "Chapter 5 - Feature Selection". *Pattern Recognition (Fourth Edition)*. Ed. by Theodoridis, S. & Koutroumbas, K. Fourth Edition. Boston: Academic Press, 2009, pp. 261 –322.

[Top12]     Topaz, C. M. et al. "Locust Dynamics: Behavioral Phase Change and Swarming". en. *PLOS Computational Biology* **8**.8 (2012), e1002642.

[Vic95]     Vicsek, T. et al. "Novel type of phase transition in a system of self-driven particles". *Physical Review Letters* **75**.6 (1995). arXiv: cond-mat/0611743, pp. 1226–1229.

[Vos02]     Voss, H et al. "Nonlinear dynamical system identification from uncertain and indirect measurements." *Int J Bif Chaos* **14** (2002), pp. 1905–1924.

[War19]     Warne, D. J. et al. "Using Experimental Data and Information Criteria to Guide Model Selection for Reaction–Diffusion Problems in Mathematical Biology". en. *Bull Math Biol* **81**.6 (2019), pp. 1760–1804.

[Yan20]     Yang, L. et al. *B-PINNs: Bayesian Physics-Informed Neural Networks for Forward and Inverse PDE Problems with Noisy Data*. 2020. arXiv: 2003.06097 [stat.ML].

[Ye15]      Ye, H. et al. "Equation-free mechanistic ecosystem forecasting using empirical dynamic modeling". en. *Proceedings of the National Academy of Sciences* **112**.13 (2015), E1569–E1576.

[Yua04]     Yuan, G et al. "Estimating the predicability of an oceanic time series using linear and nonlinear methods". *J. Geophys. Res.* **109** (2004), p. C08002.

[Zha18]     Zhang, S. & Lin, G. "Robust data-driven discovery of governing physical laws with error bars". *Proc. R. Soc. A* **474**.2217 (2018), p. 20180305.

[Zha09]     Zhang, T. "Adaptive forward-backward greedy algorithm for sparse learning with linear models". *Advances in Neural Information Processing Systems*. 2009, pp. 1921–1928.

# APPENDICES

APPENDIX

## A

# HYBRID MODELING

## A.1   Parameter distributions

Each hybrid model has its own set of parameters that need to be estimated in the context of a given inverse problem. One question that arises is how do the parameter estimates vary among different model choices? Figure A.1 shows how the parameters vary among the hybrid models across all 21 data sets. From this we see that each parameter is roughly normally distributed with relatively no outliers.

## A.2   Parameter values

To get a closer look, in the following pages we present all of the individual parameter estimates for each hybrid (and full) model applied to each data set. There are 7 experimental conditions determined by the value of $c_{pa}$ which is used to quantify the survival probability of a pupa in the presence of adult flour beetles at a given time. For each of the 7 values of $c_{pa}$, there are three replicates which yields a total of 21 unique data sets.

**Figure A.1** Histogram of each of the parameter estimates among all 7 hybrid models across all 21 data sets.

**Table A.1** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 0, replicate: 1.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|-------|-------|-------|-------|-----|----------|----------|---------|
| SSA | - | - | 87.5656 | - | - | - | - |
| SPS | - | 17.8686 | - | - | - | - | 0.2182 |
| SPA | - | 0.3122 | 126.6186 | - | - | - | 0.2182 |
| LSS | 264.9027 | - | - | 8.6070 | 0.0084 | 0.0181 | - |
| LSA | 235.4855 | - | 108.3916 | 4.7567 | 0.0095 | 0.0117 | - |
| LPS | 289.5294 | -18.0937 | - | 7.0309 | 0.0246 | 0.0150 | 0.2135 |
| LPA | 267.6326 | 4.2396 | 112.3060 | 8.7363 | 0.0191 | 0.0128 | 0.2430 |

107

**Table A.2** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 0, replicate: 2.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 105.8837 | - | - | - | - |
| SPS | - | 42.6744 | - | - | - | - | 0.1729 |
| SPA | - | 19.4947 | 80.4809 | - | - | - | 0.1989 |
| LSS | 253.5396 | - | - | 8.9462 | -0.0015 | 0.0328 | - |
| LSA | 260.9683 | - | 94.5207 | 6.8836 | 0.0008 | 0.0255 | - |
| LPS | 164.6181 | 22.9624 | - | 8.4204 | 0.0026 | 0.0259 | 0.1569 |
| LPA | 257.5946 | 0.3385 | 21.9302 | 16.0772 | 0.0263 | 0.0143 | 0.2036 |

**Table A.3** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 0, replicate: 3.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 93.2149 | - | - | - | - |
| SPS | - | 3.4573 | - | - | - | - | 0.2470 |
| SPA | - | -1.1486 | 119.4034 | - | - | - | 0.2627 |
| LSS | 270.5929 | - | - | 5.7954 | 0.0081 | 0.0131 | - |
| LSA | 269.7489 | - | 117.1913 | 4.2792 | 0.0014 | 0.0180 | - |
| LPS | 233.8927 | 11.6517 | - | 4.4630 | 0.0126 | 0.0145 | 0.0083 |
| LPA | 320.5320 | 4.6590 | 106.3479 | 8.3933 | 0.0169 | 0.0113 | 0.3285 |

**Table A.4** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 05, replicate: 1.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 93.6535 | - | - | - | - |
| SPS | - | -9.6514 | - | - | - | - | 0.2523 |
| SPA | - | -10.2065 | 99.3673 | - | - | - | 0.2096 |
| LSS | 270.2898 | - | - | 5.2134 | 0.0055 | 0.0115 | - |
| LSA | 341.4744 | - | 114.3748 | 8.6488 | 0.0154 | 0.0162 | - |
| LPS | 267.5685 | 37.7745 | - | 5.7037 | 0.0105 | 0.0111 | 0.2443 |
| LPA | 301.3494 | -43.0801 | 130.9987 | 8.5879 | 0.0081 | 0.0174 | 0.1578 |

**Table A.5** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 05, replicate: 2.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 95.4528 | - | - | - | - |
| SPS | - | -59.6550 | - | - | - | - | 0.2412 |
| SPA | - | -5.7144 | 84.5839 | - | - | - | 0.2102 |
| LSS | 225.0079 | - | - | 6.3634 | 0.0159 | 0.0122 | - |
| LSA | 245.6442 | - | 93.2592 | 7.0331 | 0.0065 | 0.0176 | - |
| LPS | 257.2530 | 4.2089 | - | 7.1117 | 0.0194 | 0.0141 | 0.3901 |
| LPA | 169.1556 | 4.4937 | 107.6186 | 8.2696 | 0.0339 | 0.0127 | 0.1209 |

**Table A.6** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 05, replicate: 3.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 91.1832 | - | - | - | - |
| SPS | - | 6.7107 | - | - | - | - | 0.2362 |
| SPA | - | 27.9214 | 98.6414 | - | - | - | 0.1860 |
| LSS | 194.6385 | - | - | 6.6468 | 0.0097 | 0.0138 | - |
| LSA | 257.2298 | - | 139.9384 | 5.7910 | 0.0169 | 0.0119 | - |
| LPS | 256.7528 | 10.8231 | - | 3.5677 | 0.0170 | 0.0095 | 0.1807 |
| LPA | 247.1146 | 16.0416 | 122.5166 | 11.3146 | 0.0102 | 0.0227 | 0.2247 |

**Table A.7** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 10, replicate: 1.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 79.1244 | - | - | - | - |
| SPS | - | -1.0929 | - | - | - | - | 0.1982 |
| SPA | - | -30.0912 | 79.3693 | - | - | - | 0.2220 |
| LSS | 284.3427 | - | - | 7.1667 | 0.0203 | 0.0141 | - |
| LSA | 198.1921 | - | 77.9394 | 7.2244 | 0.0111 | 0.0144 | - |
| LPS | 223.2605 | -0.4444 | - | 5.8579 | 0.0419 | 0.0105 | 0.2058 |
| LPA | 171.9203 | -27.0930 | 71.1554 | 16.3337 | 0.0257 | 0.0187 | 0.2154 |

**Table A.8** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 10, replicate: 2.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 73.8780 | - | - | - | - |
| SPS | - | 20.0653 | - | - | - | - | 0.2848 |
| SPA | - | 0.8115 | 99.8802 | - | - | - | 0.2241 |
| LSS | 262.1010 | - | - | 6.2641 | 0.0081 | 0.0121 | - |
| LSA | 287.7510 | - | 115.5338 | 5.9965 | 0.0074 | 0.0130 | - |
| LPS | 234.6461 | 34.8923 | - | 6.8348 | 0.0030 | 0.0138 | 0.1482 |
| LPA | 185.1986 | 8.2924 | 141.0987 | 10.2721 | 0.0139 | 0.0190 | 0.1369 |

**Table A.9** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 10, replicate: 3.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 93.7121 | - | - | - | - |
| SPS | - | 6.4716 | - | - | - | - | 0.2284 |
| SPA | - | 33.0567 | 34.7667 | - | - | - | 0.2454 |
| LSS | 225.0365 | - | - | 11.7217 | 0.0091 | 0.0187 | - |
| LSA | 208.4282 | - | 59.9405 | 27.3955 | 0.0198 | 0.0289 | - |
| LPS | 203.9533 | 7.2728 | - | 9.7588 | 0.0120 | 0.0155 | 0.1035 |
| LPA | 201.6383 | -10.4708 | 128.0393 | 13.0707 | 0.0163 | 0.0158 | 0.1621 |

**Table A.10** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 25, replicate: 1.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 117.0499 | - | - | - | - |
| SPS | - | 3.0593 | - | - | - | - | 0.1395 |
| SPA | - | 49.3746 | 75.1271 | - | - | - | 0.1948 |
| LSS | 311.8632 | - | - | 18.3808 | 0.0191 | 0.0280 | - |
| LSA | 237.8213 | - | 96.1909 | 14.8132 | 0.0176 | 0.0215 | - |
| LPS | 205.1371 | 7.9529 | - | 8.6078 | 0.0085 | 0.0207 | -0.0051 |
| LPA | 205.7767 | -13.5553 | 119.2023 | 11.1854 | 0.0108 | 0.0183 | 0.0774 |

**Table A.11** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 25, replicate: 2.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 94.3111 | - | - | - | - |
| SPS | - | 46.3850 | - | - | - | - | 0.1855 |
| SPA | - | 7.3573 | 70.5916 | - | - | - | 0.2384 |
| LSS | 300.5890 | - | - | 8.5072 | 0.0115 | 0.0121 | - |
| LSA | 243.5118 | - | 99.8626 | 13.9622 | 0.0205 | 0.0145 | - |
| LPS | 320.7861 | -27.6091 | - | 10.6516 | 0.0100 | 0.0129 | 0.1880 |
| LPA | 225.4208 | 106.1507 | 111.8993 | 6.7875 | 0.0848 | 0.0051 | -0.0000 |

**Table A.12** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 25, replicate: 3.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 74.2800 | - | - | - | - |
| SPS | - | 5.5289 | - | - | - | - | 0.1968 |
| SPA | - | -22.5522 | 99.0768 | - | - | - | 0.1987 |
| LSS | 234.1530 | - | - | 3.5936 | 0.0104 | 0.0032 | - |
| LSA | 341.6403 | - | 64.9878 | 6.6088 | 0.0213 | 0.0136 | - |
| LPS | 273.6266 | -20.6973 | - | 5.3360 | 0.0068 | 0.0092 | 0.2640 |
| LPA | 245.0015 | 30.8203 | 155.5782 | 6.5300 | 0.0078 | 0.0131 | -0.0289 |

**Table A.13** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 35, replicate: 1.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 115.5389 | - | - | - | - |
| SPS | - | -18.4751 | - | - | - | - | 0.0224 |
| SPA | - | -34.9826 | 99.1463 | - | - | - | 0.0517 |
| LSS | 261.7244 | - | - | 9.1461 | 0.0171 | 0.0155 | - |
| LSA | 235.4640 | - | 113.0659 | 4.7266 | 0.0375 | 0.0082 | - |
| LPS | 241.1091 | 45.9677 | - | 10.0909 | 0.0151 | 0.0176 | 0.0426 |
| LPA | 175.6418 | 65.0242 | 70.0661 | 2.4126 | 0.1533 | 0.0019 | -0.7355 |

**Table A.14** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 35, replicate: 2.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 117.8667 | - | - | - | - |
| SPS | - | 15.0357 | - | - | - | - | 0.1671 |
| SPA | - | 1.7302 | 105.5667 | - | - | - | 0.1495 |
| LSS | 295.0432 | - | - | 7.2311 | 0.0059 | 0.0116 | - |
| LSA | 216.9816 | - | 141.3938 | 10.8090 | 0.0407 | 0.0150 | - |
| LPS | 284.7066 | 97.4666 | - | 5.2911 | 0.0037 | 0.0101 | 0.1484 |
| LPA | 280.2145 | -0.1963 | 107.4857 | 6.9163 | 0.0173 | 0.0104 | 0.0934 |

**Table A.15** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 35, replicate: 3.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 76.9842 | - | - | - | - |
| SPS | - | 1.6093 | - | - | - | - | 0.2265 |
| SPA | - | 9.5439 | 80.1815 | - | - | - | 0.1938 |
| LSS | 318.6050 | - | - | 6.9155 | 0.0124 | 0.0110 | - |
| LSA | 214.6931 | - | 66.1158 | 7.1972 | 0.0525 | 0.0103 | - |
| LPS | 264.4815 | -9.1458 | - | 6.6379 | 0.0117 | 0.0106 | 0.1854 |
| LPA | 203.1091 | -11.1508 | 110.6879 | 8.8650 | 0.0935 | 0.0075 | 0.1933 |

**Table A.16** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 50, replicate: 1.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 128.0039 | - | - | - | - |
| SPS | - | 15.0019 | - | - | - | - | 0.1666 |
| SPA | - | 24.1235 | 94.8902 | - | - | - | 0.2008 |
| LSS | 120.0218 | - | - | 7.2009 | 0.0085 | 0.0126 | - |
| LSA | 391.7925 | - | 139.1012 | 11.9511 | 0.0110 | 0.0273 | - |
| LPS | 249.4874 | 57.5170 | - | 7.6681 | 0.0075 | 0.0147 | 0.1601 |
| LPA | 265.8448 | 51.1741 | 75.4596 | 6.2343 | 0.0673 | 0.0084 | 0.2951 |

**Table A.17** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 50, replicate: 2.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 99.8908 | - | - | - | - |
| SPS | - | 6.6477 | - | - | - | - | 0.1566 |
| SPA | - | -5.8578 | 102.9241 | - | - | - | 0.1320 |
| LSS | 225.0992 | - | - | 7.0927 | 0.0072 | 0.0120 | - |
| LSA | 240.7469 | - | 32.2565 | 11.1816 | 0.0062 | 0.0179 | - |
| LPS | 276.5789 | -80.6655 | - | 5.0354 | 0.0059 | 0.0094 | -0.0032 |
| LPA | 200.3636 | 57.9480 | 182.4789 | 8.7927 | 0.0109 | 0.0117 | 0.0198 |

**Table A.18** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 50, replicate: 3.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 109.8883 | - | - | - | - |
| SPS | - | -6.7587 | - | - | - | - | 0.2016 |
| SPA | - | -12.8388 | 102.8847 | - | - | - | 0.2092 |
| LSS | 284.9632 | - | - | 10.3289 | 0.0766 | 0.0141 | - |
| LSA | 255.4330 | - | 110.8994 | 7.0262 | 0.1958 | 0.0125 | - |
| LPS | 214.9260 | 30.3981 | - | 10.0889 | 0.0236 | 0.0140 | 0.1016 |
| LPA | 202.2075 | 13.8686 | 123.5118 | 11.4685 | 0.0322 | 0.0138 | 0.0289 |

**Table A.19** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 100, replicate: 1.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 73.4071 | - | - | - | - |
| SPS | - | -59.6402 | - | - | - | - | 0.1925 |
| SPA | - | 29.5418 | 161.9890 | - | - | - | 0.2680 |
| LSS | 263.0473 | - | - | 6.6799 | 0.0212 | 0.0102 | - |
| LSA | 307.2219 | - | 87.6893 | 19.9009 | 0.0365 | 0.0180 | - |
| LPS | 266.0186 | -11.8024 | - | 9.5711 | 0.0121 | 0.0121 | 0.2868 |
| LPA | 217.2470 | -28.3068 | 115.2743 | 21.3916 | 0.0189 | 0.0185 | 0.2453 |

**Table A.20** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 100, replicate: 2.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 137.2379 | - | - | - | - |
| SPS | - | 3.3028 | - | - | - | - | 0.2330 |
| SPA | - | 16.6531 | 19.8224 | - | - | - | 0.1920 |
| LSS | 296.7728 | - | - | 9.1518 | 0.1006 | 0.0127 | - |
| LSA | 260.3605 | - | 138.0497 | 9.1080 | 0.0383 | 0.0142 | - |
| LPS | 285.7502 | 21.4619 | - | 9.9246 | 0.0110 | 0.0128 | 0.2222 |
| LPA | 281.3465 | -0.6882 | 85.0033 | 11.7834 | 0.0074 | 0.0162 | -0.2808 |

**Table A.21** Parameter estimates for each hybrid model for the data set with $c_{pa}$: 100, replicate: 3.

| Model | $L_0$ | $P_0$ | $A_0$ | $b$ | $c_{el}$ | $c_{ea}$ | $\mu_1$ |
|---|---|---|---|---|---|---|---|
| SSA | - | - | 59.7707 | - | - | - | - |
| SPS | - | -13.8199 | - | - | - | - | 0.2208 |
| SPA | - | 10.5866 | 72.3512 | - | - | - | 0.1208 |
| LSS | 232.6881 | - | - | 9.6878 | 0.0195 | 0.0122 | - |
| LSA | 271.5676 | - | 185.1406 | 14.3152 | 0.0576 | 0.0168 | - |
| LPS | 258.6513 | 2.7035 | - | 9.6504 | 0.0203 | 0.0130 | 0.1579 |
| LPA | 233.4159 | 21.9667 | 106.5528 | 14.2631 | 0.0739 | 0.0140 | 0.0277 |

$$B$$

# NEURAL NETWORK DATA DENOISING

## B.1  Error Model Selection

This section is based on [Ban14a], and we are concerned with the selection of an appropriate error model from a given data set, $\{U_{i,j}\}_{i=1,j=1}^{M,N}$. A common and flexible error model is given by the nonconstant variance error model

$$U_{i,j} = f(x_i, t_j|\theta_0) + f^\gamma(x_i, t_j|\theta_0)\epsilon_{i,j}, \tag{B.1}$$

in which $\epsilon_{i,j} \overset{i.i.d.}{\sim} \mathcal{N}(0,\sigma^2)$. This error model is flexible in that it can account for both constant variance error models (when $\gamma = 0$) and nonconstant variance error model (when $\gamma > 0$). In the latter case, we can quantify the extent of the variance's dependence on $f(x,t)$ with estimation of $\gamma$. Computing and plotting residuals is a useful way to estimate an appropriate value of $\gamma$ from data.

From Equation (B.1), we observe that the *modified residuals*,

$$r_{i,j} = \frac{U_{i,j} - f(x_i, t_j|\theta)}{f^\gamma(x_i, t_j|\theta)}, i = 1, ..., M; j = 1, ..., N \tag{B.2}$$
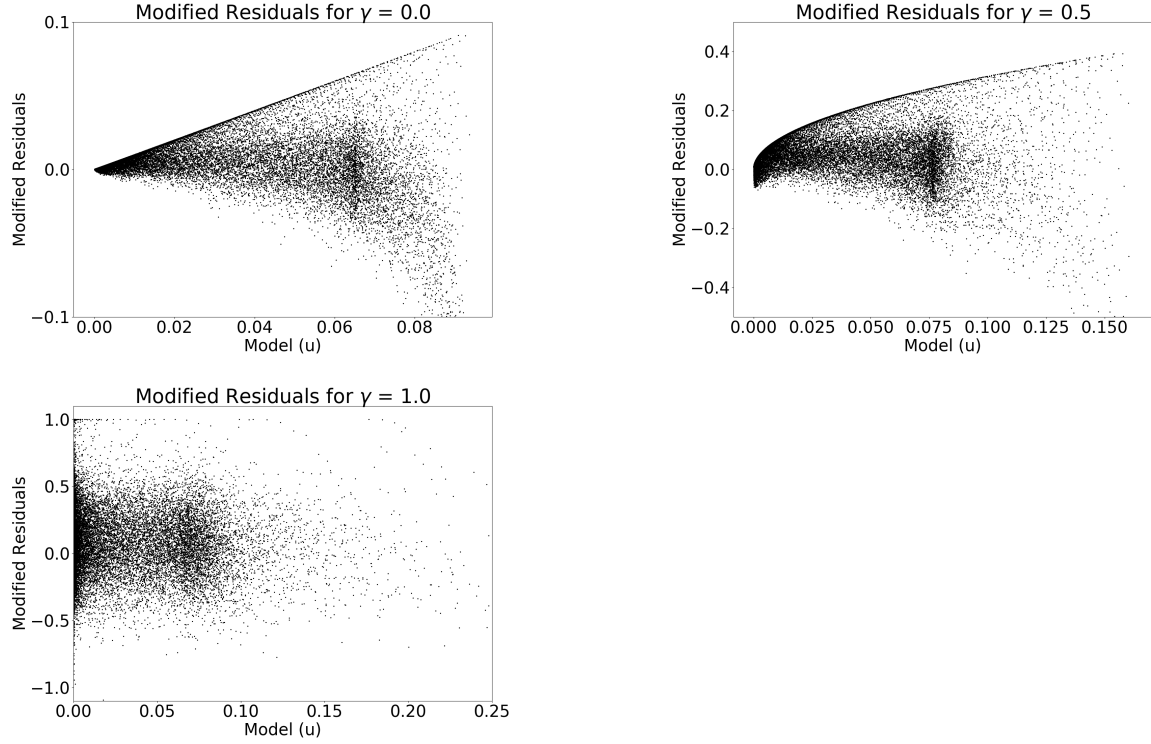
**Figure B.1** Modified Residual computations for various values of $\gamma$ from the advection-diffusion data set with $\sigma = 0.25$. Top left: results for $\gamma = 0$, Top right: results for $\gamma = 0.5$, bottom left: results for $\gamma = 1.0$.

should be $MN$ realizations of an $i.i.d.$ random variable when $\gamma$ has been chosen correctly and $\theta \approx \theta_0$. One can thus choose the correct form of Equation B.1 for a given data set and mathematical, $f(x, t; \theta)$, as follows: (i). Pick a value of $\gamma$, (ii). compute $\hat{\theta}$ by minimizing the generalized cost function, i.e., $\hat{\theta} = \arg\min_\theta \sum_{i=1, j=1}^{M,N} r_{i,j}^2$, and (iii). compute and plot each $r_{i,j}$. For the correct value of $\gamma$, the plotted modified residual computations will appear $i.i.d.$ Note this method can be used for error models different from Equation (B.1). For example, Nardini and Bortz [Nar19] used residual computations to demonstrate that a spatially-autocorrelated error model can account for numerical error arising from a PDE's discretization scheme during an inverse problem methodology.

As an example, we have plotted the modified results for the advection-diffusion dataset with $\sigma = 0.25$ in Figure B.1 for $\gamma = 0, 0.5,$ and 1.0. For each value of $\gamma$, we trained the ANN on the assumption that data was of the form given by Equation (B.1). For $\gamma = 0$ and 0.5, we observe that the residuals do not appear $i.i.d.$, as they fan out with increasing values of $u(x, t)$. At $\gamma = 1.0$, however, we see that the variance of the modified residuals appears constant, suggesting that $\gamma = 1.0$ is the appropriate value from the data.

114

## B.2    Comparing spline and bi-spline methods

We compared the accuracy in learning the correct PDE when using 1-dimensional cubic CV splines versus cubic CV bi-splines for denoising data and approximating partial derivatives (Figures B.2,B.3,B.4). We found that PDE-FIND with pruning always has a higher TPR value when using bi-spline computations as compared to 1-dimensional splines.
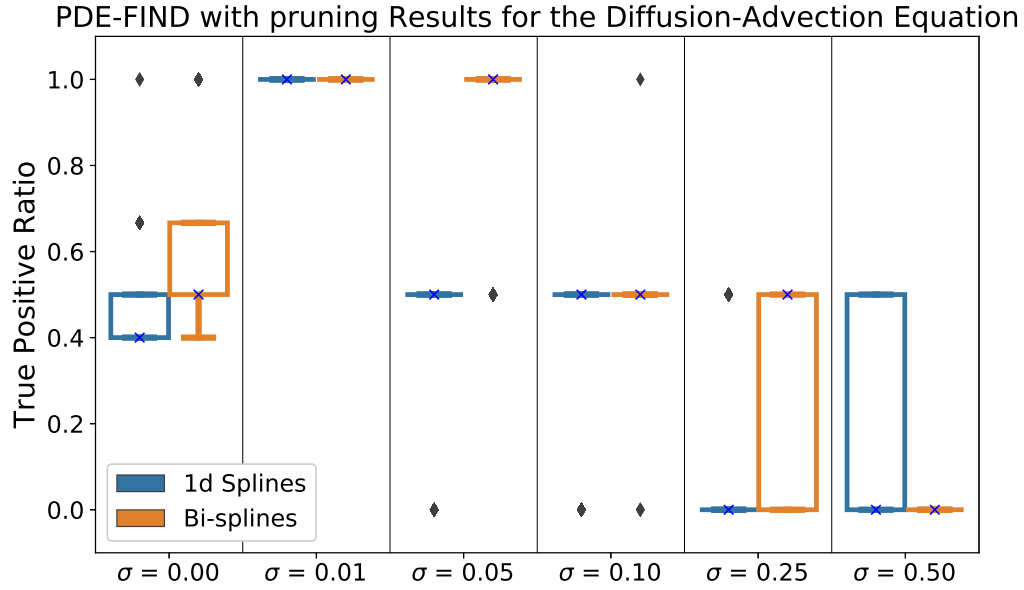


**Figure B.2** TPR values for the diffusion-advection equation when using 1-dimensional cubic splines versus cubic bi-splines for denoising data and approximating partial derivatives.

## B.3    Global Spline Calculations

We are concerned with approximating $\{U_i, j\}_{i=1,j=1}^{M,N}$ with a global spline representation of the form

$$S(x,t) = \sum_{k=1,l=1}^{K,L} c_{k,l} S_{k,l}(x,t), \tag{B.3}$$

where $S_{k,l}(x,t)$ are normalized bivariate cubic B-splines defined on the knot locations
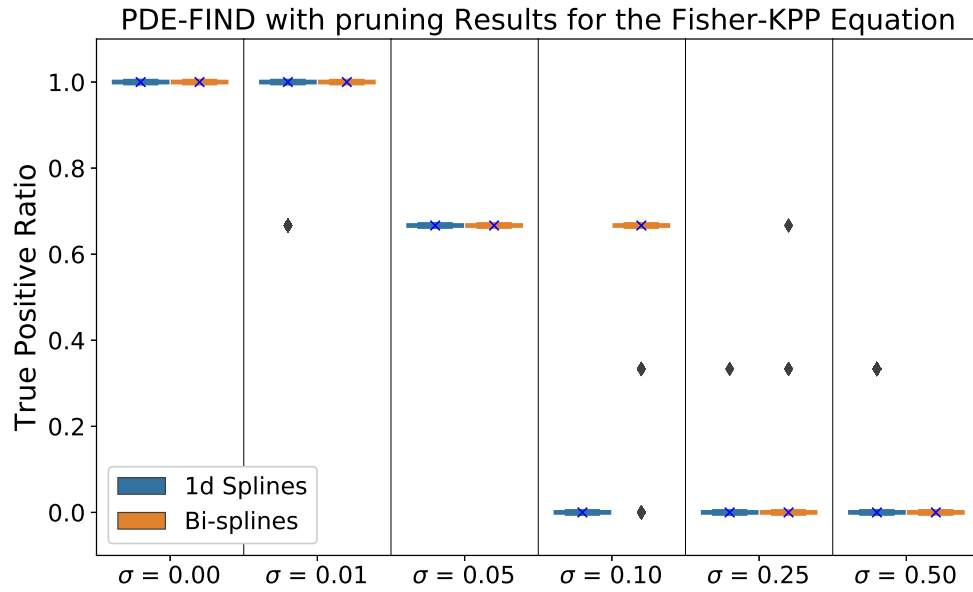
**Figure B.3** TPR values for the Fisher-KPP equation when using 1-dimensional cubic splines versus cubic bi-splines for denoising data and approximating partial derivatives.
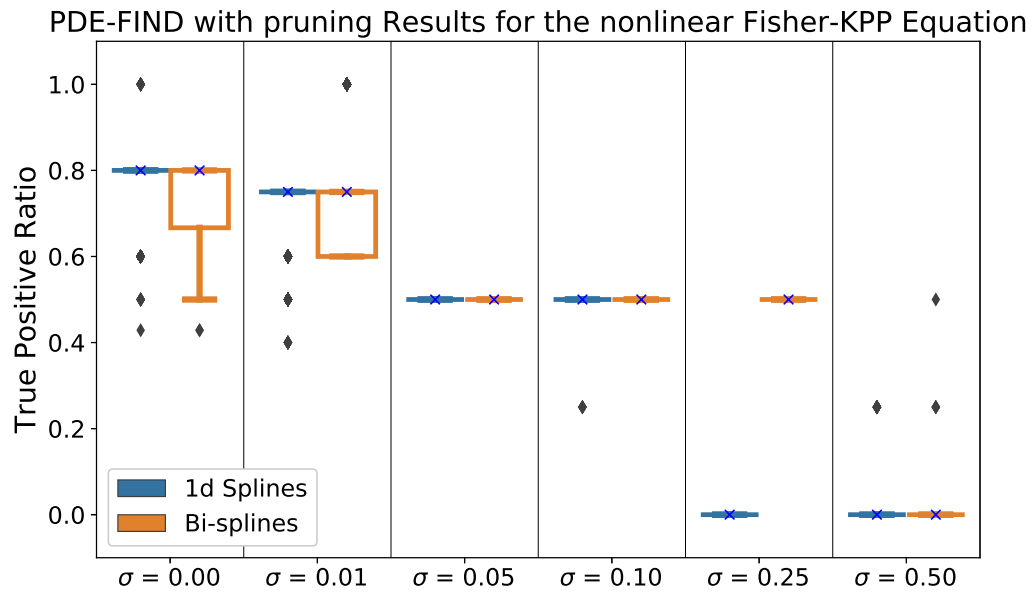


**Figure B.4** TPR values for the nonlinear Fisher-KPP equation when using 1-dimensional cubic splines versus cubic bi-splines for denoising data and approximating partial derivatives.

$(x, t)_{k,l} = \{\tilde{x}_k, ..., \tilde{x}_{k+4}\} \times \{\tilde{t}_l, ..., \tilde{t}_{l+4}\}$. In order to do so, we need to estimate a smoothing parameter, $s$, the knot locations, $\{(x, t)_{k,l}\}_{k=1,l=1}^{K,L}$, and the spline coefficients, $c = \{c_{k,l}\}_{k=1,l=1}^{K,L}$. To obtain these estimates, we split $\{U_{i,j}\}_{i=1,j=1}^{M,N}$ into a training and validation set (50%/50%). Recall that bivariate splines need to be fit to data on a rectangular grid domain, so we maintain this structure in the training and validation set by selecting all spatial points but only every other time point for the training data. The remaining points are placed in the validation data. While the (50%/50%) training and validation split here is not equivalent to the (90%/10%) split used to train the ANN, we found that implementing the global splines on a (90%/10%) took too much time for practical computation (training one data set took over a day on an Intel i7 6-core 3.5GHz desktop computer).

To find the optimal value of $s$, we set we begin with $s = MN + \sqrt{2MN}$, a previously-proposed upper bound on this smoothing parameter [Die81], and find the knot locations and coefficient values that minimize the GLS cost function $\mathscr{J}_S(\theta)$ (Equation (2.4) from in the main text) on the validation data. Further description of identification of spline locations of coefficients is given below. We then continue to divide $s$ in half and re-compute $S(x, t)$ for the updated values of $s$ until $\mathscr{J}_S(\theta)$ begins to increase on the validation data. We then compute $S(x, t)$ for a finer grid of $s$ values around whichever value minimized $\mathscr{J}_S(\theta)$ and ultimately select whichever of these value minimize $\mathscr{J}_S(\theta)$ on the validation data.

For a given value of $s$, we find the spline locations and coeffects in a similar manner to that described in [Ban14a]. We first fit the global spline model to the training data using a constant variance (OLS) error model. This first-pass spline computation allows us to estimate $u(x_i, t_j)$ for the training data and in turn estimate $\mathscr{J}(\theta)$. We now iteratively fit the cubic spline method to the training data using $\mathscr{J}(\theta)$ with $\gamma = 1.0$ until the number of knot locations (determined by the Scipy **bisplev** algorithm) does not change by more than one over five consecutive calculations. We then fix the knot locations and iteratively estimate the spline coefficients, $c_{k,l}$, by minimizing $\mathscr{J}(\theta)$ with $\gamma = 1$ until either the maximum of the relative absolute difference between consecutive estimates, *i.e.*, the inf-norm, converges within $10^{-2}$ or 100 such computations have been performed.

## B.4 PDE-FIND without pruning results

We found that using PDE-FIND without pruning results in learning the wrong equation when applied to data from biological transport models, even when no noise is added to the data. We evaluated accuracy, using the true positive ratio (TPR) as a metric, for the diffusion-advection (Figure B.5), Fisher-KPP (Figure B.6), and nonlinear Fisher-KPP equations (Figure B.7).

For the diffusion-advection equation, we found that the TPR value of the final learned equation when using ANN approximations is higher for all values of $\sigma$ when using pruning with PDE-FIND than without pruning (Figure B.5). In general, for small values of $\sigma$, we observed that pruning enables PDE-FIND to better learn the true equation when using CV local spline and finite difference computations, but it harms the ability to learn the true equation for larger values of $\sigma$. For example, the median TPR value increases after pruning when using finite difference approximations from TPR = 0.33 to 0.5 for $\sigma = 0$. However, the TPR instead decreases from TPR = 0.33 to 0 at $\sigma = 0.05$ and from TPR = 0.5 to 0 at $\sigma = 0.10$. The median TPR value when using spline approximations increases from TPR = 0.3 to 0.5 at $\sigma = 0$ and from TPR = 0.33 to 1 at $\sigma = 0.01$. At $\sigma = 0.10$, the median values decreased from TPR = 1.0 to 0.5.
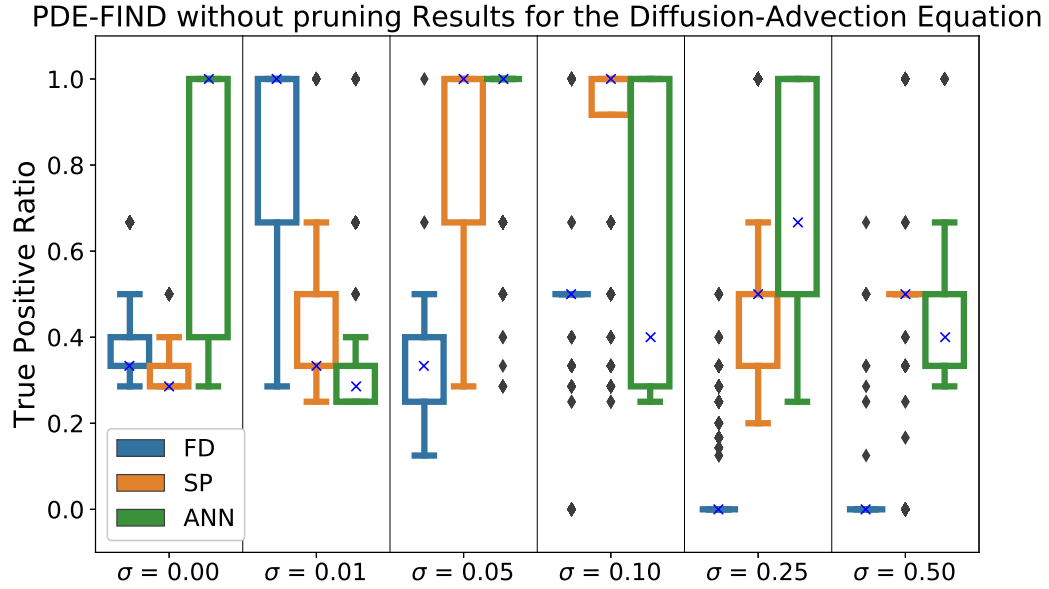


**Figure B.5** TPR values for the diffusion-advection equation.

For the Fisher-KPP equation, the median TPR value for PDE-FIND with the ANN computations always increases after using pruning (Figure B.6). The median value for PDE-FIND with finite difference computations increases for $\sigma = 0, 0.01$, but decreases from TPR = 0.5 to 0 for $\sigma = 0.05$ and from TPR = 0.45 to 0 for $\sigma = 0.10$. The median TPR value for PDE-FIND with CV local spline computations increases for $\sigma = 0, 0.01, 0.05$, and 0.10, but decreases from TPR = 0.4 to 0 at both $\sigma = 0.25$ and 0.50. Thus, pruning always helped PDE-FIND
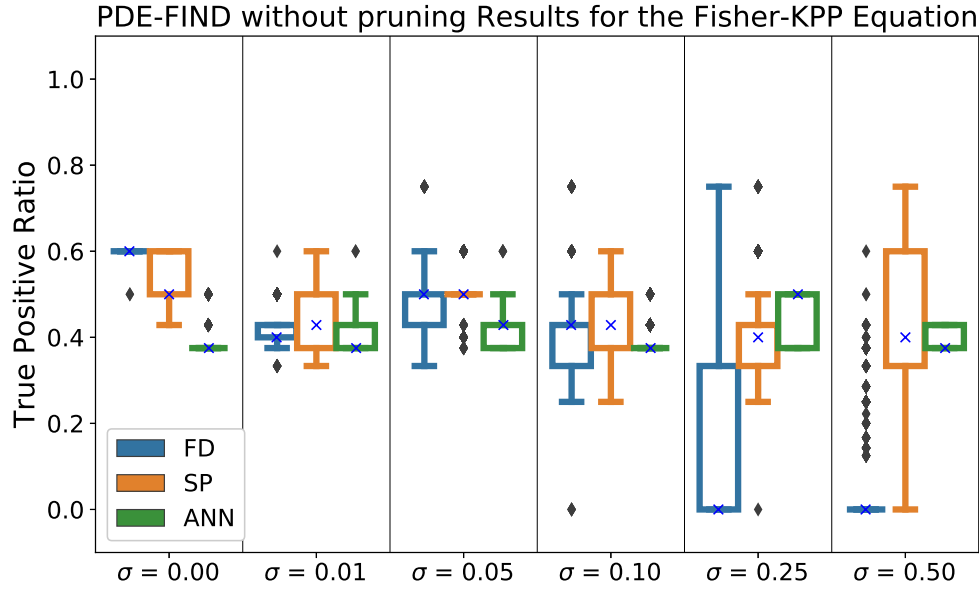
**Figure B.6** TPR values for the Fisher-KPP equation.

learn the true equation when using the ANN method, and helps the other computational methods for small noise levels.

For the nonlinear Fisher-KPP qquation, the median TPR value always improved the accuracy of the PDE-FIND method when using ANN approximations (Figure B.7). When using finite difference approximations, the median TPR value increases for $\sigma = 0, 0.01.05$. The median value decreases from TPR $= 0.3$ to $0$ at $\sigma = 0.10$ for finite difference approximations. When using CV local spline approximations, the median TPR value increases when $\sigma = 0$ and $0.01$. The median TPR value decreased from TPR $= 0.3$ to $0$ at $\sigma = 0.50$. While the median value is never TPR $= 1$ for this equation, these results suggest that pruning in general helps reduce the number of incorrect terms in the library.

## B.5 Tables of learned PDEs

This section contains tables of the final learned PDEs for data from each equation considered at a given noise level. The equation form is the one most commonly selected by the PDE-FIND method with pruning over the 1,000 different training-validation splits of $u_t$ and $\Theta$. The provided parameter values are the mean value for these parameters when the equation form was the final learned equation.

**Table B.1** Learned Equations for the diffusion-advection equation.

| $\sigma$ | Method | True Equation |
|---|---|---|
| | | $u_t = 0.01\,u_{xx} - 0.8\,u_x$ |
| $\sigma$ | **Method** | **Learned Equation** |
| 0.0 | FD | $u_t = -0.799974\,u_x + 0.010414\,u_{xx} - 0.000350\,u^2 u_x$ |
| 0.01 | FD | $u_t = -0.800226\,u_x + 0.009940\,u_{xx}$ |
| 0.05 | FD | $u_t = 0$ |
| 0.10 | FD | $u_t = 0$ |
| 0.25 | FD | $u_t = 0$ |
| 0.50 | FD | $u_t = 0$ |
| 0.0 | LCVSP | $u_t = -0.807744\,u_x + 0.011464\,u_{xx} + 0.000670\,u^2 u_x + 0.000012\,u^2 u_{xx}$ |
| 0.01 | LCVSP | $u_t = -0.793993\,u_x + 0.011877\,u_{xx}$ |
| 0.05 | LCVSP | $u_t = -0.796512\,u_x + 0.012290\,u_{xx}$ |
| 0.10 | LCVSP | $u_t = -0.774238\,u_x$ |
| 0.25 | LCVSP | $u_t = -0.709333\,u_x$ |
| 0.50 | LCVSP | $u_t = 0$ |
| 0.0 | LNCVSP | $u_t = -0.820792\,u_x + 0.011752\,u_{xx}$ |
| 0.01 | LNCVSP | $u_t = -0.819847\,u_x + 0.011726\,u_{xx}$ |
| 0.05 | LNCVSP | $u_t = -0.819094\,u_x + 0.012085\,u_{xx}$ |
| 0.10 | LNCVSP | $u_t = -0.790680\,u_x$ |
| 0.25 | LNCVSP | $u_t = -0.760747\,u_x$ |
| 0.50 | LNCVSP | $u_t = -0.729489\,u_x$ |
| 0.0 | GNCVSP | $u_t = -0.764276\,u_x$ |
| 0.01 | GNCVSP | $u_t = -0.694783\,u_x$ |
| 0.05 | GNCVSP | $u_t = -0.611367\,u_x$ |
| 0.10 | GNCVSP | $u_t = -0.581202\,u_x$ |
| 0.25 | GNCVSP | $u_t = 0$ |
| 0.50 | GNCVSP | $u_t = 0$ |
| 0.0 | ANN | $u_t = -0.809223\,u_x + 0.010963\,u_{xx}$ |
| 0.01 | ANN | $u_t = -0.802903\,u_x + 0.011107\,u_{xx} - 0.000074\,u^2 u_{xx} + 0.000765\,u_x^2$ |
| 0.05 | ANN | $u_t = -0.810360\,u_x + 0.010693\,u_{xx}$ |
| 0.10 | ANN | $u_t = -0.808996\,u_x + 0.009535\,u_{xx}$ |
| 0.25 | ANN | $u_t = -0.795770\,u_x + 0.009386\,u_{xx}$ |
| 0.50 | ANN | $u_t = -0.801987\,u_x + 0.007888\,u_{xx} + 0.000835\,u_x^2$ |

**Table B.2** Discovered Equations for the Fisher-KPP Equation

| $\sigma$ | Method | Learned Equation |
|---|---|---|
| | | **True Equation** |
| | | $u_t = 0.02\,u_{xx} + 10.0\,u - 10.0\,u^2$ |
| 0.0 | FD | $u_t = 0.020086\,u_{xx} - 9.994321\,u^2 + 9.995583\,u$ |
| 0.01 | FD | $u_t = -10.154641\,u^2 + 9.950923\,u$ |
| 0.05 | FD | $u_t = 0$ |
| 0.10 | FD | $u_t = 0$ |
| 0.25 | FD | $u_t = 0$ |
| 0.50 | FD | $u_t = 0$ |
| 0.0 | LCVSP | $u_t = 0.020474\,u_{xx} - 9.993088\,u^2 + 9.996574\,u$ |
| 0.01 | LCVSP | $u_t = 0.019600\,u_{xx} - 9.972429\,u^2 + 9.977920\,u$ |
| 0.05 | LCVSP | $u_t = -10.130441\,u^2 + 9.925709\,u$ |
| 0.10 | LCVSP | $u_t = -10.087935\,u^2 + 9.916230\,u$ |
| 0.25 | LCVSP | $u_t = 0$ |
| 0.50 | LCVSP | $u_t = 0$ |
| 0.0 | LNCVSP | $u_t = 0.020435\,u_{xx} - 9.991312\,u^2 + 9.994767\,u$ |
| 0.01 | LNCVSP | $u_t = 0.019522\,u_{xx} - 9.977385\,u^2 + 9.982515\,u$ |
| 0.05 | LNCVSP | $u_t = -10.121782\,u^2 + 9.916090\,u$ |
| 0.10 | LNCVSP | $u_t = -10.087677\,u^2 + 9.926292\,u$ |
| 0.25 | LNCVSP | $u_t = 0$ |
| 0.50 | LNCVSP | $u_t = 0$ |
| 0.0 | GNCVSP | $u_t = -10.264500\,u^2 + 10.229065\,u$ |
| 0.01 | GNCVSP | $u_t = -8.598153\,u^2 + 9.375428\,u$ |
| 0.05 | GNCVSP | $u_t = -10.346971\,u^2 + 10.122676\,u$ |
| 0.10 | GNCVSP | $u_t = -10.007866\,u^2 + 10.075741\,u$ |
| 0.25 | GNCVSP | $u_t = -9.312518\,u^2 + 9.304600\,u$ |
| 0.50 | GNCVSP | $u_t = -5.621682\,u^2 + 7.104374\,u$ |
| 0.0 | ANN | $u_t = 0.023272\,u_{xx} - 9.307794\,u^2 + 9.533177\,u$ |
| 0.01 | ANN | $u_t = 0.023017\,u_{xx} - 9.397175\,u^2 + 9.600546\,u$ |
| 0.05 | ANN | $u_t = 0.020534\,u_{xx} - 9.733768\,u^2 + 9.837442\,u$ |
| 0.10 | ANN | $u_t = 0.022343\,u_{xx} - 9.287166\,u^2 + 9.587605\,u$ |
| 0.25 | ANN | $u_t = 0.011912\,u_{xx} - 11.160631\,u^2 + 12.537031\,u$ $+ 0.071219\,u\,u_{xx} - 0.105350\,u_x^2$ |
| 0.50 | ANN | $u_t = -0.015750\,u_x + 0.013682\,u_{xx} - 8.688903\,u^2$ $+ 12.179728\,u - 0.034142\,u^2\,u_x + 0.077472\,u\,u_{xx} - 0.109284\,u_x^2$ |

**Table B.3** Discovered Equations for the nonlinear Fisher-KPP Equation.

| | | True Equation |
|---|---|---|
| | | $u_t = 0.020000\,u\,u_{xx} + 0.020000\,u_x^2 + 10.000000\,u - 10.000000\,u^2$ |
| $\sigma$ | **Method** | **Learned Equation** |
| 0.0 | FD | $u_t = 0.000178\,u_{xx} - 9.996233\,u^2 + 9.996516\,u + 0.019671\,u\,u_{xx} + 0.019729\,u_x^2$ |
| 0.01 | FD | $u_t = -10.268294\,u^2 + 10.210714\,u$ |
| 0.05 | FD | $u_t = -9.531702\,u^2 + 9.731417\,u$ |
| 0.10 | FD | $u_t = 0$ |
| 0.25 | FD | $u_t = 0$ |
| 0.50 | FD | $u_t = 0$ |
| 0.0 | LCVSP | $u_t = 0.000760\,u_{xx} - 10.013375\,u^2 + 10.013144\,u + 0.018798\,u\,u_{xx} + 0.018324\,u_x^2$ |
| 0.01 | LCVSP | $u_t = 0.005632\,u_{xx} - 9.820662\,u^2 + 9.789693\,u + 0.017619\,u_x^2$ |
| 0.05 | LCVSP | $u_t = -10.393255\,u^2 + 10.287537\,u$ |
| 0.10 | LCVSP | $u_t = -10.264929\,u^2 + 10.194679\,u$ |
| 0.25 | LCVSP | $u_t = -10.146329\,u^2 + 10.082739\,u$ |
| 0.50 | LCVSP | $u_t = 0$ |
| 0.0 | LNCVSP | $u_t = 0.000738\,u_{xx} - 10.011882\,u^2 + 10.011901\,u + 0.018874\,u\,u_{xx} + 0.018394\,u_x^2$ |
| 0.01 | LNCVSP | $u_t = 0.005585\,u_{xx} - 9.828360\,u^2 + 9.796878\,u + 0.017437\,u_x^2$ |
| 0.05 | LNCVSP | $u_t = -10.396293\,u^2 + 10.285593\,u$ |
| 0.10 | LNCVSP | $u_t = -10.333053\,u^2 + 10.259315\,u$ |
| 0.25 | LNCVSP | $u_t = -9.875062\,u^2 + 9.794690\,u$ |
| 0.50 | LNCVSP | $u_t = 0$ |
| 0.0 | GNCVSP | $u_t = -0.025827\,u_{xx} - 10.407031\,u^2 + 10.439599\,u$ |
| 0.01 | GNCVSP | $u_t = -0.010336\,u_{xx} - 10.663675\,u^2 + 10.575536\,u$ |
| 0.05 | GNCVSP | $u_t = -0.009571\,u_{xx} - 10.670700\,u^2 + 10.563083\,u$ |
| 0.10 | GNCVSP | $u_t = -9.792682\,u^2 + 9.824648\,u$ |
| 0.25 | GNCVSP | $u_t = -0.017669\,u_{xx} - 10.233510\,u^2 + 10.177806\,u$ |
| 0.50 | GNCVSP | $u_t = -0.023890\,u_{xx} - 9.392645\,u^2 + 9.576821\,u$ |
| 0.0 | ANN | $u_t = 0.009869\,u_{xx} - 9.295491\,u^2 + 9.237594\,u - 0.032329\,u\,u_{xx} + 0.016924\,u_x^2$ |
| 0.01 | ANN | $u_t = -9.398164\,u^2 + 9.389853\,u + 0.024833\,u_x^2$ |
| 0.05 | ANN | $u_t = 0.009080\,u_{xx} - 9.311857\,u^2 + 9.245655\,u - 0.032236\,u\,u_{xx} + 0.016673\,u_x^2$ |
| 0.10 | ANN | $u_t = -0.006456\,u_{xx} - 8.965042\,u^2 + 9.140129\,u + 0.027012\,u_x^2$ |
| 0.25 | ANN | $u_t = -0.010203\,u_{xx} - 7.927777\,u^2 + 8.551556\,u + 0.034169\,u_x^2$ |
| 0.50 | ANN | $u_t = 0.285757 - 0.026084\,u_{xx} - 5.419017\,u^2 + 6.724382\,u + 0.044730\,u_x^2$ |

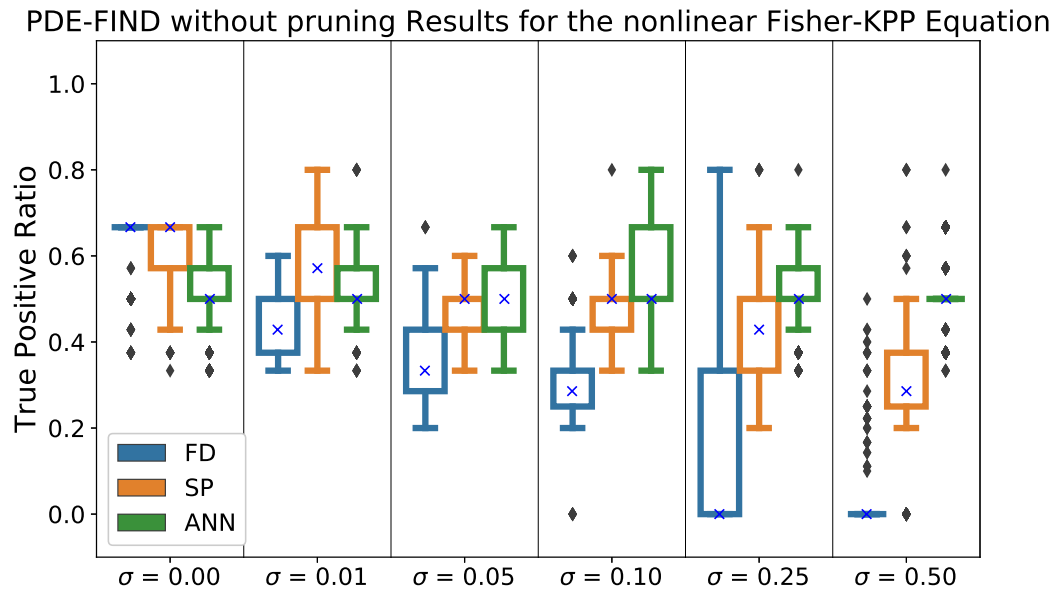**Figure B.7** TPR values for the nonlinear Fisher-KPP equation.

# C

# BIOLOGICALLY-INFORMED NEURAL NETWORKS
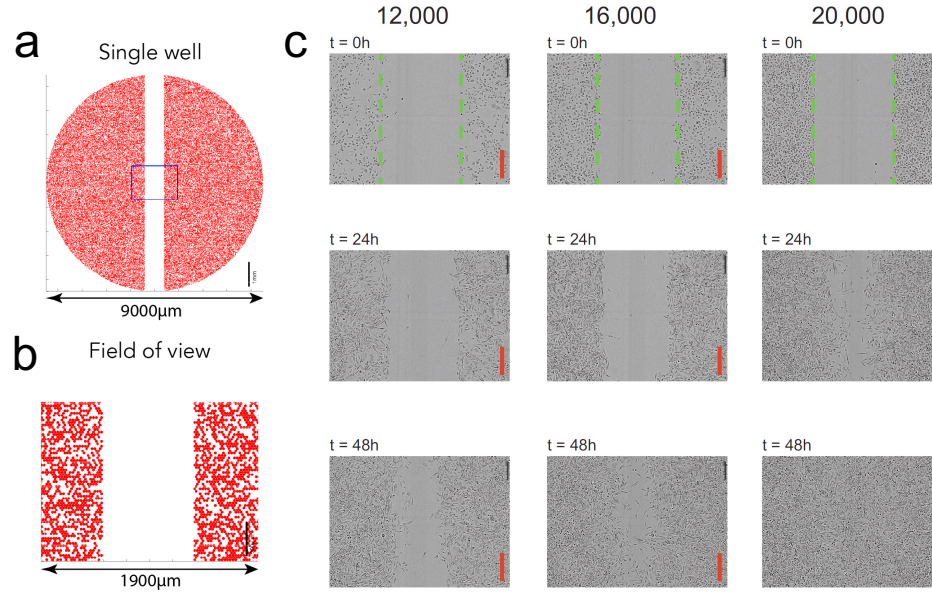
# C.1 Supporting information



**Figure C.1** (a) An illustration of an experiment with the IncuCyte ZOOM™ system (Essen Bio-Science, MI USA). Full details of the experiment and image processing can be found in [Jin16a]. Cells are seeded uniformly within each well in a 96-well plate at a pre-specified density of between 10,000 and 20,000 cells per well. A WoundMaker™ (Essen BioScience) is used to create a uniform vertical scratch along the middle of the well. (b) Microscopy images are collected from a rectangular region of the well. (c) Example images corresponding to experiments initiated with 12,000, 16,000, or 20,000 cells per well. A PC-3 prostate cancer cell line was used. The image recording time is indicated on each subfigure and the scale bar corresponds to 300 $\mu$m. The green dashed lines in the images in the top row show the approximate location of the leading edge created by the scratch. Each image is divided into equally-spaced vertical columns, and the number of cells in each column divided by the column area is calculated to yield an estimate of the 1-D cell density.
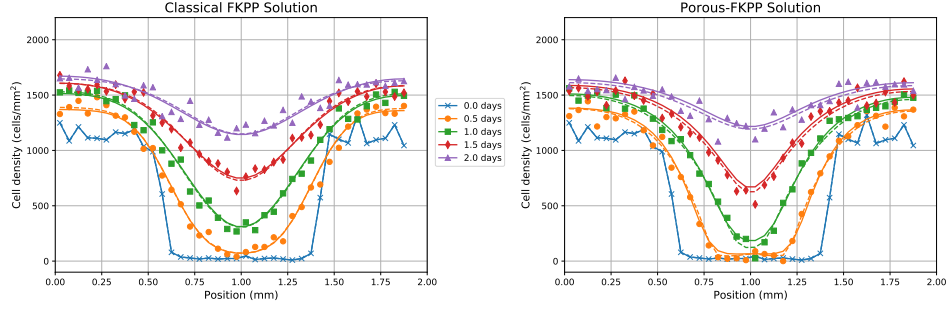
**Figure C.2** Predicted cell density profiles using BINNs with the governing reaction-diffusion PDE in Equation (4.15). The left subplot corresponds to the set of simulated data using the classical FKPP equation and the rig!ht subplot corresponds to the Generalized Porous-FKPP equation. Solid lines represent the numerical solution to Equation (4.15) using $D_{\text{MLP}}$, and $G_{\text{MLP}}$. Dashed lines represent the noiseless numerical simulations of the classical FKPP and Generalized Porous-FKPP equations. The markers represent the numerical simulations of the classical FKPP and Generalized Porous-FKPP equations with artificial noise generated by the statistical error model in Equation (4.4).
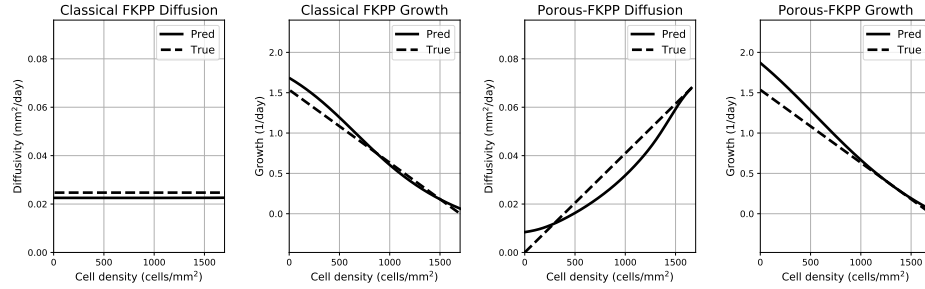


**Figure C.3** The learned diffusivity and growth functions $D_{\text{MLP}}$ and $G_{\text{MLP}}$ evaluated over cell density $u$. Starting from the left, the first two subplots correspond to the learned diffusivity and growth functions from simulated data using the classical FKPP equation. The last two subplots correspond to the learned diffusivity and growth functions from simulated data using the Generalized Porous-FKPP equation. Solid lines represent the parameter networks $D_{\text{MLP}}$ and $G_{\text{MLP}}$ and dashed lines represent the true diffusivity and growth functions used to simulate the data.
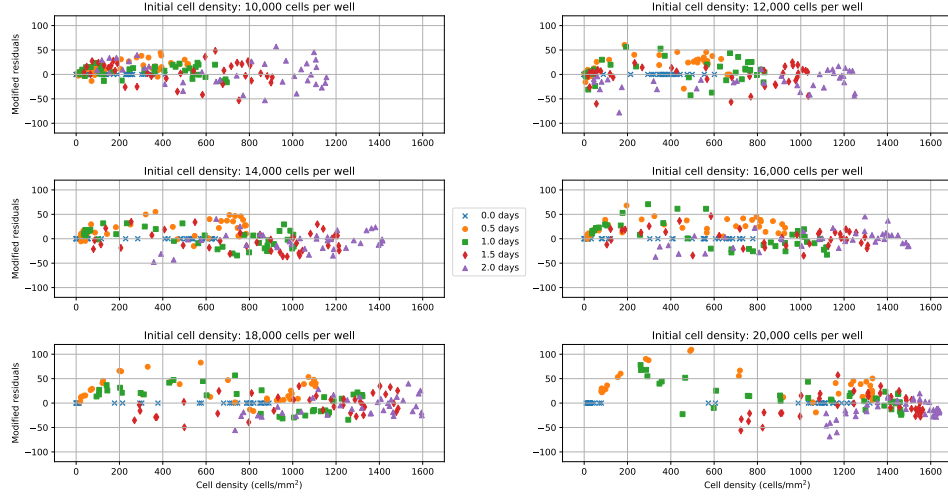
**Figure C.4** Modified residuals using BINNs with the governing reaction-diffusion PDE in Equation (4.15). Each subplot corresponds to an experiment with a different initial cell density.
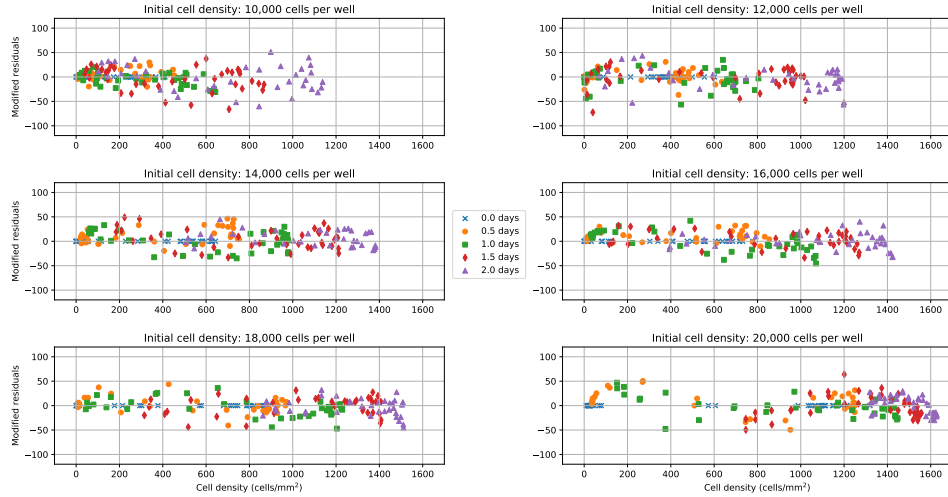


**Figure C.5** Modified residuals using BINNs with the governing delay-reaction-diffusion PDE in Equation (4.16). Each subplot corresponds to an experiment with a different initial cell density.
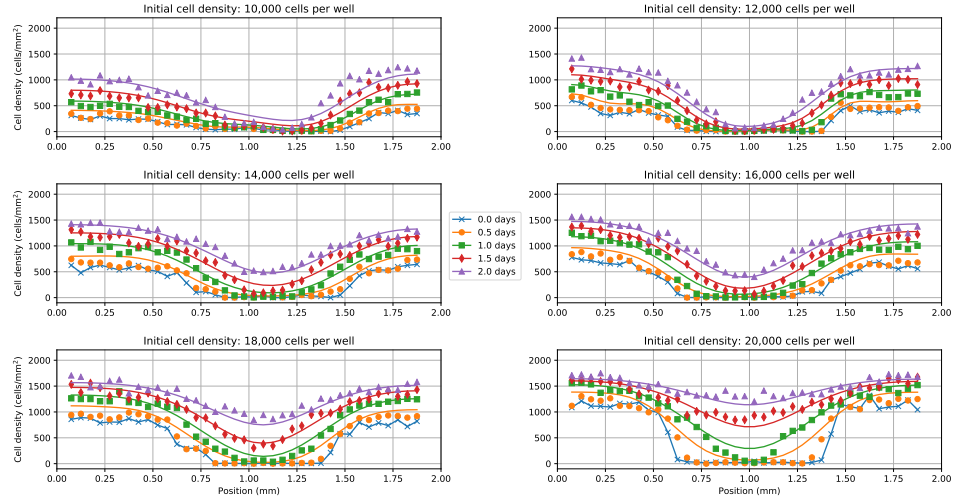
**Figure C.6** Predicted cell density profiles using the classical FKPP model in Equation (4.19). Each subplot corresponds to an experiment with a different initial cell density.
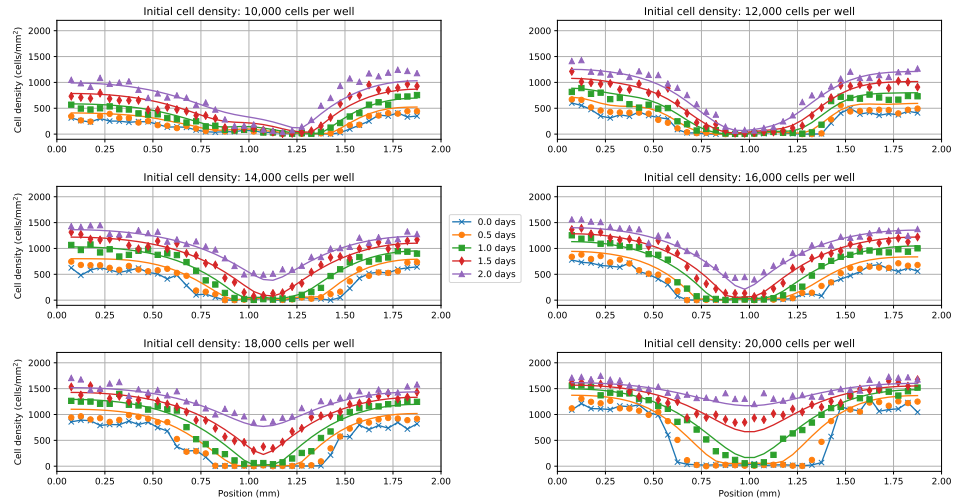


**Figure C.7** Predicted cell density profiles using the Generalized Porous-FKPP model in Equation (4.20). Each subplot corresponds to an experiment with a different initial cell density.

**Table C.1** Each column corresponds to an experiment with different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well).

| | Initial cell density | | | | | |
|---|---|---|---|---|---|---|
| Parameter | 10,000 | 12,000 | 14,000 | 16,000 | 18,000 | 20,000 |
| $D$ ($\mu$m$^2$/hr) | 309.7 | 253.8 | 681.8 | 540.9 | 735.7 | 978.5 |
| $r$ (1/hr) | 0.0437 | 0.0438 | 0.0483 | 0.0490 | 0.0540 | 0.0649 |

**Table C.2** Each column corresponds to an experiment with different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well).

| | Initial cell density | | | | | |
|---|---|---|---|---|---|---|
| Parameter | 10,000 | 12,000 | 14,000 | 16,000 | 18,000 | 20,000 |
| $D$ ($\mu$m$^2$/hr) | 1851.8 | 465.0 | 2993.4 | 2370.8 | 2377.8 | 2017.1 |
| $m$ (unitless) | 0.9704 | 0.3001 | 0.9923 | 0.9879 | 0.8265 | 0.6235 |
| $r$ (1/hr) | 0.0435 | 0.0436 | 0.0481 | 0.0484 | 0.0526 | 0.0639 |

**Figure C.8** The learned diffusivity $D_{\text{MLP}}$, growth $G_{\text{MLP}}$, and delay $T_{\text{MLP}}$ functions extracted from the corresponding BINNs with governing reaction-diffusion PDE in Equation (4.15) (first row) and delay-reaction-diffusion PDE in Equation (4.16) (second row). Each line corresponds to an experiment with a different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). Note that $D_{\text{MLP}}$ and $G_{\text{MLP}}$ have different lengths since they are evaluated between the minimum and maximum observed cell densities corresponding to each data set.
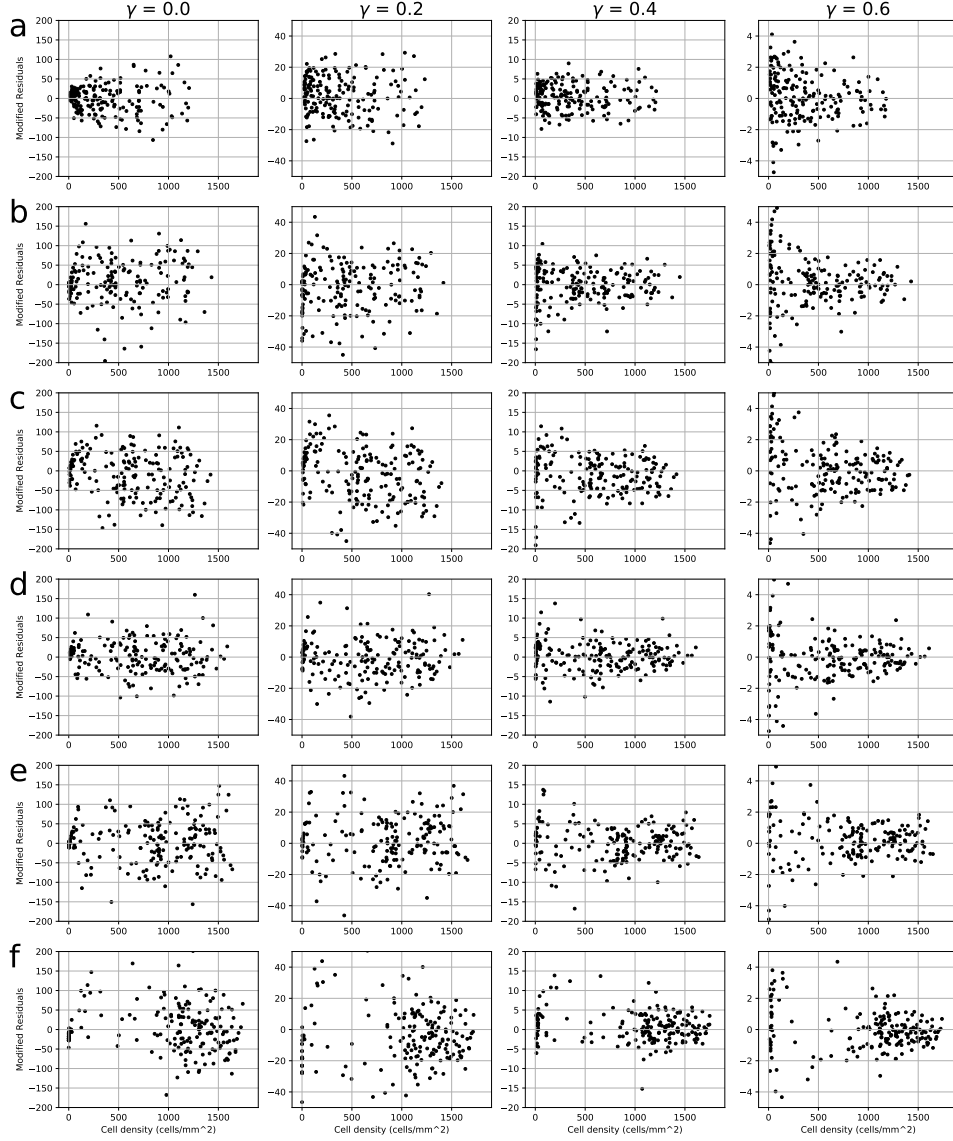
**Figure C.9** Statistical error model selection. The function-approximating deep neural network $u_{\mathrm{MLP}}$ is trained using $\mathscr{L}_{\mathrm{GLS}}$ for different values of $\gamma$ across each data set. Each subplot shows the modified residuals (see Equation (4.6)) as a function of the predicted cell density $u$. The columns correspond to different levels of proportionality (i.e. $\gamma = 0.0, 0.2, 0.4, 0.6$) where $\gamma = 0.0$ represents the constant variance (ordinary least squares) case. Each row (a-f) corresponds to an experiment with different initial cell density (i.e. 10,000, 12,000, 14,000, 16,000, 18,000, and 20,000 cells per well). The proportionality constant that results in the most i.i.d. residuals across each data set was chosen to calibrate the statistical error model in Equation (4.4).
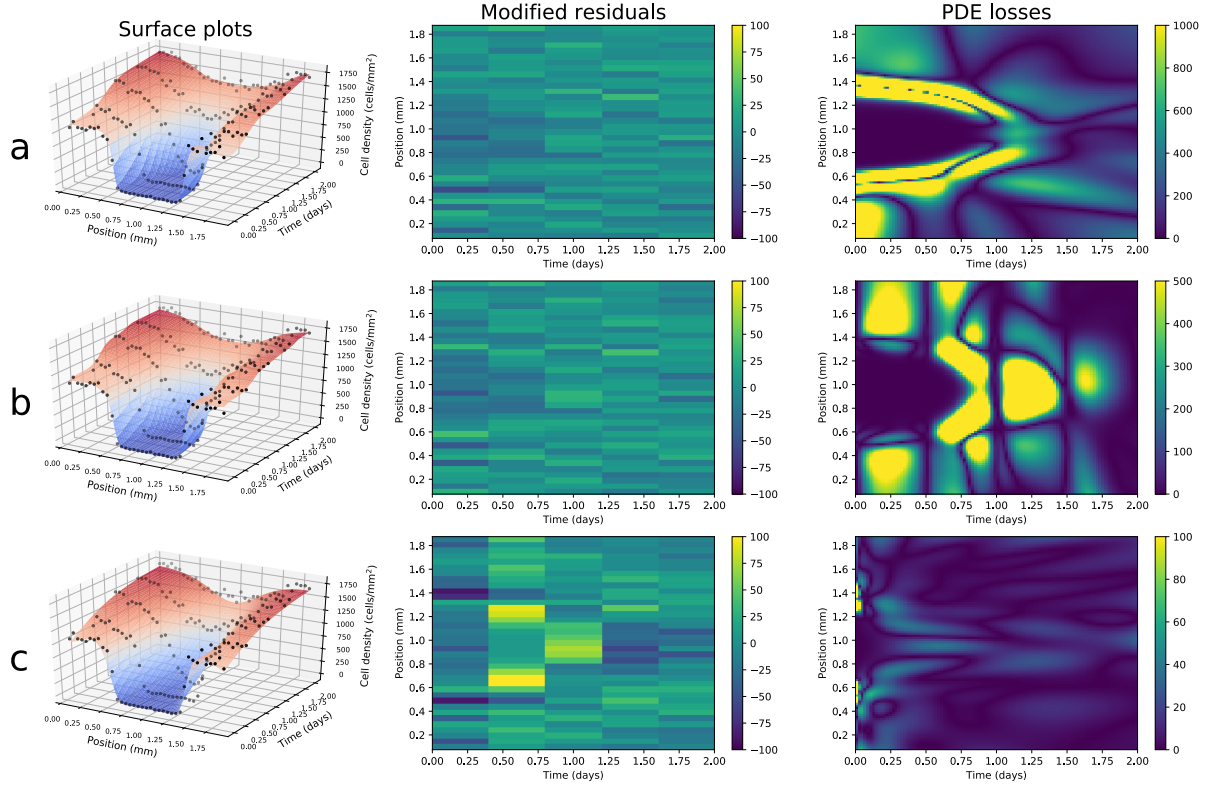
**Figure C.10** PDE random sampling validation. The BINNs framework is trained using $\mathscr{L}_{\text{Total}}$ with three ways of including the PDE error term $\mathscr{L}_{\text{Total}}$: (a) no PDE regularization, (b) PDE regularization at the data locations, and (c) PDE regularization at 10,000 randomly sampled points at each training iteration. The first column shows the scratch assay data with initial cell density 20,000 cells per well (black dots) with the corresponding BINNs approximation to the governing PDE $u_{\text{MLP}}$ (surface plot). The second column shows heatmaps of the modified residual errors (see Equation (4.6)) at each data point. The third column shows heatmaps of the PDE errors (see Equation (4.7)) evaluated on a $100 \times 100$ meshgrid over the input domain.