



Published in final edited form as:

*J Am Stat Assoc.* 2017 ; 112(520): 1697–1707. doi:10.1080/01621459.2016.1229197.

## Network Reconstruction From High-Dimensional Ordinary Differential Equations

Shizhe Chen<sup>a</sup>, Ali Shojaie<sup>b</sup>, and Daniela M. Witten<sup>b</sup>

<sup>a</sup>Department of Biostatistics, University of Washington, WA

<sup>b</sup>Departments of Biostatistics and Statistics, University of Washington, WA

### Abstract

We consider the task of learning a dynamical system from high-dimensional time-course data. For instance, we might wish to estimate a gene regulatory network from gene expression data measured at discrete time points. We model the dynamical system nonparametrically as a system of additive ordinary differential equations. Most existing methods for parameter estimation in ordinary differential equations estimate the derivatives from noisy observations. This is known to be challenging and inefficient. We propose a novel approach that does not involve derivative estimation. We show that the proposed method can consistently recover the true network structure even in high dimensions, and we demonstrate empirical improvement over competing approaches. Supplementary materials for this article are available online.

### Keywords

Additive model; Group lasso; High dimensionality; Ordinary differential equation; Variable selection consistency

### 1. Introduction

Ordinary differential equations (ODEs) have been widely used to model dynamical systems in many fields, including chemical engineering (Biegler, Damiano, and Blau 1986), genomics (Chou and Voit 2009), neuroscience (Izhikevich 2007), and infectious diseases (Wu 2005). A system of ODEs takes the form

---

**CONTACT:** Shizhe Chen, shizhe.chen@gmail.com; szchen@uw.edu Department of Biostatistics, University of Washington, Box 357232, Seattle, WA 98195-7232.

Color versions of one or more of the figures in the article can be found online at [www.tandfonline.com/r/JASA](http://www.tandfonline.com/r/JASA).

Supplementary materials for this article are available online. Please go to [www.tandfonline.com/r/JASA](http://www.tandfonline.com/r/JASA).

#### Supplementary Materials

The supplementary material contains proofs and details on data generation used in the main article.

$$X'(t; \theta) \equiv \begin{bmatrix} \frac{dX_1(t; \theta)}{dt} \\ \vdots \\ \frac{dX_p(t; \theta)}{dt} \end{bmatrix} = \begin{bmatrix} f_1(X(t; \theta), \theta) \\ \vdots \\ f_p(X(t; \theta), \theta) \end{bmatrix} \equiv f(X(t; \theta), \theta); \quad t \in [0, 1], \quad (1)$$

where  $X(t; \theta) = (X_1(t; \theta), \dots, X_p(t; \theta))^T$  denotes a set of variables, and the form of the functions  $f = (f_1, \dots, f_p)^T$  may be known or unknown. In (1),  $t$  indexes time. Typically, there is also an initial condition of the form  $X(0; \theta) = C$ , where  $C$  is a  $p$ -vector. In practice, the system (1) is often observed on discrete time points subject to measurement errors. Let  $Y_i \in \mathbb{R}^p$  be the measurement of the system at time  $t_i$  such that

$$Y_i = X(t_i; \theta^*) + \varepsilon_i, \quad i = 1, \dots, n, \quad (2)$$

where  $\theta^*$  denotes the true set of parameter values and the random  $p$ -vector  $\varepsilon_i$  represents independent measurement errors. In what follows, for notational simplicity, we sometimes suppress the dependence of  $X(t; \theta)$  on  $\theta$ , that is,  $X(t) \equiv X(t; \theta)$  in (1) and  $X^*(t) \equiv X(t; \theta^*)$  in (2).

In the context of high-dimensional time-course data arising from biology, it can be of interest to recover the structure of a system of ODEs—that is, to determine which features regulate each other. If  $f_j$  in (1) is a function of  $X_k$ , then we say that  $X_k$  *regulates*  $X_j$  in the sense that  $X_k$  controls the changes of  $X_j$  through its derivative  $X_j'$ . For instance, biologists might want to infer gene regulatory networks from noisy time-course gene expression data. In this case, the number of variables  $p$  exceeds the number of time points  $n$ ; we refer to this as the high-dimensional setting.

In high-dimensional statistics, sparsity-inducing penalties such as the lasso (Tibshirani 1996) and the group lasso (Yuan and Lin 2006) have been well-studied. Such penalties have also been extensively used to recover the structure of probabilistic graphical models (e.g., Yuan and Lin 2007; Friedman, Hastie, and Tibshirani 2008; Meinshausen and Bühlmann 2010; Voorman, Shojaie, and Witten 2014). However, model selection in high-dimensional ODEs remains a relatively open problem, with the exception of some notable recent work (Lu et al. 2011; Henderson and Michailidis 2014; Wu et al. 2014). In fact, the tasks of parameter estimation and model selection in ODEs from noisy data are very challenging, even in the classical statistical setting where  $n > p$  (see, e.g., Ramsay et al. 2007; Brunel 2008; Liang and Wu 2008; Qi and Zhao 2010; Xue, Miao, and Wu 2010; Gugushvili and Klaassen 2012; Hall and Ma 2014; Zhang, Cao, and Carroll 2015). Moreover, the problem of high-dimensionality is compounded if the form of the function  $f$  in (1) is unknown, leading to both statistical and computational issues.

In this article, we propose an efficient procedure for structure recovery of an ODE system of the form (1) from noisy observations of the form (2), in the setting where the functional

form of  $f$  is unknown. In Section 2, we review existing methods. In Section 3, we propose a new structure recovery procedure. In Section 4, we study the theoretical properties of our proposal. In Section 5, we apply our procedure to simulated data. In Section 6, we apply it to *in silico* gene expression data generated by GeneNetWeaver (Schaffter, Marbach, and Floreano 2011) and to calcium imaging data. We conclude with a discussion in Section 7. Proofs and additional details are provided in the supplementary material.

## 2. Literature Review

In this section, we review existing statistical methods for parameter estimation and/or model selection in ODEs. Most of the methods reviewed in this section are proposed for the low-dimensional setting. Even though they may not be directly applicable to the high-dimensional setting, they lay the foundation for the development of model selection procedures in high-dimensional additive ODEs.

### 2.1. Notation

Without loss of generality, assume that  $0 = t_1 < t_2 < \dots < t_n = 1$ . We let  $Y_{ij}$  indicate the observation of the  $j$ th variable at the  $i$ th time point,  $t_i$ . We use  $\mathcal{X}(h)$  to denote a nonparametric class of functions on  $[0, 1]$  indexed by some smoothing parameter(s)  $h$ . We use  $Z(\cdot)$  to represent an arbitrary function belonging to  $\mathcal{X}(\cdot)$ . We use  $\|\cdot\|_2$  to denote the  $\ell_2$ -norm of a vector or a matrix, and  $\|f\|$  to denote the  $\ell_2$ -norm of a function  $f$  on the interval  $[0, 1]$ , that is,  $\|f\|^2 \equiv \int_0^1 f^2(t) dt$ . We use an asterisk to denote true values—for instance,  $\theta^*$  denotes the true value of  $\theta$  in (1). We use  $\Lambda_{\min}(A)$  and  $\Lambda_{\max}(A)$  to denote the minimum and maximum eigenvalues of a square matrix  $A$ , respectively.

### 2.2. Methods that Assume a Known Form of $f$

**2.2.1. Gold Standard Approach**—To begin, we suppose that the function  $f$  in (1) takes a known form. Benson (1979) and Biegler, Damiano, and Blau (1986) proposed to estimate the unknown parameter  $\theta^*$  in (2) by solving the problem

$$\hat{\theta}^{\text{gold}} = \arg \min_{\theta} \sum_{i=1}^n \|Y_i - X(t_i; \theta)\|_2^2 \quad (3a)$$

$$\text{subject to } X'(t; \theta) = f(X(t; \theta), \theta), \quad t \in [0, 1]. \quad (3b)$$

Note that  $X(\cdot; \theta)$  in (3) is a fixed function given  $\theta$ , although an analytic expression may not be available. The resulting estimator  $\hat{\theta}^{\text{gold}}$  has appealing theoretical properties: for instance, when the measurement errors  $\epsilon_j$  in (2) are Gaussian, then  $\hat{\theta}^{\text{gold}}$  is the maximum likelihood estimator, and is  $\sqrt{n}$ -consistent. In this sense, (3) can thus be considered the *gold standard* approach. However, solving (3) is often computationally challenging.

**2.2.2. Two-Step Collocation Methods**—To overcome the computational challenges associated with solving (3), *collocation* methods have been employed by a number of authors (Varah 1982; Ellner, Seifu, and Smith 2002; Ramsay et al. 2007; Brunel 2008; Cao and Zhao 2008; Liang and Wu 2008; Cao, Wang, and Xu 2011; Lu et al. 2011; Gugushvili and Klaassen 2012; Brunel, Clairon, and d’Alché Buc 2014; Hall and Ma 2014; Henderson and Michailidis 2014; Wu et al. 2014; Dattner and Klaassen 2015; Zhang, Cao, and Carroll 2015).

The two-step collocation procedure first proposed by Varah (1982) involves fitting a smoothing estimate  $\hat{X}(\cdot; h)$  to the observations  $Y_1, \dots, Y_n$  in (2) with a smoothing parameter  $h$ , and then plugging  $\hat{X}(\cdot; h)$  and its derivative with respect to  $t$  into (1) to estimate  $\theta$ . This amounts to solving the optimization problem

$$\hat{\theta}^{\text{TS}} = \arg \min_{\theta} \int_0^1 \|\hat{X}'(t; h) - f(\hat{X}(t; h), \theta)\|_2^2 dt, \quad (4a)$$

where

$$\hat{X}(\cdot; h) = \arg \min_{Z(\cdot) \in \mathcal{X}(h)} \sum_{i=1}^n \|Y_i - Z(t_i)\|_2^2. \quad (4b)$$

The two-step procedure (4) has a clear advantage over the gold standard approach (3) because the former decouples the estimation of  $\theta$  and  $X$ . However, this advantage comes at a cost: due to the presence of  $\hat{X}'$  in (4a), the properties of the estimator  $\hat{\theta}^{\text{TS}}$  in (4) rely heavily on the smoothing estimates obtained in (4b), and  $\sqrt{n}$ -consistency has only been shown for certain values of the smoothing parameter  $h$  that are hard to choose in practice (Brunel 2008; Liang and Wu 2008; Gugushvili and Klaassen 2012).

Dattner and Klaassen (2015) proposed an improvement to (4) for a special case of (1). To be more specific, they assume that  $f_j(X(t), \theta)$  in (1) is a linear function of  $\theta$ , which leads to

$$X'(t) \equiv \begin{bmatrix} \frac{dX_1(t)}{dt} \\ \vdots \\ \frac{dX_p(t)}{dt} \end{bmatrix} = \begin{bmatrix} g_1^T(X(t))\theta \\ \vdots \\ g_p^T(X(t))\theta \end{bmatrix} \equiv g(X(t))\theta; \quad t \in [0, 1], \quad (5)$$

where  $g(X(t))$  is a known function of  $X(t)$ . Integrating both sides of (5) gives

$$X(t) = \left\{ \int_0^t g(X(u)) du \right\} \theta + C, \quad (6)$$

where  $C \equiv X(0; \theta)$ . The unknown parameter  $\theta^*$  is estimated by solving

$$\hat{\theta}^{\text{LM}} = \arg \min_{\theta} \int_0^1 \left\| \hat{X}(t; h) - \left\{ \int_0^t g(\hat{X}(u; h)) du \right\} \theta - C \right\|_2^2 dt, \quad (7a)$$

where

$$\hat{X}(\cdot; h) = \arg \min_{Z(\cdot) \in \mathcal{X}(h)} \sum_{i=1}^n \|Y_i - Z(t_i)\|_2^2. \quad (7b)$$

The optimization problem (7a) has an analytical solution, given the smoothing estimates from (7b). Compared with the two-step procedure (4), this approach requires an estimate of the integral,  $\int_0^t g(\hat{X}(u; h)) du$  in (7a), rather than an estimate of the derivative,  $\hat{X}'(t; h)$ . This has profound effects on the asymptotic behavior of the estimator  $\hat{\theta}^{\text{LM}}$ .  $\sqrt{n}$ -consistency of  $\hat{\theta}^{\text{LM}}$  has been established under mild conditions, and it has been found that the choice of smoothing parameter  $h$  is less crucial than for other methods (Gugushvili and Klaassen 2012).

Recently, Brunel, Clairon, and d'Alché Buc (2014) and Hall and Ma (2014) had considered alternatives to the loss function in (4a). Let  $C^1(0, 1)$  be the set of functions that are first-order differentiable on  $(0, 1)$  and equal zero on the boundary points 0 and 1. Then (1) implies that, for any  $\phi \in C^1(0, 1)$ ,

$$\int_0^1 f(X(t), \theta) \phi(t) dt + \int_0^1 X(t) \phi'(t) dt = 0. \quad (8)$$

Equation (8) is referred to as the *variational formulation* of the ODE. A least-square loss based on (8) takes the form

$$\hat{\theta}^{\text{V}} = \arg \min_{\theta} \frac{1}{L} \sum_{l=1}^L \left\| \int_0^1 f(\hat{X}(t; h), \theta) \phi_l(t) dt + \int_0^1 \hat{X}(t; h) \phi_l'(t) dt \right\|_2^2, \quad (9)$$

where  $\hat{X}(t; h)$  is defined in (4b) and  $\{\phi_l, l=1, \dots, L\}$  is a finite set of functions in  $C^1(0, 1)$  (Brunel, Clairon, and d'Alché Buc 2014). In Hall and Ma (2014), the loss function is the sum of the loss functions in (4b) and (9), so that  $\theta$  and the optimal bandwidth  $h$  are estimated simultaneously. It is immediately clear that the derivative  $X'(\cdot; \theta)$  is not needed in (9), which can lead to substantial improvement compared to the two-step procedure in (4). A minor drawback of (9) is that the variational formulation (8) is enforced on a finite set of functions  $\{\phi_l, l=1, \dots, L\}$  rather than on the whole class  $C^1(0, 1)$ . Under suitable

assumptions, the estimator  $\hat{\theta}^V$  is  $\sqrt{n}$ -consistent (Brunel, Clairon, and d'Alché Buc 2014; Hall and Ma 2014).

**2.2.3. The Generalized Profiling Method**—Another collocation-based method is the generalized profiling method of Ramsay et al. (2007). Instead of the smoothing estimate  $\hat{X}(\cdot; h)$  in (4b), the generalized profiling method uses a smoothing estimate  $\check{X}(\cdot; h, \theta)$  that minimizes the weighted sum of a data-fitting loss and a model-fitting loss for any given  $\theta$ . In greater detail,

$$\hat{\theta}_\lambda^{\text{GP}} = \arg \min_{\theta} \sum_{i=1}^n \|Y_i - \check{X}(t_i; h, \theta)\|_2^2, \quad (10a)$$

where

$$\check{X}(\cdot; h, \theta) = \arg \min_{Z(\cdot) \in \mathcal{X}(h)} \frac{1}{n} \sum_{i=1}^n \|Y_i - Z(t_i)\|_2^2 + \lambda \int_0^1 \|Z'(t) - f(Z(t), \theta)\|_2^2 dt. \quad (10b)$$

In Ramsay et al. (2007), the authors solve (10a) iteratively for a nondecreasing sequence of  $\lambda$ 's in (10b).  $\sqrt{n}$ -consistency of the limiting estimator was later established by Qi and Zhao (2010). Zhang, Cao, and Carroll (2015) proposed a model selection procedure by applying an ad hoc lasso procedure (Wang and Leng 2007) to the estimates from (10).

### 2.3. Methods that do not Assume the Form of $f$

A few authors have recently considered modeling large-scale dynamical systems from biology using ODEs (Henderson and Michailidis 2014; Wu et al. 2014), under the assumption that the right-hand side of (1) is additive,

$$X_j'(t) = \theta_{j0} + \sum_{k=1}^p f_{jk}(X_k(t)), \quad \theta_{j0} \in \mathbb{R}. \quad (11)$$

Henderson and Michailidis (2014) and Wu et al. (2014) approximated the unknown  $f_{jk}$  with a truncated basis expansion. Consider a finite basis,  $\psi(x) = (\psi_1(x), \dots, \psi_M(x))^T$ , such that

$$f_{jk}(a_k) = \psi(a_k)^T \theta_{jk} + \delta_{jk}(a_k), \quad \theta_{jk} \in \mathbb{R}^M, \quad (12)$$

where  $\delta_{jk}(a_k)$  denotes the residual. Using (12), a system of additive ODEs of the form (11) can be written as

$$X'_j(t) = \theta_{j0} + \sum_{k=1}^p \psi(X_k(t))^T \theta_{jk} + \sum_{k=1}^p \delta_{jk}(X_k(t)), \quad j = 1, \dots, p. \quad (13)$$

Henderson and Michailidis (2014) and Wu et al. (2014) considered the problem of estimating and selecting the nonzero elements  $\theta_{jk}$  in (13). Roughly speaking, they proposed to solve optimization problems of the form

$$\begin{aligned} \hat{\theta}_j^{\text{NP}} = & \arg \min_{\theta_{j0} \in \mathbb{R}, \theta_{jk} \in \mathbb{R}^M} \int_0^1 \left\| \hat{X}'_j(t; h) - \theta_{j0} - \sum_{k=1}^p \psi(\hat{X}_k(t; h))^T \theta_{jk} \right\|_2^2 dt \quad (14a) \\ & + \lambda_n \sum_{k=1}^p \left[ \int_0^1 \{ \psi(\hat{X}_k(t; h))^T \theta_{jk} \}^2 dt \right]^{1/2}, \end{aligned}$$

for  $j = 1, \dots, p$ , where

$$\hat{X}(\cdot; h) = \arg \min_{Z(\cdot) \in \mathcal{X}(h)} \sum_{i=1}^n \|Y_i - Z(t_i)\|_2^2. \quad (14b)$$

In (14a), a standardized group lasso penalty forces all elements in  $\theta_{jk}$  to be either zero or nonzero when  $\lambda_n$  is large, thereby providing variable selection.

The proposals by Henderson and Michailidis (2014) and Wu et al. (2014) are slightly more involved than (14): an extra  $\ell_2$ -penalty is applied to the  $\theta_{jk}$ 's in (14a) in Henderson and Michailidis (2014), whereas in Wu et al. (2014) (14a) is followed by tuning parameter selection using Bayesian information criterion (BIC), an adaptive group lasso regression, and a regular lasso. We refer the reader to Henderson and Michailidis (2014) and Wu et al. (2014) for further details.

### 3. Proposed Approach

We consider the problem of model selection in high-dimensional ODEs. As in Henderson and Michailidis (2014) and Wu et al. (2014), we assume an additive ODE model (11). We use a finite basis  $\psi(\cdot)$  to approximate the additive components  $f_{jk}$  as in (12), leading to an ODE system that is linear in the unknown parameters (13). Following the example by Dattner and Klaassen (2015), we exploit this linearity by integrating both sides of (13), which yields

$$X_j(t) = X_j(0) + \theta_{j0}t + \sum_{k=1}^p \Psi_k(t)^T \theta_{jk} + \sum_{k=1}^p \int_0^t \delta_{jk}(X_k(u)) du, \quad (15)$$

where  $\Psi_k(t)$  denotes the integrated basis such that

$$\begin{aligned}\Psi_k(t) &= (\Psi_{k1}(t), \dots, \Psi_{kM}(t))^T \\ &= \int_0^t \psi(X_k(u)) du, k = 1, \dots, p,\end{aligned}\quad (16)$$

and  $\Psi_0(t) = t$ . Our method, called *graph reconstruction via additive differential equations* (GRADE), then solves the following problem for  $j = 1, \dots, p$ :

$$\begin{aligned}\hat{\theta}_j &= \arg \min_{C_{j0} \in \mathbb{R}, \theta_{j0} \in \mathbb{R}, \theta_{j1}, \dots, \theta_{jp} \in \mathbb{R}^M} \frac{1}{2n} \\ &\times \sum_{i=1}^n \left\{ Y_{ij} - C_{j0} - \theta_{j0} \hat{\Psi}_0(t_i) - \sum_{k=1}^p \theta_{jk}^T \hat{\Psi}_k(t_i) \right\}^2 \\ &+ \lambda_{n,j} \sum_{k=1}^p \left[ \frac{1}{n} \sum_{i=1}^n \{ \theta_{jk}^T \hat{\Psi}_k(t_i) \}^2 \right]^{1/2},\end{aligned}\quad (17a)$$

where

$$\hat{X}(\cdot; h) = \arg \min_{Z(\cdot) \in \mathcal{X}(h)} \sum_{i=1}^n \|Y_i - Z(t_i)\|_2^2, \quad (17b)$$

and

$$\hat{\Psi}_0(t) = t; \hat{\Psi}_k(t) = \int_0^t \psi(\hat{X}_k(u; h)) du, k = 1, \dots, p. \quad (17c)$$

In (17a),  $\lambda_{n,j}$  is a nonnegative sparsity-inducing tuning parameter. We may sometimes use  $\lambda_{n,j} \equiv \lambda_n$  for  $j = 1, \dots, p$  for simplicity. If the true function  $f_{jk}^*$  in (11) is nonzero, we say that the  $k$ th variable  $X_k^*$  is a true regulator of  $X_j^*$ . We let  $S_j \equiv \{k: \|f_{jk}^*\|_2 \neq 0, k = 1, \dots, p\}$  denote the set of true regulators. We let the estimated index set of regulators be  $\hat{S}_j \equiv \{k: \|\hat{\theta}_{jk}\|_2 \neq 0, k = 1, \dots, p\}$ . We then reconstruct the network using  $\hat{S}_j, j = 1, \dots, p$ .

Both (17a) and (17b) can be implemented efficiently using existing software (see, e.g., software methods in Loader 1999; Meier, van de Geer, and Bühlmann 2008). In our theoretical analysis in Section 4, we use local polynomial regression to obtain the smoothing estimate in (17b). We use generalized cross-validation (GCV) on the loss (17b) to select the smoothing tuning parameter  $h$ . We use BIC to select the number of bases  $M$  for  $\psi$  and  $\hat{\Psi}$  in (17c), and the sparsity tuning parameter  $\lambda_n$  in (17a).



In some studies, time-course data are collected from multiple samples, or experiments. Let  $R$  denote the total number of experiments, and  $Y^{(r)}$  the observations in the  $r$ th experiment. We assume that the same ODE system (13) applies across all experiments with the same true parameter  $\theta_{jk}^*$ . We allow a different set of initial values for each experiment. Assume that each experiment consists of measurements on the same set of time points. This leads us to modify (17) as follows:

$$\begin{aligned} \hat{\theta}_j = & \arg \min_{C_{j0}^{(r)} \in \mathbb{R}, \theta_{j0} \in \mathbb{R}, \theta_{j1}, \dots, \theta_{jp} \in \mathbb{R}^M} \frac{1}{2Rn} \sum_{r=1}^R \sum_{i=1}^n \quad (18) \\ & \times \left\{ Y_{ij}^{(r)} - C_{j0}^{(r)} - \theta_{j0} \widehat{\Psi}_0(t_i) - \sum_{k=1}^p \theta_{jk}^T \widehat{\Psi}_k^{(r)}(t_i) \right\}^2 \\ & + \lambda_n \sum_{k=1}^p \left[ \frac{1}{Rn} \sum_{r=1}^R \sum_{i=1}^n \{ \theta_{jk}^T \widehat{\Psi}_k^{(r)}(t_i) \}^2 \right]^{1/2}, \end{aligned}$$

where

$$\begin{aligned} \widehat{X}^{(r)}(\cdot; h) &= \arg \min_{Z(\cdot) \in \mathcal{X}(h)} \sum_{i=1}^n \|Y_i^{(r)} - Z(t_i)\|_2^2, r = 1, \dots, R, \\ \widehat{\Psi}_0(t) = t; \widehat{\Psi}_k^{(r)}(t) &= \int_0^t \psi(\widehat{X}_k^{(r)}(u; h)) du, k = 1, \dots, p. \end{aligned}$$

In Sections 4, 5.1, and 5.2, we will assume that only one experiment is available, so that our proposal takes the form (17). In Sections 5.3 and 6, we will apply our proposal to data from multiple experiments using (18).

### Remark 1

To facilitate the comparison of GRADE (17) with other methods, we introduce an intermediate variable,

$$\widetilde{X}_j(t; h, \theta) \equiv C_{j0} + \theta_{j0}t + \sum_{k=1}^p \theta_{jk}^T \widehat{\Psi}_k(t), \quad (19)$$

following from (15). Plugging (19) into the loss function in (17a) yields

$\sum_{i=1}^n \{Y_{ij} - \widetilde{X}_j(t_i; h, \theta)\}^2$ . In the gold standard (3), the ODE system (1) is strictly satisfied due to the constraint in (3b). In the two-step procedure (4a) and (14a), the smoothing estimate  $\widehat{X}(\cdot; h)$  does not satisfy (1). GRADE stands in between: the initial estimate  $\widehat{X}(\cdot; h)$  in (17b) is solely based on the observations, while the intermediate estimate  $\widetilde{X}(\cdot; h, \theta)$  is calculated by plugging  $\widehat{X}(\cdot; h)$  into the additive ODE (13).

## 4. Theoretical Properties

In this section, we establish variable selection consistency of the GRADE estimator (17). Technical proofs of the statements in this section are available in Section A in the supplementary material. We use  $s_j$  to denote the cardinality of  $S_j$ , and set  $s = \max_j \{s_j\}$ . For ease of presentation, we let  $S_j^0 = \{0\} \cup S_j$ , so that  $\Psi_{S_j^0}(t) = (\Psi_0(t), \Psi_{S_j}^T(t))^T = (t, \Psi_{S_j}^T(t))^T$  is an  $(s_j M + 1)$ -vector.

The proposed method (17) differs from the standard sparse additive model (Ravikumar et al. 2009) in that the regressors  $\widehat{\Psi}_k(t)$  in (17c) are estimated from smoothing estimates  $\widehat{X}(\cdot; h)$  (17b) instead of the true trajectories  $X^*$  in (2). We use local polynomial regression to compute  $\widehat{X}(\cdot; h)$  in (17b) (see, e.g., eq. (1.67) of Tsybakov 2009 for details on parameterization). To establish variable selection consistency, it is necessary to obtain a bound for the difference between  $\widehat{X}(\cdot; h)$  and  $X^*$ . This is addressed in Theorem 1. Using the bound in Theorem 1, we then establish variable selection consistency of the estimator in (17) for high-dimensional ODEs in Theorem 2.

In this study, we assume that the measurement errors in (2) are normally distributed. Generalizations to bounded or sub-Gaussian errors are straightforward.

### Assumption 1

The measurement errors in (2) are independent, and  $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$ ,  $i = 1, \dots, n, j = 1, \dots, p$ .

We also require the true trajectories  $X_j^*$  in (2) to be smooth.

### Assumption 2

Assume that the solutions  $X_j^*$ ,  $1 \leq j \leq p$ , belong to a Hölder class  $\Sigma(\beta_1, L_1)$ , where  $\beta_1 \geq 3$ .

In addition, we need some regularity assumptions to hold for the smoothing estimation (17b). These assumptions are common and not crucial to this study, and are hence deferred to Section A.2 in the supplementary material (or see sec. 1.6.1 in Tsybakov 2009). We arrive at the following concentration inequality for  $\|\widehat{X} - X^*\|$ .

### Theorem 1

Suppose that Assumptions 1–2 and S1–S3 in the supplementary material are satisfied. Let  $\widehat{X}_j$  in (17b) be the local polynomial regression estimator of order  $\lfloor \beta_1 \rfloor$  with bandwidth

$$h_n \propto n^{-(\alpha - 1)/(2\beta_1 + 1)} \quad (20)$$

for some positive  $\alpha < 1$ . Then, for each  $j = 1, \dots, p$ ,

$$\|\hat{X}_j - X_j^*\|^2 \leq C_2 n^{\frac{2\beta_1}{2\beta_1+1}(\alpha-1)} \quad (21)$$

holds with probability at least  $1 - 2 \exp\{-n^\alpha/(2C_3\sigma^2)\}$ , for some constants  $C_2$  and  $C_3$ .

The concentration inequality in Theorem 1 is derived using concentration bounds for Gaussian errors (Boucheron, Lugosi, and Massart 2013). Using Theorem 1, we see that the bound (21) holds uniformly for  $j = 1, \dots, p$  with probability at least  $1 - 2p \exp\{-n^\alpha/(2C_3\sigma^2)\}$ . The bound in Theorem 1 thus holds uniformly for  $j = 1, \dots, p$  with probability converging to unity if  $p = o(\exp\{n^\alpha/(2C_3\sigma^2)\})$ .

For the methods outlined in (14) (Henderson and Michailidis 2014; Wu et al. 2014), variable selection consistency depends on the convergence of  $\|\hat{X}' - (X^*)'\|$  and  $\|\hat{X} - X^*\|$ . In contrast, our method depends only on the convergence rate of  $\|\hat{X} - X^*\|$ . It is known that the convergence of  $\|\hat{X}' - (X^*)'\|$  is slower than that of  $\|\hat{X} - X^*\|$ , see, for example, Gugushvili and Klaassen (2012). As a result, the rate of convergence of  $\hat{\theta}_{jk}$  from (14) is slower than that of our proposed method (17).

To establish the main result, we need the following additional assumptions. Recall the definition of  $\Psi_j(t)$  from (16); for convenience, we suppress the dependence of  $\Psi(t)$  on  $t$  in what follows.

### Assumption 3

For  $j = 1, \dots, p$ ,  $(X_j^*)'$  is an additive function of  $X_k^*$ ,  $k = 1, \dots, p$ . In other words,

$$(X_j^*)'(t) = \theta_{j0}^* + \sum_{k=1}^p f_{jk}^*(X_k^*(t)), \quad \theta_{j0}^* \in \mathbb{R}, \quad j = 1, \dots, p, \quad (22)$$

where  $\int_0^1 f_{jk}^*(X_k^*(t))dt = 0$  for all  $j, k$ . Furthermore, the functions  $f_{jk}^*$  ( $1 \leq j, k \leq p$ ) belong to a Sobolev class  $W(\beta_2, L_2)$  on a finite interval with  $\beta_2 \geq 3$ .

### Assumption 4

The eigenvalues of  $\int_0^1 \Psi_{S_j^0} \Psi_{S_j^0}^T dt$  are bounded from above by  $C_{\max}$  and bounded from below by a positive number  $C_{\min}$ , and for  $k \notin S_j^0$ , the eigenvalues of  $\int_0^1 \Psi_k \Psi_k^T dt$  are bounded from below by  $C_{\min}$ . In other words,

$$\begin{aligned}
 0 < C_{\min} &\leq \Lambda_{\min} \left( \int_0^1 \Psi_{S_j^0} \Psi_{S_j^0}^T dt \right) \quad (23) \\
 &\leq \Lambda_{\max} \left( \int_0^1 \Psi_{S_j^0} \Psi_{S_j^0}^T dt \right) \leq C_{\max},
 \end{aligned}$$

and

$$C_{\min} \leq \Lambda_{\min} \left( \int_0^1 \Psi_k \Psi_k^T dt \right), \text{ for } k \notin S_j^0. \quad (24)$$

**Assumption 5**

Assume that

$$\max_{k \notin S_j^0} \left\| \left( \int_0^1 \Psi_k \Psi_{S_j^0}^T dt \right) \left( \int_0^1 \Psi_{S_j^0} \Psi_{S_j^0}^T dt \right)^{-1} \right\|_2 \leq \xi. \quad (25)$$

The first part of Assumption 4 ensures identifiability among the  $s_j + 1$  elements in the set  $\{t, X_{S_j}^*\}$ , and the second part ensures that  $\Psi_k$  is nondegenerate for  $k \notin S_j^0$ . Assumption 5 restricts the association between the elements in the set  $\{t, X_{S_j}^*\}$  and the elements in the set  $X_{S_j^c}^*$ . Note that in order for the parameters in an additive model such as (13) to be

identifiable, there must be no concurrency among the variables (Buja, Hastie, and Tibshirani 1989). This is guaranteed by Assumptions 4 and 5, which appear often in the literature of lasso regression (Meinshausen and Bühlmann 2006; Zhao and Yu 2006; Ravikumar et al. 2009; Wainwright 2009; Lee, Sun, and Taylor 2013). We refer the readers to Miao et al. (2011) for a detailed discussion of the identifiability of the parameters in an ODE model.

The next assumption characterizes the relationships between the quantities in Assumptions 4 and 5 and the sparsity tuning parameter  $\lambda_n$  in (17a). Similar assumptions have been made in lasso-type regression (Meinshausen and Bühlmann 2006; Zhao and Yu 2006; Ravikumar et al. 2009; Wainwright 2009; Lee, Sun, and Taylor 2013).

**Assumption 6**

Assume that

$$f_{\min} > \lambda_n \frac{4\sqrt{2sC_{\max}}}{C_{\min}} \text{ and } \xi < \frac{1}{4} \sqrt{\frac{C_{\min}}{sC_{\max}}},$$

where  $f_{\min} \equiv \min_{k \in S_j} \left\{ \int_0^1 [f_{jk}^*(X_k^*(t))]^2 dt \right\}^{1/2}$  is the minimum regulatory effect.

Furthermore, we impose some regularity conditions on the bases  $\psi(\cdot)$ ; these are deferred to Assumption S4 in the supplementary material.

We arrive at the following theorem.

### Theorem 2

Suppose that Assumptions 1–6 and S1–S4 in the supplementary material hold, and let

$$\begin{aligned} h_n &\propto n^{(\alpha-1)/(2\beta_1+1)}, \quad M \propto n^{\frac{2\beta_1(1-\alpha)}{(2\beta_1+1)(2\beta_2+1)}}, \\ \lambda_n &\propto n^{-\frac{\beta_1(2\beta_2-1)(1-\alpha)}{(2\beta_1+1)(2\beta_2+1)} + 2\gamma}, \end{aligned}$$

where  $0 < \alpha < 1$ ,  $0 < \gamma < H_1(\beta_1, \beta_2, \alpha)$ , and  $H_1(\beta_1, \beta_2, \alpha)$  is a constant that depends only on  $\beta_1$ ,  $\beta_2$  and  $\alpha$ . Then as  $n$  increases, the proposed procedure (17) correctly recovers the true graph, that is,  $\hat{S}_j = S_j$  for all  $j = 1, \dots, p$ , with probability converging to 1, if  $s = O(n^\gamma)$  and  $pn \exp(-C_4 n^\alpha / \sigma^2) = o(1)$  for some constant  $C_4$ .

Because the regressors  $\hat{\Psi}$  are estimated, establishing variable selection consistency requires extra attention. To prove Theorem 2, we must first establish variable selection consistency of group lasso regression with errors in variables. This generalizes the recent work on errors in variables for lasso regression (Loh and Wainwright 2012). Theorem 2 ensures that the proposed method is able to recover the true graph exactly, given sufficiently dense observations in a finite time interval if the graph is sparse. The number of variables in the system can grow exponentially fast with respect to  $n$ , which means that the result holds for the “large  $p$ , small  $n$ ” scenario.

Theorem 2 does not provide us with practical guidance for selecting the bandwidth  $h_n$  for the local polynomial regression estimator  $\hat{X}_j$ . The next result mirrors Theorem 2 for the

bandwidths selected by cross-validation or GCV, which converge to  $h_n \propto n^{-1/(2\beta_1+1)}$  asymptotically (see Xia and Li 2002; Tsybakov 2009 for details).

### Proposition 1

Suppose that Assumptions 1–6 and S1–S4 in the supplementary material hold, and let

$$\begin{aligned} h_n &\propto n^{-1/(2\beta_1+1)}, \quad M \propto n^{\frac{1}{2\beta_2+1} \left( \frac{2\beta_1}{2\beta_1+1} - \alpha \right)}, \quad \text{and} \\ \lambda_n &\propto n^{-\frac{2\beta_2-1}{4\beta_2+2} \left( \frac{2\beta_1}{2\beta_1+1} - \alpha \right) + 2\gamma}, \end{aligned}$$

where  $0 < \alpha < \frac{2\beta_1}{2\beta_1+1}$ ,  $0 < \gamma < H_2(\beta_1, \beta_2, \alpha)$ , and  $H_2(\beta_1, \beta_2, \alpha)$  is a constant that depends only on  $\beta_1, \beta_2$ , and  $\alpha$ . Then as  $n$  increases, the proposed procedure (17) correctly recovers the true graph, that is,  $\hat{S}_j = S_j$  for all  $j = 1, \dots, p$ , with probability converging to 1, if  $s = O(n^\gamma)$  and  $pn \exp(-C_4 n^\alpha / \sigma^2) = o(1)$  for some constant  $C_4$ .

We note that selecting the values of  $M$  and  $\lambda_n$  that yield the rate specified in Proposition 1 is challenging in practice. The rate of convergence of the sparsity tuning parameter  $\lambda_n$  is slower in Proposition 1 compared to Theorem 2. This results in an increase in the minimum regulatory effect  $f_{\min}$  because of the relation between  $f_{\min}$  and  $\lambda_n$  in Assumption 6.

## 5. Numerical Experiments

We study the empirical performance of our proposal in three different scenarios in the following subsections. In what follows, given a set of initial conditions and a system of ODEs, numerical solutions of the ODEs are obtained using the Euler method with step size 0.001. Observations are drawn from the solutions at an evenly spaced time grid  $\{iT/n; i = 1, \dots, n\}$  with independent  $N(0, 1)$  measurement errors, unless specified otherwise. To facilitate the comparison of GRADE with other methods, we fit the smoothing estimates  $\hat{X}$  in (17b) using smoothing splines with bandwidth chosen by GCV. We use cubic splines with two internal knots as the basis functions in (17c) in Sections 5.1 and 5.3. Linear basis functions are used in Section 5.2. The integral  $\hat{\Psi}_k(t) = \int_0^t \psi(\hat{X}_k(u; h)) du$  in (17c) is calculated numerically with step size 0.01.

### 5.1. Variable Selection in Additive ODEs

In this simulation, we compare GRADE with NeRDS (Henderson and Michailidis 2014) and SA-ODE (Wu et al. 2014) described in (14). We consider the following system of additive ODEs, for  $k = 1, \dots, 5$ :

$$\begin{cases} X'_{2k-1}(t) = \theta_{2k-1,0} + \psi(X_{2k-1}(t))^T \theta_{2k-1,2k-1} + \psi(X_{2k}(t))^T \theta_{2k-1,2k} \\ X'_{2k}(t) = \theta_{2k,0} + \psi(X_{2k-1}(t))^T \theta_{2k,2k-1} + \psi(X_{2k}(t))^T \theta_{2k,2k} \end{cases}, t \in [0, 20] \quad (26)$$

where  $\psi(x) = (x, x^2, x^3)^T$  is the cubic monomial basis. The parameters and initial conditions are chosen so that the solution trajectories are identifiable under an additive model (Buja, Hastie, and Tibshirani 1989). Detailed specification of (26) can be found in Section C of the supplementary material.

After generating data according to (26) and introducing noise, we apply GRADE, NeRDS, and SA-ODE to recover the directed graph encoded in (26). Both NeRDS and SA-ODE are implemented using code provided by the authors. NeRDS and SA-ODE use smoothing splines to estimate  $\hat{X}$  and  $\hat{X}'$  in (14b), and cubic splines with two internal knots as the basis  $\psi$  in (14a). As mentioned briefly in Section 2, NeRDS applies an additional smoothing penalty that amounts to an  $\ell_2$  penalty on  $\theta_{jk}$  in (14a), controlled by a parameter selected

using GCV (Henderson and Michailidis 2014). We apply GRADE using the same smoothing estimates and basis functions as NeRDS and SA-ODE. To facilitate a direct comparison to NeRDS, we apply GRADE both with and without an additional  $\ell_2$ -type penalty on the  $\theta_{jk}$ 's in (17a). We apply all methods for a range of values of the sparsity-inducing tuning parameter (e.g.,  $\lambda_n$  in (17a)), to yield a recovery curve of varying sparsity.

We summarize the simulation results in Figure 1, where the numbers of true edges selected are displayed against the total numbers of selected edges over a range of sparsity tuning parameters. We see that GRADE outperforms the other two methods, which corroborates our theoretical findings in Section 4 that our proposed method is more efficient than methods such as NeRDS and SA-ODE, which involve derivative estimation (see, e.g., comments below Theorem 1).

## 5.2. Variable Selection in Linear ODEs

In this simulation, we compare GRADE to two recent proposals by Brunel, Clairon, and d'Alché Buc (2014) and Hall and Ma (2014). Recall from Section 2.2.2 that Brunel, Clairon, and d'Alché Buc (2014) and Hall and Ma (2014) proposed to estimate a few unknown parameters in an ODE system of known form. Hence, we consider a simple linear ODE system, for  $k = 1, \dots, 4$ ,

$$\begin{cases} X'_{2k-1}(t) = 2k\pi X_{2k}(t) \\ X'_{2k}(t) = -2k\pi X_{2k-1}(t) \end{cases}, t \in [0, 1]. \quad (27)$$

For each  $k = 1, \dots, 4$ , we set the initial condition to be  $(X_{2k-1}(0), X_{2k}(0)) = (\sin(y_k), \cos(y_k))$  where  $y_k \sim \mathcal{N}(0, 1)$ . The solutions to (27) take the form of sine and cosine functions of frequencies ranging from  $2\pi$  to  $8\pi$ . The graph corresponding to (27) is sparse, with only eight directed edges out of 64 possible edges. We fit the model

$$X'(t) = \Theta X(t) + C, \quad (28)$$

where  $\Theta$  is an unknown  $8 \times 8$  matrix and  $C$  is an 8-vector. We apply the method in Brunel, Clairon, and d'Alché Buc (2014) using the code provided by the authors. We implement the method in Hall and Ma (2014) in R based on the authors' code in Fortran. Because the loss function in Hall and Ma (2014) is not convex, we use five sets of random initial values and report the best performance. Since both Brunel, Clairon, and d'Alché Buc (2014) and Hall and Ma (2014) yielded dense estimates for  $\Theta$  in (28), to examine how well these methods recover the true graph, we threshold the estimates at a range of values to obtain a variable selection path. We apply GRADE using the linear basis function  $\psi(x) = x$ .

Results are shown in Figure 2. We can see that GRADE outperforms the methods in Brunel, Clairon, and d'Alché Buc (2014) and Hall and Ma (2014). This is likely because GRADE exploits the sparsity of the true graph with a sparsity-inducing penalty. In principle, Brunel, Clairon, and d'Alché Buc (2014) and Hall and Ma (2014) could be generalized to include penalties on the parameters. We leave this to future research.

### 5.3. Robustness of GRADE to the Additivity Assumption

The GRADE method assumes that the true underlying model is additive (Assumption 3). However, in many systems, the additivity assumption is violated; for instance, multiplicative effects may be present in gene regulatory networks (Ma et al. 2009). In this subsection, we investigate the performance of GRADE in a setting where the true model is nonadditive. We consider the following system of ODEs, for  $k = 1, \dots, 5$ ,

$$\begin{cases} X'_{2k-1}(t) = f_{2k-1}(X_{2k-1}(t), X_{2k}(t)) \equiv 2X_{2k-1}(t) - vX_{2k-1}(t)X_{2k}(t) \\ X'_{2k}(t) = f_{2k}(X_{2k-1}(t), X_{2k}(t)) \equiv vX_{2k-1}(t)X_{2k}(t) - 2X_{2k}(t) \end{cases}, t \in [0, 5], \quad (29)$$

where  $v$  is a positive constant. For each  $k = 1, \dots, 5$ , the pair of Equation (29) is a special case of the Lotka–Volterra equations (Volterra 1928), which represent the dynamics between predators ( $X_{2k}$ ) and prey ( $X_{2k-1}$ ). The parameter  $v$  defines the interaction between the two populations. For  $v = 0$ , both  $X'_{2k-1}$  and  $X'_{2k}$  are nonadditive functions of  $X_{2k-1}$  and  $X_{2k}$ . We define two types of directed edges, where  $\mathcal{E}_1 \equiv \{(X_j, X_j), j = 1, \dots, 10\}$  and  $\mathcal{E}_2 \equiv \{(X_{2k-1}, X_{2k}), (X_{2k}, X_{2k-1}), k = 1, \dots, 5\}$  represent the self-edges and nonself-edges, respectively. Figure 3(a) contains an illustration of the graph and edge types for each pair of equations. In what follows, we investigate how well GRADE recovers these two types of edges as we change the parameter  $v$ , that is, as the additivity assumption is violated.

Since measurement error is not essential to the current discussion, we generate data according to (29) without adding noise. To ensure that the trajectories are identifiable, we generate  $R = 2$  sets of random initial values drawn from  $\mathcal{N}_{10}(0, 2I_{10})$ , where  $I_{10}$  is a  $10 \times 10$  identity matrix. To quantify the amount of signal in an edge that GRADE can detect, we introduce the quantity

$$D_{j,k}(v) = \mathbb{E} \left[ R \int_0^T \left\{ \frac{\partial f_j}{\partial X_k}(t; X(0)) \right\}^2 dt \right], \quad (30)$$

where the expectation is taken with respect to the random initial values  $X(0)$  and  $R$  is the number of initial values. The measure  $D_{j,k}$  in (30) is a loose analogy to

$\left\{ \int_0^1 [f_{jk}^*(X_k^*(t))]^2 dt \right\}^{1/2}$  used in Assumption 6. Note that if no edge is present from  $X_k$  to  $X_j$ ,

then  $f_j / X_k \equiv 0$  and hence  $D_{j,k}(v) = 0$ . One immediately notes that, as  $R$  increases, the regulatory effect for a true edge increases proportionally to  $R$ , while the regulatory effect of a nonedge remains zero. For the self-edges in  $\mathcal{E}_1$  and the nonself-edges in  $\mathcal{E}_2$ , we can define  $D^{(1)}(v)$  and  $D^{(2)}(v)$  as

$$\begin{aligned} D^{(1)}(v) &= \min_{k=1, \dots, 10} D_{k,k}(v), \quad \text{and} \\ D^{(2)}(v) &= \min_{k=1, \dots, 5} \{D_{2k-1, 2k}(v), D_{2k, 2k-1}(v)\}, \end{aligned} \quad (31)$$



where we use the minimum because variable selection is limited by the minimum regulatory effect (see Assumption 6). With a slight abuse of definition, we refer to (31) as the minimum regulatory effects in a nonadditive model.

We apply GRADE using the formulation in (18). The sparsity parameter  $\lambda$  is chosen so that there are 20 directed edges in the estimated network. We record the number of estimated edges that are in  $\mathcal{E}_1$  and  $\mathcal{E}_2$ . The edge recovery performance is shown in Figure 3(b). In Figure 3(c), we display the minimum regulatory effects defined in (31). Edge recovery and minimum regulatory effects show a similar trend as a function of  $r$  in (29). This suggests that (31), and thus (30), is a reasonable measure of the additive components of the regulatory effect of the edges. The slight deviation between the trends reflects the fact that the measure

defined in (30) is not a direct counterpart of  $\left\{ \int_0^1 [f_{jk}^*(X_k^*(t))]^2 dt \right\}^{1/2}$  in a nonadditive model.

The edge recovery improves when a larger value of  $R$  is used, though these results are omitted due to space constraints. Our results indicate that GRADE can recover the true graph even when the additivity assumption is violated, provided that the regulatory effects (30) for the true edges are sufficiently large.

## 6. Applications

### 6.1. Application to in Silico Gene Expression Data

GeneNetWeaver (GNW) provides an in silico benchmark for assessing the performance of network recovery methods (Schaffter, Marbach, and Floreano 2011), and was used in the third DREAM challenge (Marbach et al. 2009). GNW is based upon real gene regulatory networks of yeast and *E. coli*. It extracts sub-networks from the yeast or *E. coli* gene regulatory networks, and assigns a system of ODEs to the extracted network. This system of ODEs is nonadditive, and includes unobserved variables (Marbach et al. 2010). Therefore, the assumptions of GRADE are violated in the GNW data.

To mimic real-world laboratory experiments, GNW provides several data generation mechanisms. In this study, we consider data from the *perturbation* experiments. The perturbation experiments are similar to the data-generating mechanisms used in Section 5.3, where initial conditions of the ODE system are perturbed to emulate the diversity of trajectories from multiple independent experiments.

We investigate 10 networks from GNW that have been previously studied by Henderson and Michailidis (2014), of which five have 10 nodes and five have 100 nodes. For each network, GNW provides one set of noiseless gene expression data consisting of  $R$  perturbation experiments where the trajectories are measured at  $n = 21$  evenly spaced time points in  $[0, 1]$ . Here  $R = 10$  for the five 10-node networks and  $R = 100$  for the five 100-node networks. As in Henderson and Michailidis (2014), we add independent  $N(0, 0.025^2)$  measurement errors to the data at each timepoint.

We apply NeRDS as described in Henderson and Michailidis (2014). We apply GRADE using the formulation (18) to handle observations from multiple experiments, with the smoothing estimates  $\hat{X}$  in (17b) fit using smoothing splines with bandwidth chosen by GCV,

and using cubic splines with two internal knots as the basis functions in (17c). The integral  $\widehat{\Psi}_k(t) = \int_0^t \psi(\widehat{X}_k(u; h)) du$  in (17c) is calculated numerically with step size 0.01. Finally, we apply an additional  $\ell_2$ -type penalty to the  $\theta_{jk}$ 's in (18) to match the setup of NeRDS. The tuning parameter for this penalty is set to be 0.1.

Results are shown in Table 1. Recall that the data-generating mechanism violates crucial assumptions for both NeRDS and GRADE. We see in Table 1 that NeRDS outperforms GRADE in one network, while GRADE outperforms NeRDS in the other nine networks. This suggests that GRADE is a competitive exploratory tool for reconstructing gene regulatory networks.

## 6.2. Application to Calcium Imaging Recordings

In this section, we consider the task of learning regulatory relationships among populations of neurons. We investigate the calcium imaging recording data from the Allen Brain Observatory project conducted by the Allen Institute for Brain Science (available at <http://observatory.brain-map.org>). Here, we investigate one of the experiments in the project. In this experiment, calcium fluorescence levels (a surrogate for neuronal activity) are recorded at 30 Hz on a region of the primary visual cortex while the subject mouse is shown 40 visual stimuli. The 40 visual stimuli are combinations of eight spatial orientations and five temporal frequencies. Each stimulus lasts for 2 sec and is repeated 15 times. The recorded videos are processed by the Allen Institute to identify individual neurons. In this particular experiment, there are 575 neurons. Each neuron's activity is defined as the average calcium fluorescence level of the pixels that it covers in the video.

It is known that the activities of individual neurons are noisy and sometimes misleading (Cunningham and Byron 2014). As an alternative, neuronal populations can be studied (see, e.g., Part Three of Gerstner, Kistler, Naud, and Paninski 2014). We define 25 neuronal populations by dividing the recording region into a  $5 \times 5$  grid, where each population contains roughly 20 neurons. We use GRADE to capture the functional connectivity among the 25 neuronal populations. Note that functional connectivity is distinct from physical connectivity. Functional connectivity involves the relationships among neuronal populations that can be observed through neuron activities and may change across stimuli, whereas physical connectivity consists of synaptic interactions.

We estimate the functional connectivity corresponding to three different but related stimuli, consisting of frequencies of 1 Hz, 2 Hz, and 4 Hz, each at a spatial orientation of  $90^\circ$ . For each stimulus, we have calcium fluorescence levels of the  $p = 25$  neuronal populations for each of  $R = 15$  repetitions. Since each repetition spans 2 sec and the calcium fluorescence is recorded at 30 Hz, there are 60 timepoints per repetition. We apply GRADE using the formulation in (18) to reconstruct the functional connectivity under each of the three stimuli. We use smoothing splines with bandwidth  $h$  selected with GCV to estimate  $\widehat{X}$  in (17b), and use cubic splines with four internal knots as the basis functions  $\psi(\cdot)$  in (17c). The sparsity parameter  $\lambda_{j,n}$  for each nodewise regression in (18) is selected using BIC for each  $j = 1, \dots, 25$ . For ease of visualization, we prefer a sparse network, and so we fit GRADE using tuning

parameter values  $\alpha(\lambda_{1,n}, \dots, \lambda_{p,n})$ , where the scalar  $\alpha$  is selected so that each of the estimated networks contains approximately 25 edges.

Estimated functional connectivities are shown in Figure 4. We see that, in all three networks, the 24th neuronal population regulates many other neuronal populations, indicating that this region may contain neurons that are sensitive to this spatial orientation. Furthermore, we see that the adjacent connectivity networks in Figure 4 are somewhat similar to each other, whereas the networks at 1 Hz and 4 Hz have few similarities. This agrees with the observation in neuroscience that neurons in the mouse primary visual cortex are responsive to a somewhat narrow range of temporal frequencies near their peak frequencies (see, e.g., Gao, DeAngelis, and Burkhalter 2010).

## 7. Discussion

In this article, we propose a new approach, GRADE, for estimating a system of high-dimensional additive ODEs. GRADE involves estimation of an integral rather than a derivative. We show that estimating the integral is superior to estimating the derivatives both theoretically and empirically. We leave an extension of our work to nonadditive ODEs to future research.

In this article, we have not addressed the issue of experimental design. Given a finite set of resources, one may choose to design an experiment to measure  $n$  observations on a very dense time grid, or on a coarse time grid. Alternatively, one might choose to measure  $n/R$  observations for  $R$  distinct experiments from a single ODE system (1), each with a different initial condition. This presents a trade-off that is especially interesting in the context of ODEs: using a dense time grid improves the quality of the smoothing estimates  $\hat{X}$ , as seen in Sections 5.1 and 5.2, while running multiple experiments enhances the identifiability of the true structure, as seen in Section 5.3. We leave a more detailed treatment of these issues to future work.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

The authors thank the associate editor and two anonymous reviewers for helpful comments. The authors thank the authors of Brunel, Clairon, and d'Alché Buc (2014), Hall and Ma (2014), Henderson and Michailidis (2014), and Wu et al. (2014) for sharing their code for their proposals, and for responding to their inquiries. The authors thank the Allen Institute for Brain Science for providing the dataset analyzed in Section 6.2.

### Funding

A.S. was supported by NSF grant DMS-1561814 and NIH grants 1K01HL124050-01A1 and 1R01GM114029-01A1, and D.W. was supported by NIH Grant DP5OD009145, NSF CAREER Award DMS-1252624, and an Alfred P. Sloan Foundation Research Fellowship.

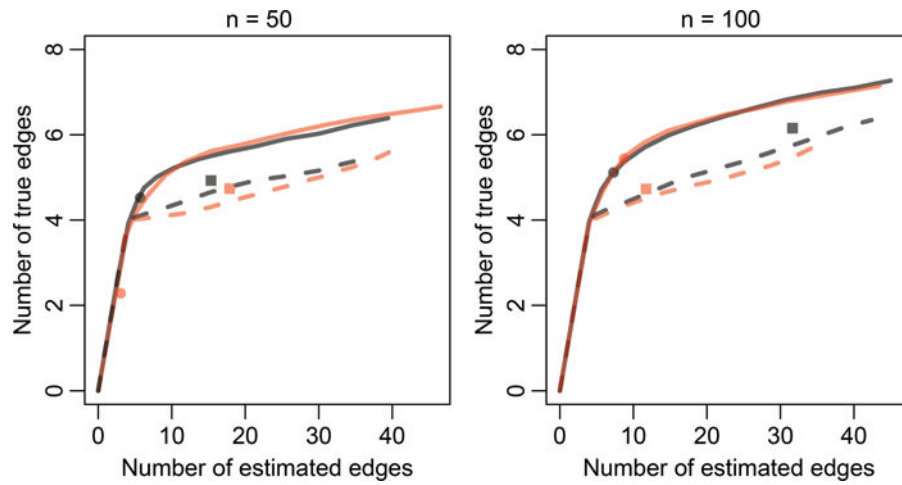
## References

Benson M. Parameter Fitting in Dynamic Models. *Ecological Modelling*. 1979; 6:97–115.

- Biegler LT, Damiano JJ, Blau GE. Nonlinear Parameter Estimation: A Case Study Comparison. *AIChE Journal*. 1986; 32:29–45.
- Boucheron, S., Lugosi, G., Massart, P. Concentration Inequalities: A Nonasymptotic Theory of Independence, With a Foreword by Michel Ledoux. Oxford, UK: Oxford University Press; 2013.
- Brunel NJ-B. Parameter Estimation of ODE's via Nonparametric Estimators. *Electronic Journal of Statistics*. 2008; 2:1242–1267.
- Brunel NJ-B, Clairon Q, d'Alché Buc F. Parametric Estimation of Ordinary Differential Equations with Orthogonality Conditions. *Journal of the American Statistical Association*. 2014; 109:173–185.
- Buja A, Hastie TJ, Tibshirani RJ. Linear Smoothers and Additive Models. *Annals of Statistics*. 1989; 17:453–555.
- Cao J, Wang L, Xu J. Robust Estimation for Ordinary Differential Equation Models. *Biometrics*. 2011; 67:1305–1313. [PubMed: 21401565]
- Cao J, Zhao H. Estimating Dynamic Models for Gene Regulation Networks. *Bioinformatics*. 2008; 24:1619–1624. [PubMed: 18505754]
- Chou IC, Voit EO. Recent Developments in Parameter Estimation and Structure Identification of Biochemical and Genomic Systems. *Mathematical Biosciences*. 2009; 219:57–83. [PubMed: 19327372]
- Cunningham JP, Byron MY. Dimensionality Reduction for Large-Scale Neural Recordings. *Nature Neuroscience*. 2014; 17:1500–1509. [PubMed: 25151264]
- Dattner I, Klaassen CAJ. Optimal Rate of Direct Estimators in Systems of Ordinary Differential Equations Linear in Functions of the Parameters. *Electronic Journal of Statistics*. 2015; 9:1939–1973.
- Ellner SP, Seifu Y, Smith RH. Fitting Population Dynamic Models to Time-Series Data by Gradient Matching. *Ecology*. 2002; 83:2256–2270.
- Friedman JH, Hastie TJ, Tibshirani RJ. Sparse Inverse Covariance Estimation with the Graphical Lasso. *Biostatistics*. 2008; 9:432–441. [PubMed: 18079126]
- Gao E, DeAngelis GC, Burkhalter A. Parallel Input Channels to Mouse Primary Visual Cortex. *The Journal of Neuroscience*. 2010; 30:5912–5926. [PubMed: 20427651]
- Gerstner, W., Kistler, WM., Naud, R., Paninski, L. *Neuronal Dynamics: From Single Neurons to Networks and Models of Cognition*. Cambridge, UK: Cambridge University Press; 2014.
- Gugushvili S, Klaassen CAJ. n- Consistent Parameter Estimation for Systems of Ordinary Differential Equations: Bypassing Numerical Integration via Smoothing. *Bernoulli*. 2012; 18:1061–1098.
- Hall P, Ma Y. Quick and Easy One-Step Parameter Estimation in Differential Equations. *Journal of the Royal Statistical Society, Series B*. 2014; 76:735–748.
- Henderson J, Michailidis G. Network Reconstruction Using Nonparametric Additive Ode Models. *PLoS ONE*. 2014; 9:e94003. [PubMed: 24732037]
- Izhikevich, EM. *Dynamical Systems in Neuroscience: The Geometry of Excitability and Bursting*, Computational Neuroscience. Cambridge, MA: MIT Press; 2007.
- Lee JD, Sun Y, Taylor JE. On Model Selection Consistency of Regularized M-Estimators. *Electronic Journal of Statistics*. 2013; 9:608–642.
- Liang H, Wu H. Parameter Estimation for Differential Equation Models using a Framework of Measurement Error in Regression Models. *Journal of the American Statistical Association*. 2008; 103:1570–1583. [PubMed: 19956350]
- Loader, C. *Statistics and Computing*. New York: Springer; 1999. Local Regression and Likelihood; p. 1-290.
- Loh PL, Wainwright MJ. High-Dimensional Regression With Noisy and Missing Data: Provable Guarantees with Nonconvexity. *Annals of Statistics*. 2012; 40:1637–1664.
- Lu T, Liang H, Li H, Wu H. High-Dimensional ODEs Coupled With Mixed-Effects Modeling Techniques for Dynamic Gene Regulatory Network Identification. *Journal of the American Statistical Association*. 2011; 106:1242–1258. [PubMed: 23204614]
- Ma W, Trusina A, El-Samad H, Lim WA, Tang C. Defining Network Topologies that can Achieve Biochemical Adaptation. *Cell*. 2009; 138:760–773. [PubMed: 19703401]

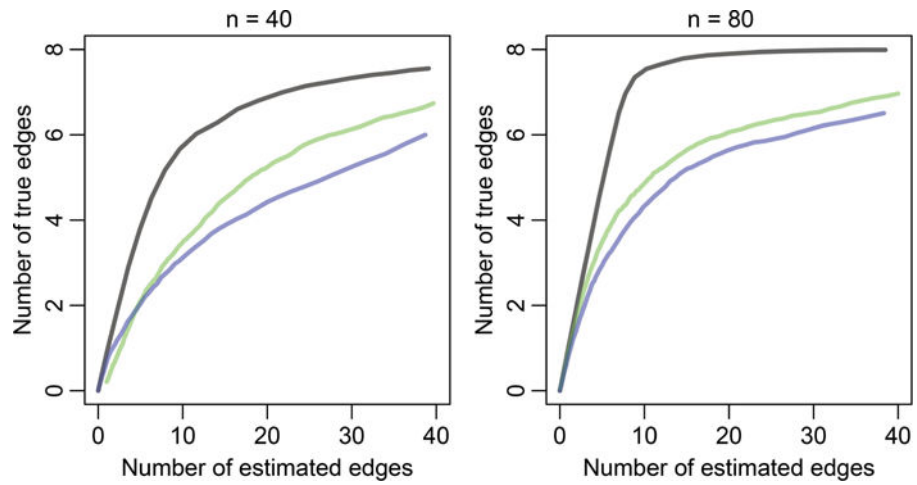
- Marbach D, Prill RJ, Schaffter T, Mattiussi C, Floreano D, Stolovitzky G. Revealing Strengths and Weaknesses of Methods for Gene Network Inference. *Proceedings of the National Academy of Sciences*. 2010; 107:6286–6291.
- Marbach D, Schaffter T, Mattiussi C, Floreano D. Generating Realistic *in silico* Gene Networks for Performance Assessment of Reverse Engineering Methods. *Journal of Computational Biology*. 2009; 16:229–239. [PubMed: 19183003]
- Meier L, van de Geer S, Bühlmann P. The Group Lasso for Logistic Regression. *Journal of the Royal Statistical Society, Series B*. 2008; 70:53–71.
- Meinshausen N, Bühlmann P. High-Dimensional Graphs and Variable Selection with the Lasso. *Annals of Statistics*. 2006; 34:1436–1462.
- Meinshausen N, Bühlmann P. Stability Selection. *Journal of the Royal Statistical Society, Series B*. 2010; 72:417–473.
- Miao H, Xia X, Perelson A, Wu H. On Identifiability of Nonlinear ODE Models and Applications in Viral Dynamics. *SIAM Review*. 2011; 53:3–39. [PubMed: 21785515]
- Qi X, Zhao H. Asymptotic Efficiency and Finite-Sample Properties of the Generalized Profiling Estimation of Parameters in Ordinary Differential Equations. *Annals of Statistics*. 2010; 38:435–481.
- Ramsay JO, Hooker G, Campbell D, Cao J. Parameter Estimation for Differential Equations: A Generalized Smoothing Approach (with discussions and a reply by the authors). *Journal of the Royal Statistical Society, Series B*. 2007; 69:741–796.
- Ravikumar PK, Lafferty J, Liu H, Wasserman L. Sparse Additive Models. *Journal of the Royal Statistical Society, Series B*. 2009; 71:1009–1030.
- Schaffter T, Marbach D, Floreano D. Genenetweaver: In Silico Benchmark Generation and Performance Profiling of Network Inference Methods. *Bioinformatics*. 2011; 27:2263–2270. [PubMed: 21697125]
- Tibshirani RJ. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society, Series B*. 1996; 58:267–288.
- Tsybakov, AB. *Introduction to Nonparametric Estimation* (Springer Series in Statistics). New York: Springer; 2009. (Revised and extended from the 2004 French original, Translated by Vladimir Zaiats)
- Varah JM. A Spline Least Squares Method for Numerical Parameter Estimation in Differential Equations. *SIAM Journal on Scientific Computing*. 1982; 3:28–46.
- Volterra V. Variations and Fluctuations of the Number of Individuals in Animal Species Living Together. *Journal of Marine and Freshwater Research*. 1928; 3:3–51.
- Voorman AL, Shojaie A, Witten DM. Graph Estimation With Joint Additive Models. *Biometrika*. 2014; 101:85–101. [PubMed: 25013234]
- Wainwright MJ. Sharp Thresholds for High-Dimensional and Noisy Sparsity Recovery using  $\ell_1$ -Constrained Quadratic Programming (Lasso). *IEEE Transactions on Information Theory*. 2009; 55:2183–2202.
- Wang H, Leng C. Unified LASSO Estimation by Least Squares Approximation. *Journal of the American Statistical Association*. 2007; 102:1039–1048.
- Wu H. Statistical Methods for HIV Dynamic Studies in AIDS Clinical Trials. *Statistical Methods in Medical Research*. 2005; 14:171–192. [PubMed: 15807150]
- Wu H, Lu T, Xue H, Liang H. Sparse Additive Ordinary Differential Equations for Dynamic Gene Regulatory Network Modeling. *Journal of the American Statistical Association*. 2014; 109:700–716. [PubMed: 25061254]
- Xia Y, Li WK. Asymptotic Behavior of Bandwidth Selected by the Cross-Validation Method for Local Polynomial Fitting. *Journal of Multivariate Analysis*. 2002; 83:265–287.
- Xue H, Miao H, Wu H. Sieve Estimation of Constant and Time-Varying Coefficients in Nonlinear Ordinary Differential Equation Models by Considering both Numerical Error and Measurement Error. *Annals of Statistics*. 2010; 38:2351–2387. [PubMed: 21132064]
- Yuan M, Lin Y. Model Selection and Estimation in Regression With Grouped Variables. *Journal of the Royal Statistical Society, Series B*. 2006; 68:49–67.

- Yuan M, Lin Y. Model Selection and Estimation in the Gaussian Graphical Model. *Biometrika*. 2007; 94:19–35.
- Zhang X, Cao J, Carroll RJ. On the Selection of Ordinary Differential Equation Models with Application to Predator-Prey Dynamical Models. *Biometrics*. 2015; 71:131–138. [PubMed: 25287611]
- Zhao P, Yu B. On Model Selection Consistency of Lasso. *Journal of Machine Learning Research*. 2006; 7:2541–2563.



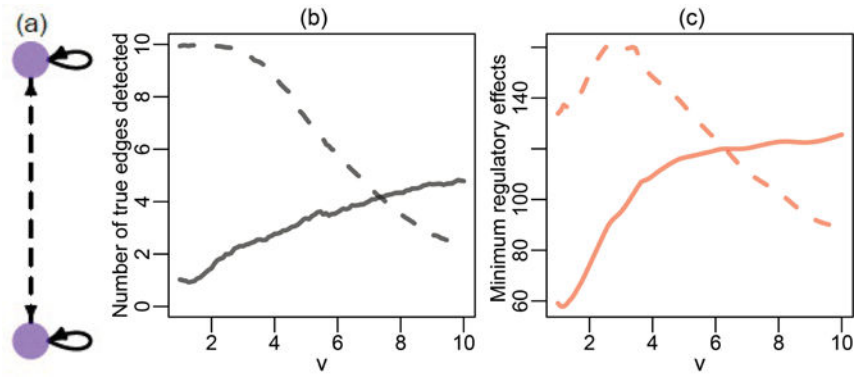
**Figure 1.**

Performance of network recovery methods on the system of additive ODEs in (26), averaged over 400 simulations. The four curves represent SA-ODE (dashed, red line), NeRDS (dashed, gray line), and GRADE without (solid, red line) and with (solid, gray line) the additional smoothing penalty in (17a) used by NeRDS. Each point on the curves corresponds to average performance for a given sparsity tuning parameter  $\lambda_n$  in (14a) or (17a). The symbols indicate the sparsity tuning parameter  $\lambda_n$  selected using BIC (SA-ODE, red square, and GRADE, red circle and gray circle) or GCV (NeRDS, gray square).



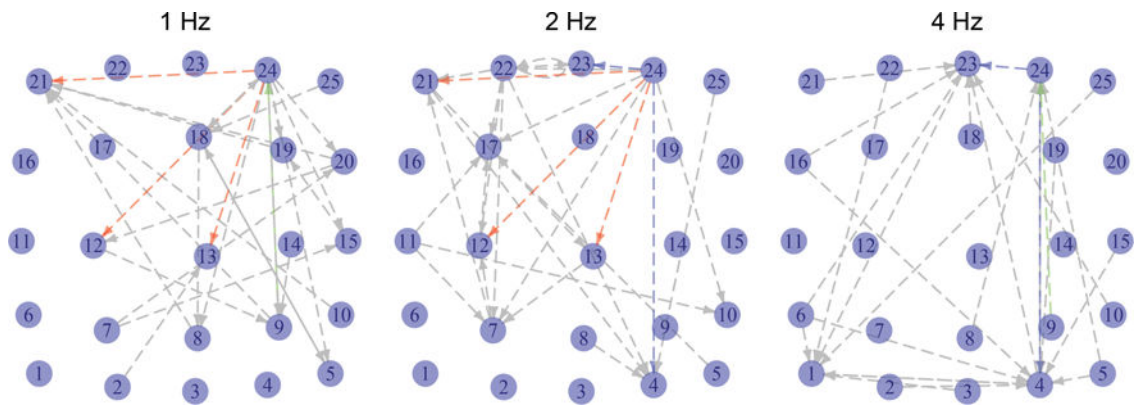
**Figure 2.** Network recovery on the system of linear ODEs (27), averaged over 200 simulated datasets. The three curves represent GRADE (gray line), Hall and Ma (2014) (blue line), Brunel, Clairon, and d'Alché Buc (2014) (green line).





**Figure 3.**

(a) The graph encoded by a pair of Lotka-Volterra equations as given in (29). Self-edges (solid, gray line) and nonself-edges (dashed, gray line) are shown. (b) Self-edge (solid, gray line) and nonself-edge (dashed, gray line) recovery of GRADE, averaged over 200 simulated datasets. (c) Minimum signals defined in (31), for self-edges,  $D^{(1)}(\cdot)$  (solid, red line), and nonself-edges,  $D^{(2)}(\cdot)$  (dashed, red line).



**Figure 4.**

Estimated functional connectivities among neuronal populations from the calcium imaging data described in Section 6.2. Each node is positioned near the center of the neuronal population it represents, with jitter added for ease of display. The three red edges are shared between the estimated networks at 1 Hz and 2 Hz; the two blue edges are shared between estimated networks at 2 Hz and 4 Hz; the single green edge is shared between the estimated networks at 1 Hz and 4 Hz. For reference, given two Erdős-Rényi graphs consisting of 25 nodes and 25 edges, the probability of having three or more shared edges is 0.07, and the probability of having two or more shared edges is 0.26.

**Table 1**

Area under ROC curves for NeRDS and GRADE.

	$p = 10$		$p = 100$	
	NeRDS	GRADE	NeRDS	GRADE
Ecol11	0.450 (0.438, 0.462)	<b>0.545</b> (0.534, 0.557)	0.624 (0.622, 0.627)	<b>0.670</b> (0.667, 0.673)
Ecol12	0.512 (0.502, 0.523)	<b>0.643</b> (0.634, 0.653)	0.637 (0.635, 0.640)	<b>0.653</b> (0.650, 0.656)
Yeast1	0.486 (0.476, 0.495)	<b>0.679</b> (0.666, 0.691)	0.610 (0.607, 0.612)	<b>0.636</b> (0.635, 0.638)
Yeast2	0.525 (0.518, 0.532)	<b>0.607</b> (0.600, 0.613)	0.568 (0.566, 0.569)	<b>0.584</b> (0.582, 0.585)
Yeast3	0.467 (0.460, 0.474)	<b>0.576</b> (0.566, 0.587)	<b>0.617</b> (0.616, 0.619)	0.567 (0.566, 0.568)

NOTES: The average area under the curves and 90% confidence intervals, over 100 simulated datasets. Networks and data-generating mechanisms are described in Section 6.1. Boldface indicates the method with larger AUC.