

RESEARCH ARTICLE

Open Access

Set membership experimental design for biological systems

Skyler W Marvel and Cranos M Williams*

Abstract

Background: Experimental design approaches for biological systems are needed to help conserve the limited resources that are allocated for performing experiments. The assumptions used when assigning probability density functions to characterize uncertainty in biological systems are unwarranted when only a small number of measurements can be obtained. In these situations, the uncertainty in biological systems is more appropriately characterized in a bounded-error context. Additionally, effort must be made to improve the connection between modelers and experimentalists by relating design metrics to biologically relevant information. Bounded-error experimental design approaches that can assess the impact of additional measurements on model uncertainty are needed to identify the most appropriate balance between the collection of data and the availability of resources.

Results: In this work we develop a bounded-error experimental design framework for nonlinear continuous-time systems when few data measurements are available. This approach leverages many of the recent advances in bounded-error parameter and state estimation methods that use interval analysis to generate parameter sets and state bounds consistent with uncertain data measurements. We devise a novel approach using set-based uncertainty propagation to estimate measurement ranges at candidate time points. We then use these estimated measurements at the candidate time points to evaluate which candidate measurements furthest reduce model uncertainty. A method for quickly combining multiple candidate time points is presented and allows for determining the effect of adding multiple measurements. Biologically relevant metrics are developed and used to predict when new data measurements should be acquired, which system components should be measured and how many additional measurements should be obtained.

Conclusions: The practicability of our approach is illustrated with a case study. This study shows that our approach is able to 1) identify candidate measurement time points that maximize information corresponding to biologically relevant metrics and 2) determine the number at which additional measurements begin to provide insignificant information. This framework can be used to balance the availability of resources with the addition of one or more measurement time points to improve the predictability of resulting models.

Background

Costly materials, limited resources, and lengthy experiments are constraints that hinder our ability to acquire quantifiable measurements from biological systems. Experimental design approaches are computational techniques for extracting the most useful information from experiments yet to be performed [1]. These techniques are needed for the study of biological systems to conserve the limited resources that are allocated for performing experiments. Application of these techniques to

biological systems has introduced novel mathematical algorithms and models to life sciences, while also requiring the development of new mathematical theories and programming tools [2]. An important aspect of experimental design for biological systems is model calibration, which requires the estimation of parameters such as kinetic and diffusivity constants [3]. The development of accurate biological models is constrained by the financial costs and time required to perform biological experiments, often leading to a collection of sparse datasets with which to estimate the parameters of proposed model structures. Experimental design provides a method to yield the best estimates from data given the

* Correspondence: cmwilli5@ncsu.edu
Department of Electrical and Computer Engineering, North Carolina State University, Raleigh, NC 27695, USA

limitations in data collection, component observability and limited system excitability.

The development and application of experimental design has a rich history spread across a wide range of fields. An excellent review article by Pronzato has condensed the underlying concepts behind the most widely used techniques of experimental design for nonparametric and parametric models [1]. The reader is referred to the review article and the works cited therein for a thorough understanding of statistical methods for experimental design.

Typically, parameter estimation problems begin by claiming that observations \hat{y} are perturbed from ideal model outputs $\mathbf{g}(\mathbf{x}, \theta^*)$ by an error ε , such that

$$\hat{y}_i = \mathbf{g}(\mathbf{x}_i, \theta^*) + \varepsilon_i, i = 1, \dots, k, \quad (1)$$

where \mathbf{x}_i are the model states at k different times or experimental conditions, θ^* are the true parameter values and the errors, ε_i , are statistically independent with zero mean and variance $E(\varepsilon_i^2) = \sigma^2(\mathbf{x}_i)$. It is assumed that the errors can be defined by probability density functions, often assumed to be independent and identically distributed Gaussian random variables with zero mean and variance σ^2 for mathematical convenience. The unknown parameter vector can then be determined by the maximum likelihood estimate $\hat{\theta}_{ML}^k$. As $k \rightarrow \infty$ the difference between $\hat{\theta}_{ML}^k$ and θ^* can be described by a normal distribution with zero mean and covariance matrix, Σ , which is bounded from below by the inverse of the Fisher Information Matrix (FIM) according to the Cramér-Rao inequality [1].

Experimental design aims to maximize information, or minimize uncertainty, about unknown model parameters by exploring experimental configurations such as the sampling times where new measurements should be acquired, the desired number of measurements to add, which system components should be measured, etc. The criteria used to evaluate the information of a design are derived from scalar functions of the FIM [1]. A-optimal design, for example, minimizes $\text{trace}(\text{FIM}^{-1})$, or equivalently minimizes the sum of squared lengths of the axes of asymptotic confidence ellipsoids for θ . E-optimality refers to designs where the longest axis of asymptotic confidence ellipsoids for θ is minimized, which is equivalent to maximizing the minimum eigenvalue of the FIM. D-optimal design maximizes $\det(\text{FIM})$ and corresponds to minimizing the volume of asymptotic confidence ellipsoids for θ .

Although there is a large body of work dedicated to experimental design using statistical methods [1], several problems arise when using these approaches for the modeling of biological systems [4]. Kreutz and Timmer state several of the difficulties in using experimental

design for biological systems: i) models are often large and the number of measurements is very limited, ii) relative noise levels of 10% or more are standard for biochemical data, iii) little prior knowledge exists. These considerations make it difficult to correctly characterize the distribution of uncertainty in the model, which is the primary pillar upon which FIM approaches for experimental design are based. Even if the correct distribution is obtained, accurate parameter estimations using the FIM are usually valid only when a large number of data points are available, which is not often the case for biological systems [5]. Rather, the finite range of values that system component concentrations can take on at a given time more appropriately characterizes the uncertainty in biological systems. This bound can be inferred based on the experimental technology, the characteristics of limited replicates, and/or first principles knowledge. Therefore, a set membership framework is more appropriate for the development of experimental design for many biological systems, where the error is bounded with no other hypothesis given regarding its distribution [6].

A key aspect of experimental design for bounded-error models is how to characterize the set of parameter values that are consistent with all data measurements. Initial methods for constructing this set use conservative bounding approaches based on ellipsoids to characterize the parameter sets. More precise parameter set estimations can be obtained using interval analysis [7,8], *but these interval techniques have not previously been applied to experimental design approaches*. Apart from the method used to bound the parameter set, proper experimental design metrics are important because they provide a logical link between physical resources and mathematical constructs. Traditional experimental design criteria for bounded-error models minimize the volume of parameter sets that are consistent with the data [6,9-11]. However, the information provided by this metric may not be useful to a biologist. Other metrics that are related directly to the uncertainty of specific parameters or the effects on unmeasurable model states may be of more interest. Such biologically relevant information can be obtained from simple criteria functions previously not used in experimental design for bounded-error models. Set membership experimental design methods have recently regained attention. Hasebauer et al. have developed a set-based experimental design method using semidefinite programming with V-optimality as the only design metric [12]. The expected information content from additional measurements is determined using a Monte-Carlo approach to simulate different parameters, input sequences and measurement errors. While this method demonstrates the usefulness of bounded-error techniques, there is a lack

of connection between the design metric and biological interpretation. Additionally, the use of a Monte-Carlo approach to simulate the effect of additional measurements requires a large number of simulations and can be very time consuming. Bounded approaches, such as the one we outline in this paper, allow for the impact of uncertainty to be assessed without needing to perform Monte-Carlo simulations.

In this work, we develop an experimental design framework that utilizes interval analysis to generate the set of parameters and state bounds consistent with all data measurements. This approach leverages many of the recent advances in bounded-error parameter and state estimation methods [7,8], including increased accuracy through the use of interval analysis instead of bounded ellipsoids, as the base of our experimental design framework. Our novel framework uses parameter and state estimations based on initial data measurements, which may provide data for only a subset of the model states, to estimate measurement bounds at candidate time points of interest to the experimenter (times when measurements have not been taken). We then use these estimated measurements at the candidate time points to evaluate which candidate measurements furthest reduce model uncertainty. We propose a method for combining candidate time points to determine the effect of adding multiple measurements. We present biologically relevant design metrics to evaluate candidate designs in order to address issues associated with making a better connection between modelers and experimentalists. These contributions comprise a bounded-error experimental design framework that can be applied to nonlinear continuous-time systems when few data measurements are available. This framework can be used to balance the availability of resources with the addition of one or more measurement points to improve the predictability of resulting models.

Methods

In this section, we define a specific experimental design problem and outline how our framework is used to determine the number of additional measurements that are warranted and at what time points these measurements should be taken. The relevant interval arithmetic algorithms for parameter and state estimation used throughout this process are briefly presented. We show how to select a set of candidate time points based on the estimated state bounds of a proposed model given initial data measurements and provide a method to estimate the corresponding candidate measurement bounds. Techniques for determining the effect of adding multiple candidate time points on parameter and state estimations are discussed. We define several biologically relevant metrics, which are scalar functions of the

parameter and state estimations after incorporating estimated candidate time point measurements. These metrics can convey information such as the activity of specific enzyme kinetic parameters or bounding values for the estimation of unmeasured component concentrations.

Problem statement

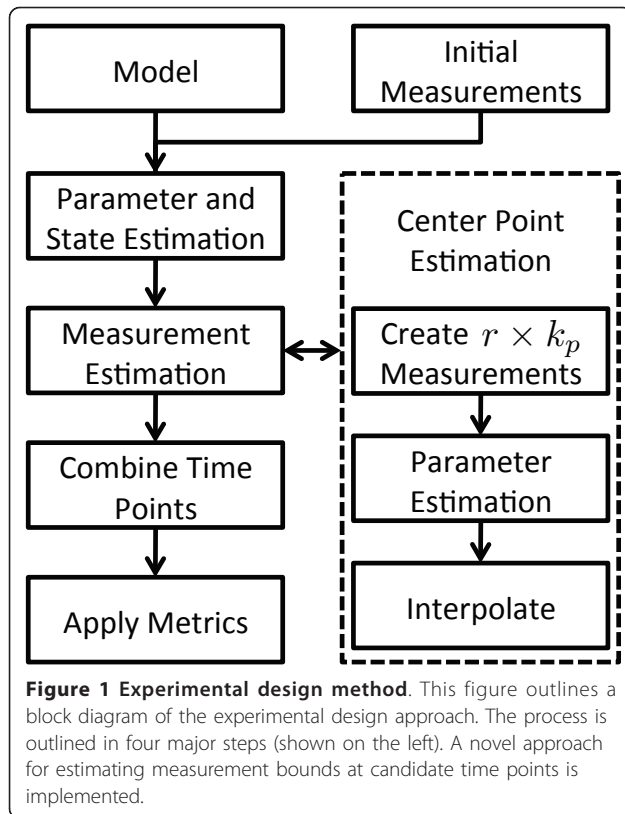
Consider the following ordinary differential equation (ODE) model of a biological system:

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{f}(\mathbf{x}(t), \theta) \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}(t), \theta),\end{aligned}\tag{2}$$

Where $\mathbf{x} \in \mathbb{R}^n$ is an n -dimensional vector of component concentrations, $\mathbf{y} \in \mathbb{R}^m$ is an m -dimensional vector of measurements, and $\theta \in \mathbb{R}^p$ are the p model parameters. An initial set of bounded data measurements has been obtained at k different times:

$\mathcal{Y} := \left\{ \hat{\mathbf{y}}_i | \underline{\mathbf{y}}_i \leq \mathbf{y}(t_i) \leq \bar{\mathbf{y}}_i; i = 1, \dots, k \right\}$, where i is the index corresponding to time t_i and $\underline{\mathbf{y}}_i$ and $\bar{\mathbf{y}}_i$ are the lower and upper measurement bounds, respectively. The problem under study is to determine at what time points to collect new data measurements for minimizing or maximizing specific parameter and/or state information metrics.

We use the method outlined in the left half of Figure 1 to solve this problem using a set membership approach by applying biologically relevant information metrics to evaluate candidate time points. First, we perform bounded parameter estimation using the initial bounded-error measurements. Estimated state bounds are then generated from the resulting parameter space. Second, a set of candidate time points is selected from locations where relatively large uncertainties exist in the estimated model states. We propose a novel approach to estimate the measurement bounds at candidate time points using a set-based approach that incorporates the initial bounded-error measurements adjacent to each candidate time point. Third, we perform bounded parameter and state estimations that incorporate the candidate measurements to predict the possible effects of adding a measurement at the corresponding time point. We also assess the impact of adding multiple measurements on the resulting estimates. As a proof of concept, we compare the performance of our estimated measurements and true measurements at each candidate time point, assessing the ability of estimated measurements to predict which candidate time point most reduces a given uncertainty metric. We assess this for single and combinations of candidate time points. We also use our estimated measurements at each candidate time point to identify the 'point of diminishing return' where additional measurements no longer provide



additional information, leading to no further decrease in estimate uncertainty.

Bounded estimation

These methods use interval analysis to computationally guarantee a valid bounded-error solution to the system of ODEs by employing interval box enclosures that bound the states during integration steps. Methods have been introduced in the literature to address overestimation due to wrapping effect [13-15] and to help reduce the computational burden for estimating parameters of complex, higher dimensional models [16], which are typical for biological processes.

Uncertainty propagation

Interval analysis is a form of guaranteed computing and can be used to generate solutions to ODEs through the use of interval boxes and inclusion functions [17]. Consider the model function \mathbf{g} , which maps a state interval box $[\mathbf{x}]$ to the corresponding image in the data space $\mathbf{g}([\mathbf{x}])$. Here the interval box $[\mathbf{x}]$ represents the Cartesian product of n scalar intervals $[\mathbf{x}] = [x_1] \times [x_2] \times \dots \times [x_n]$, where $[x_i]$ represent the interval $\underline{x}_i \leq x_i \leq \bar{x}_i$. A non-minimal inclusion function, \mathcal{G} , is a non-unique mapping from state space to data space and contains the smallest interval box that encloses the image $\mathbf{g}([\mathbf{x}])$.

Computing the solution of ODEs for $t_0 \leq t \leq t_N$ with time step h is done using Taylor expansions [17-19]. This method involves an inflation step where the bounds of the remainder term for the k^{th} -order Taylor expansion of the model ODEs are inflated by $1 \pm \alpha$. Evaluation of the Taylor expansion is performed using the Extended Mean Value (EMV) algorithm proposed by Rihm [19] using mean value forms [20] and matrix preconditioning. Whenever the EMV algorithm generates state values, $[\mathbf{x}]$, at a time where data measurements, $[\hat{\mathbf{y}}]$, are available, Set Inversion Via Interval Analysis (SIVIA) [21] is used to compare the two.

Set inversion

SIVIA is able to determine solution sets for unknown quantities \mathbf{u} from a functional relationship $\mathbf{q}(\mathbf{u}) = [\mathbf{y}]$. An *a priori* search space for \mathbf{u} is recursively explored using SIVIA to determine a guaranteed enclosure of the solution space. The resulting solution space is comprised of feasible and indeterminate boxes. These boxes, $[\mathbf{u}]$, are determined from the following relations: if $\mathbf{q}([\mathbf{u}]) \subseteq [\mathbf{y}]$ then $[\mathbf{u}]$ is *feasible*; if $\mathbf{q}([\mathbf{u}]) \cap [\mathbf{y}] = \varphi$ then $[\mathbf{u}]$ is *unfeasible*; else $[\mathbf{u}]$ is *indeterminate*. Indeterminate boxes are bisected and tested again until its widest dimension reaches a user specified threshold $\varepsilon > 0$.

Parameter and state estimation

The methods presented in this paper leverage the works of Jaulin for state estimation [7] and Raïssi et al. for parameter estimation [8]. Parameter estimation combines the EMV and SIVIA algorithms to systematically evaluate candidate boxes in the parameter space. Our framework uses these two algorithms to build our set-based experimental design approach. We perform parameter estimation by evaluating hypercubes in the partitioned parameter space to identify if each hypercube or box produces trajectories that are consistent with the measurements obtained from the system. A parameter box that produces trajectories that are inconsistent with any data measurement is classified as unfeasible and discarded. Any parameter box that produces a trajectory determined by SIVIA to be completely contained within all data measurements is labeled as feasible. All other parameter boxes are labeled as indeterminate. These indeterminate boxes are bisected and retained for further evaluation. We apply this bisection process recursively to any indeterminate box where the width of the widest dimension is larger than a user-defined length, $\varepsilon > 0$. We implemented the augmented estimation method presented by Marvel and Williams to enable its application for systems where few data measurements are available (Marvel S, Williams C: Set Membership and Parameter Estimation for Nonlinear Differential Equations Using Discrete Measurements, Submitted). We estimate

bounds on the resulting component concentrations consistent with the data measurements by executing the EMV algorithm using the parameter boxes classified as feasible and indeterminate. This state estimation will not only produce bounds between data measurements of measured states, but also provide bounds for unmeasurable states. We parallelized this method using the Message Passing Interface (MPI) protocol to distribute the boxes across multiple processors to effectively distribute computations across available processing resources [22].

Estimating candidate measurements

The measurements at a given time are characterized by an upper and lower bound such that $\underline{y} \leq \mathbf{y} \leq \bar{y}$. Mathematically, this measurement can be defined by three values: 1) the time t_j at which the measurement was observed, 2) the center point C_j , and 3) its range R_j such that $|C_j - \mathbf{y}(t_j)| \leq R_j/2$. We estimate the center points and ranges of candidate measurements using the bounds of adjacent data measurements and the estimated bounds on component concentration trajectories generated by the EMV algorithm. Once estimated, each candidate measurement is added to the original k data measurements to assess the impact of the additional measurement information on our ability to estimate the parameters and unmeasured states. We describe below how t_j , C_j , and R_j are estimated for candidate measurements.

To simplify notation in this subsection, we assume that one or more of the states can be directly measured ($\mathbf{y} = \mathbf{x}$). This will allow for direct comparison between estimated state bounds and measurement values. This is a common assumption made for biological systems [7,8]. In a more general case, comparisons would require use of the inclusion function \mathcal{G} to compare $\mathcal{G}(\mathbf{x})$ and \mathbf{y} via SIVIA.

Time point and range estimation

For a given state, time points for candidate measurements are chosen by first identifying all times t between measurements at t_i and t_{i+1} , whose estimated range (generated by the EMV algorithm) is greater than or equal to both of the measurement uncertainties at times t_i or t_{i+1} . This presents a worst case scenario because we are selecting candidate time points with the most possible uncertainty. Alternative time points can be selected based on practical experimental limitations or first principles knowledge. The set of time intervals, \mathbb{T} , for a corresponding state can be written as

$$T := \left\{ t \mid \bar{x}(t) - \underline{x}(t) \geq \max \left[\bar{x}(t^-) - \underline{x}(t^-), \bar{x}(t^+) - \underline{x}(t^+) \right] \right\} \quad (3)$$

where $t^- = \max_{t_i}(t_i < t)$ and $t^+ = \min_{t_i}(t_i > t)$. Selecting candidate time points from the intervals in \mathbb{T} is an empirical task. For example, a total of k_p candidate time points could be selected from within the interval set \mathbb{T} based on a collection of physically feasible time slots where measurements can be observed. The set of candidate time points is denoted as $\mathcal{T} := \{t_j; j = 1, \dots, k_p\}$, with $\mathcal{T} \subset \mathbb{T}$. The corresponding candidate time point ranges are determined based on a conservative premise that uses uncertainty information contained in adjacent measurements. Here, we set the range of candidate measurements to be

$$R_j = \max \left[\left(\bar{x}(t_j^-) - \underline{x}(t_j^-) \right), \left(\bar{x}(t_j^+) - \underline{x}(t_j^+) \right) \right], \quad (4)$$

where $t_j^- = \max_{t_i}(t_i < t_j)$ and $t_j^+ = \min_{t_i}(t_i > t_j)$. The amount of information available for determining appropriate range values is limited when no probabilistic assumptions are imposed on the uncertainty. Here, R_j is a relatively conservative estimate that assumes the uncertainty of the system at a new candidate time point is not less than that of data measurements taken near the same time.

Center point selection

Center point estimation is conservatively implemented to reduce the chance of erroneously eliminating valid kinetic parameters and component concentrations. We introduce a novel approach for estimating the corresponding center point of each candidate time point. This approach estimates the position of the center point C_j that maximizes the resulting parameter estimate volume at given time t_j and range R_j . The three main steps in this process are shown in the right half of Figure 1. First, r measurements are simulated at each $t_j \in \mathcal{T}$ by shifting R_j from the lower bound to the upper bound on the estimated state bounds. For example, if $r = 3$ and the estimate of state x at time t_4 is bounded between the range [3,6] with $R_4 = 1.5$, the resulting shifted candidate measurements at time t_4 would have bounds [3,4.5], [3.75,5.25] and [4.5,6]. Second, bounded parameter estimation is performed for each of the r shifted candidate measurements for each of the k_p candidate time points. Curve fits for each set of r parameter volumes are used to determine the center point, C_j , that maximizes the parameter volume for each candidate time point t_j . This allows us to fully construct conservative measurement estimations for candidate time point t_j using C_j and R_j .

Combining measurements

The ability to investigate the effects of adding multiple measurements is often desirable when designing

biological experiments. Employing a brute-force method for assessing the impact of all combinations of candidate measurements at t_1, t_2, \dots, t_{k_p} on estimated kinetic parameters and component concentrations is a large computational burden. A brute-force approach for exploring combinations of up to k_c candidate time points would require computing parameter and state estimations for $\sum_{m=2}^{k_c} \binom{k_p}{m}$ measurement sets. This is potentially problematic even for systems of low dimension and few unknown parameters. We hypothesize that there exists a level of independence among candidate time points that can be exploited to speed up our ability to evaluate the impact of multiple measurements on parameter uncertainty.

The estimated parameter space for a combination of candidate time points, \mathcal{P}_c , can be obtained by intersecting the parameter estimates of the individual time points, e.g. $\mathcal{P}_c = \mathcal{P}_1 \cap \mathcal{P}_4 \cap \mathcal{P}_9$. Computing the intersection of a set of non identical boxes is not an obvious task. We developed a simple approach for forming the union of sets of nonuniform shaped boxes by bisecting the larger feasible boxes until all boxes have widths less than ϵ . This approach allows boxes to be directly compared between estimated parameter sets. More sophisticated approaches can be applied that will preserve the largest possible feasible boxes during the intersection process. The estimated state bounds resulting from this combination of additional candidate measurements, \mathbf{x}_c , is then determined using the resulting intersected parameter boxes.

Metrics

Scalar functions of the estimated parameter set and state bounds are used as metrics to predict the impact of adding measurements at candidate times t_j on kinetic parameter and component concentration estimates. The metrics in this section can be conceptually related to traditional stochastic experimental design criteria functions (e.g. D-optimality, E-optimality, A-optimality [23]). However, the computation of these bounded-error metrics require no assumptions about underlying stochastic distributions of the model parameters or system states and relate directly to the physical components of the system. Thus, the biological interpretation of the bounded-error metrics is straightforward since they can be directly related to biological concepts instead of the mathematical construct of the FIM.

Parameter volume

We will evaluate the parameter volume as a means to compare our new metrics to traditional V- and D-optimality design criteria [6]. This metric will predict

the candidate time points that minimize the volume of the estimated parameter space. The parameter volume, \mathcal{P}_V , can easily be calculated by summing the volumes of the interval boxes,

$$\mathcal{P}_V = \sum_i \prod_{j=1}^p \Delta p_i^j, \quad (5)$$

where Δp_i^j is the width of the j^{th} dimension of the i^{th} parameter box. A drawback of this metric is the inability to detect large uncertainties in potentially important parameters if they are masked by less important but well known parameters. To combat this, the parameters could be weighted based on biological importance, giving more weight to parameter dimensions deemed important by the experimenter.

Parameter bounds

This metric can be customized for predicting candidate time points based on the uncertainty of a single parameter or a subset of parameters. Single parameter values are compared using the width of the uncertainty for the parameter of interest, e.g. $\mathcal{P}_{p_i} = \bar{p}_i - \underline{p}_i$. Multiple parameters are compared using the Euclidean norm to produce a scalar value from the widths of uncertainty for the selected parameters, e.g. $\mathcal{P}_{\|p_i, p_j\|_2} = \left[(\bar{p}_i - \underline{p}_i)^2 + (\bar{p}_j - \underline{p}_j)^2 \right]^{1/2}$.

State bounds

This metric utilizes estimated state bound information and allows the experimenter to see how estimated ranges of unmeasured states are affected by additional measurements. This may be of interest when constraining the range of state values is more important than parameter information. Also, the information provided by this metric is biologically meaningful because it provides a predicted limit on state values such as component concentrations. This metric is computed similarly to the parameter bounds metric but with the parameter uncertainties replaced by the maximum ranges of estimated states. Other custom metrics are also possible; for example, designing a metric to select the time points that minimizes the maximum value of a specific state.

Results and discussion

In this section, the proposed experimental design method is applied to an example problem. We evaluate our set-based experimental design approach by performing a proof of concept on a model that has been used in the literature to evaluate several other set-based approaches [7,8,14,15]. Our problem set-up is more stringent than the approach outlined in [8] because we assume only a small set of data measurements from a

single state is available as opposed to assuming data measurements are available at every time step. We use our approach to predict at what time additional measurements should be made in order to identify the candidate measurements that maximize information corresponding to previously defined metrics and to determine the number at which additional measurements begin to provide insignificant information.

Problem setup

The model under examinations is the Lotka-Volterra predator prey model, which is a canonical biological ODE model [24] and serves as a key model for testing algorithms in this field. This is a two-state model and is described by the following differential equations:

$$\begin{aligned} \dot{x}_1 &= x_1(p_1 - p_2x_2) \\ \dot{x}_2 &= -x_2(p_3 - p_4x_1), \end{aligned} \quad (6)$$

where x_1 is the prey population, x_2 is the predator population, p_1 is the prey birth rate, p_2 is the decrease in prey population due to encounters with predators, p_3 is the predator death rate, and p_4 is the increase in predator population due to encounters with prey. This model was used by Raïssi et al. to demonstrate their bounded parameter estimation algorithm when data measurements of the prey population are available for all $N = 1,400$ time points between $t_0 = 0$ and $t_N = 7$.

Initial data measurements were simulated by first generating model state values using exact inputs to the EMV algorithm and then adding uncertainty. The underlying state values, \mathbf{x}^* , were generated using the same initial state values, model parameters and EMV algorithm settings as those used by Raïssi et al.: $x_1(t_0) = 50$, $x_2(t_0) = 50$, $p_1 = 1$, $p_2 = 0.01$, $p_3 = 1$, $p_4 = 0.02$, $\alpha = 0.005$, $h = 0.005$ and $k = 4$ for $0 \leq t \leq 7$. Three initial data measurements were generated by adding random uncertainty to the true state values in order to create interval bounds at discrete time points, far fewer than the $N \geq 1,000$ measurement time points used in prior literature involving this model [8,15]. The second state, x_2 , was assumed to be unmeasurable while for the first state, x_1 , measurements were generated by adding error intervals as follows: $\hat{x}_1(t_i) = x_1^*(t_i) + \epsilon_i$, where $\epsilon_i = [-8.2190, 13.6065]$, $[-11.3067, 14.9691]$ and $[-7.6254, 10.5414]$ at $t_i = \{2, 4 \text{ and } 6\}$, respectively.

The assigned task is to determine at what times additional measurements would provide useful information with regards to the previously defined metrics and how many measurements would be beneficial. It was assumed that the initial conditions of both populations and parameters p_1 and p_3 were exactly known. We first wish to estimate the set of parameters p_2 and p_4 , along with the range of the unmeasured state x_2 for $0 \leq t \leq 7$,

that are consistent with the uncertain measurements of x_1 .

Initial parameter and state estimation

Bounded estimates of parameters p_2 and p_4 and states x_1 and x_2 were calculated using the initial measurements $\mathcal{Y} := \{\hat{x}_1(t_i); i = 1, 2, 3\}$. Parameter estimation was performed assuming an *a priori* search area of $[-1, 1]$ for both p_2 and p_4 and indeterminate boxes were bisected until a minimum box width of $\epsilon = 10^{-5}$ was obtained. This resulted in the generation of ~ 20 k indeterminate and feasible boxes shown in Figure 2 where no distinction is made between the two box types. Each box was then used in the EMV algorithm to produce the estimated state bounds, \mathbf{x}_{est} , shown in Figure 3 where \mathbf{x}^* are the grey waveforms, \hat{x}_i are the intervals and \mathbf{x}_{est} are the black dashed waveforms.

Estimating candidate measurements

The initial data measurements were compared to the estimated state bounds for x_1 to generate the interval set \mathbb{T} from which the candidate time points will be selected. Here, $k_p = 10$ candidate time points were chosen from within \mathbb{T} , namely $\mathcal{T} = \{1.25, 1.5, 1.75, 2.25, 2.5, 2.75, 3, 3.25, 3.5, 3.75\}$ as indicated in Figure 3. The corresponding values of R_j and C_j were estimated for each candidate time point t_j using the approach described above. The ranges R_j were shifted along the estimated state bounds for each corresponding t_j using $r = 15$ steps, where r was determined empirically to obtain curve fits with large R^2

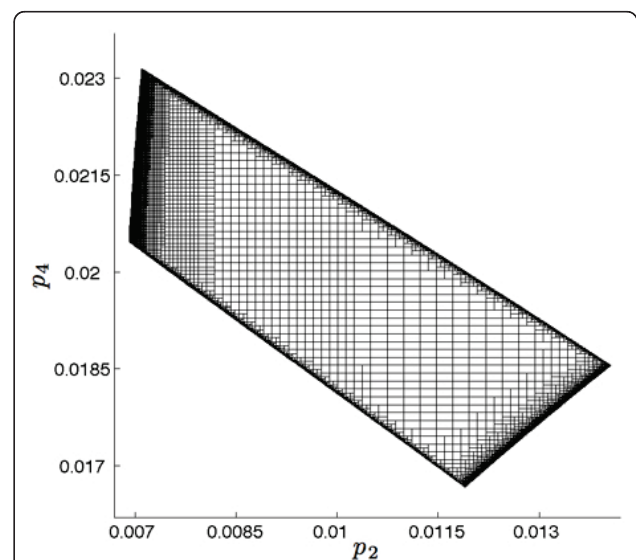


Figure 2 Initial parameter estimate. This figure shows the feasible and infeasible boxes in the parameter space that result from the SIMA algorithm. No distinction between feasible and infeasible is shown.

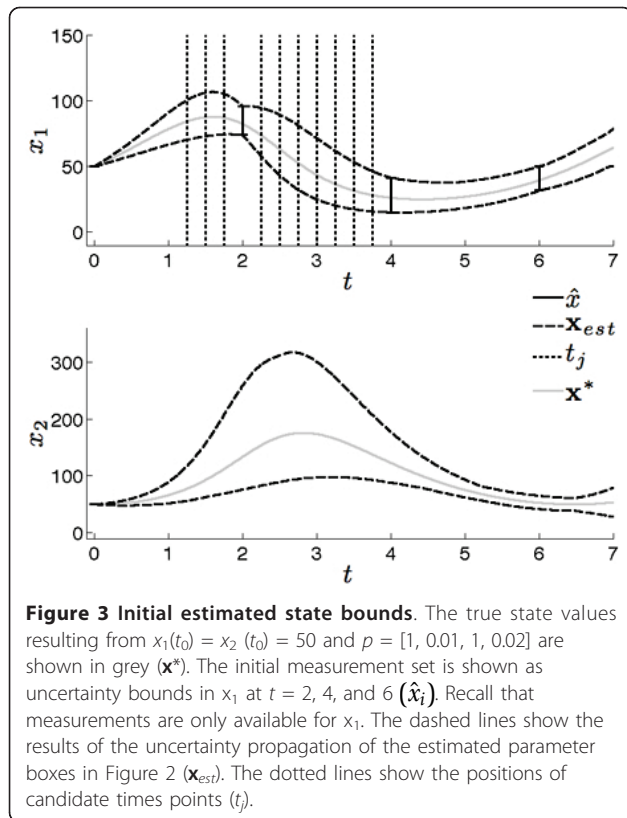


Figure 3 Initial estimated state bounds. The true state values resulting from $x_1(t_0) = x_2(t_0) = 50$ and $p = [1, 0.01, 1, 0.02]$ are shown in grey (\mathbf{x}^*). The initial measurement set is shown as uncertainty bounds in x_1 at $t = 2, 4,$ and 6 ($\hat{\mathbf{x}}_i$). Recall that measurements are only available for x_1 . The dashed lines show the results of the uncertainty propagation of the estimated parameter boxes in Figure 2 (\mathbf{x}_{est}). The dotted lines show the positions of candidate times points (t_j).

values. Bounded parameter estimations were performed for the $k_p \times r = 150$ shifted candidate measurements. The estimated parameter volumes were fit to quadratic curves with resulting R^2 values greater than 0.99. We were then able to identify an estimate of the center point that maximized this curve.

Combining time points

We were able to establish independence between candidate time points by showing that the brute-force estimates using all possible permutations and the intersected parameter sets cover identical parameter regions. The brute-force combinations and the intersections of parameter sets for all combinations of two candidate time points were compared and found to produce both the same parameter volumes and parameter bounds with a tolerance of 10^{-12} . Parameter intersections were then computed for combinations of up to $k_c = 5$ candidate time points. An example parameter intersection is shown in Figure 4 where the parameter estimates of $t_2 = 1.5$ and $t_6 = 2.75$ were combined. The parameter box colors correspond as follows: dark grey for \mathcal{P}_2 , which corresponds to t_2 , light grey for \mathcal{P}_6 , which corresponds to t_6 , and black for the brute-force combination which is used to depict the intersected parameter space.

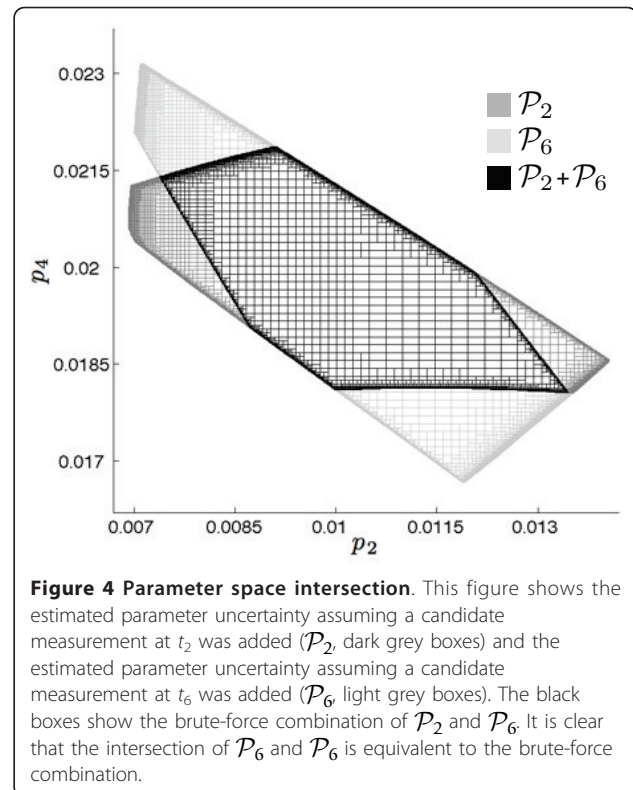


Figure 4 Parameter space intersection. This figure shows the estimated parameter uncertainty assuming a candidate measurement at t_2 was added (\mathcal{P}_2 , dark grey boxes) and the estimated parameter uncertainty assuming a candidate measurement at t_6 was added (\mathcal{P}_6 , light grey boxes). The black boxes show the brute-force combination of \mathcal{P}_2 and \mathcal{P}_6 . It is clear that the intersection of \mathcal{P}_2 and \mathcal{P}_6 is equivalent to the brute-force combination.

Estimates of state bounds were computed from the intersected parameter sets. An example estimate of state bounds is shown in Figure 5 for the parameter intersection of t_2 and t_6 . The underlying state values \mathbf{x}^* are the solid grey waveforms, the combined estimated state bounds \mathbf{x}_c are the solid black waveforms and the estimated state bounds \mathbf{x}_2 and \mathbf{x}_6 , corresponding to the results obtained from adding candidate measurements at t_2 and t_6 , respectively, are the dashed black and dashed grey waveforms, respectively. The decrease in uncertainty for state x_2 during $1 \leq t \leq 4$ is caused by the removal of the non-overlapping parameter regions.

Applying metrics

We tested whether the estimated candidate measurements generated by our algorithm could effectively be used to predict where the most appropriate measurements should be placed to reduce model uncertainty. With this in mind, we generated a set of true measurements at each candidate time point using the underlying state values, \mathbf{x}^* , as the true center points, C^* . We consider estimates obtained from measurements characterized by the true center points to be *ground truth*, corresponding to the best estimate of the measurement at a specific candidate time point. The metric results for estimates using the true center points C^* are used as a

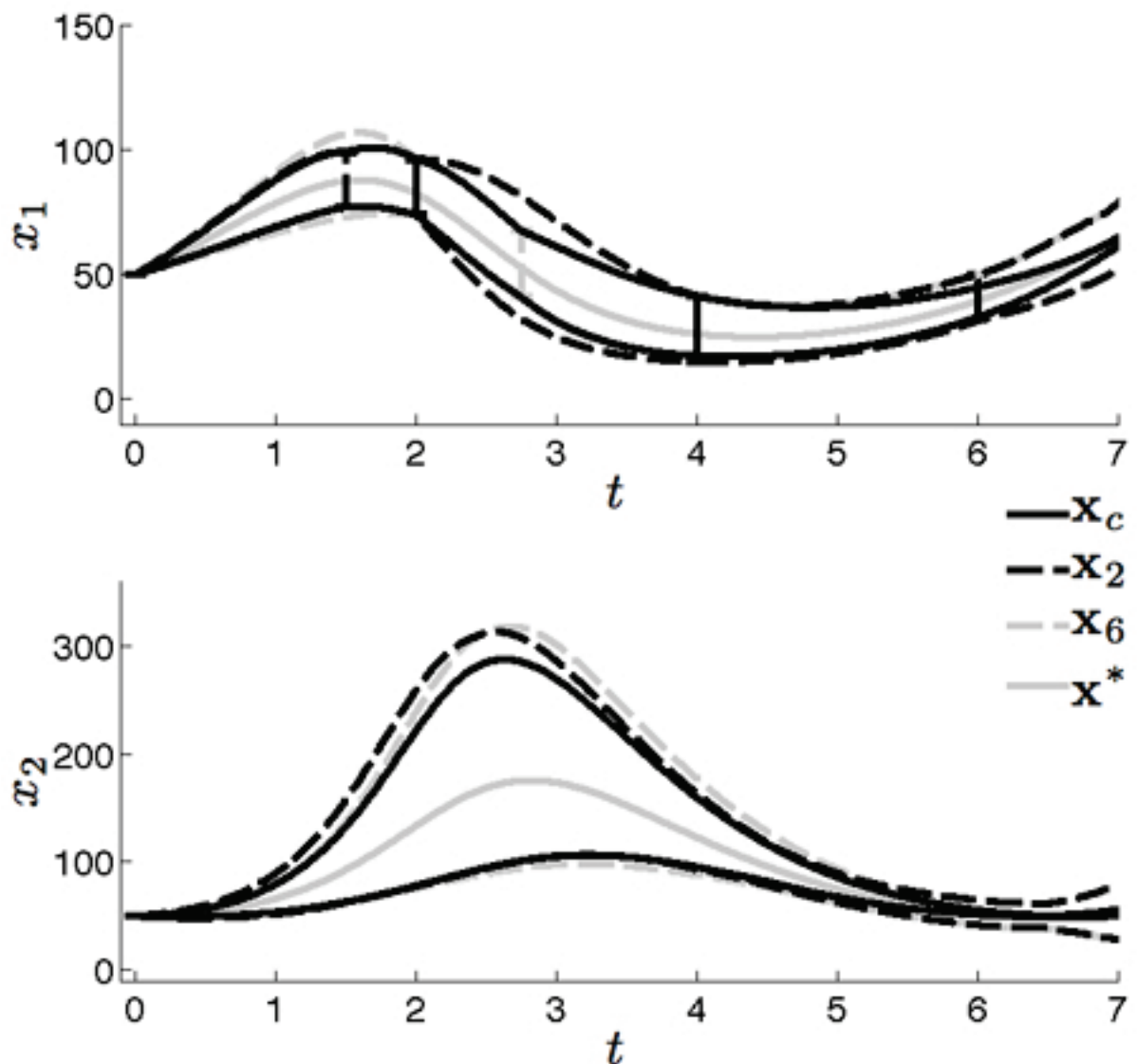


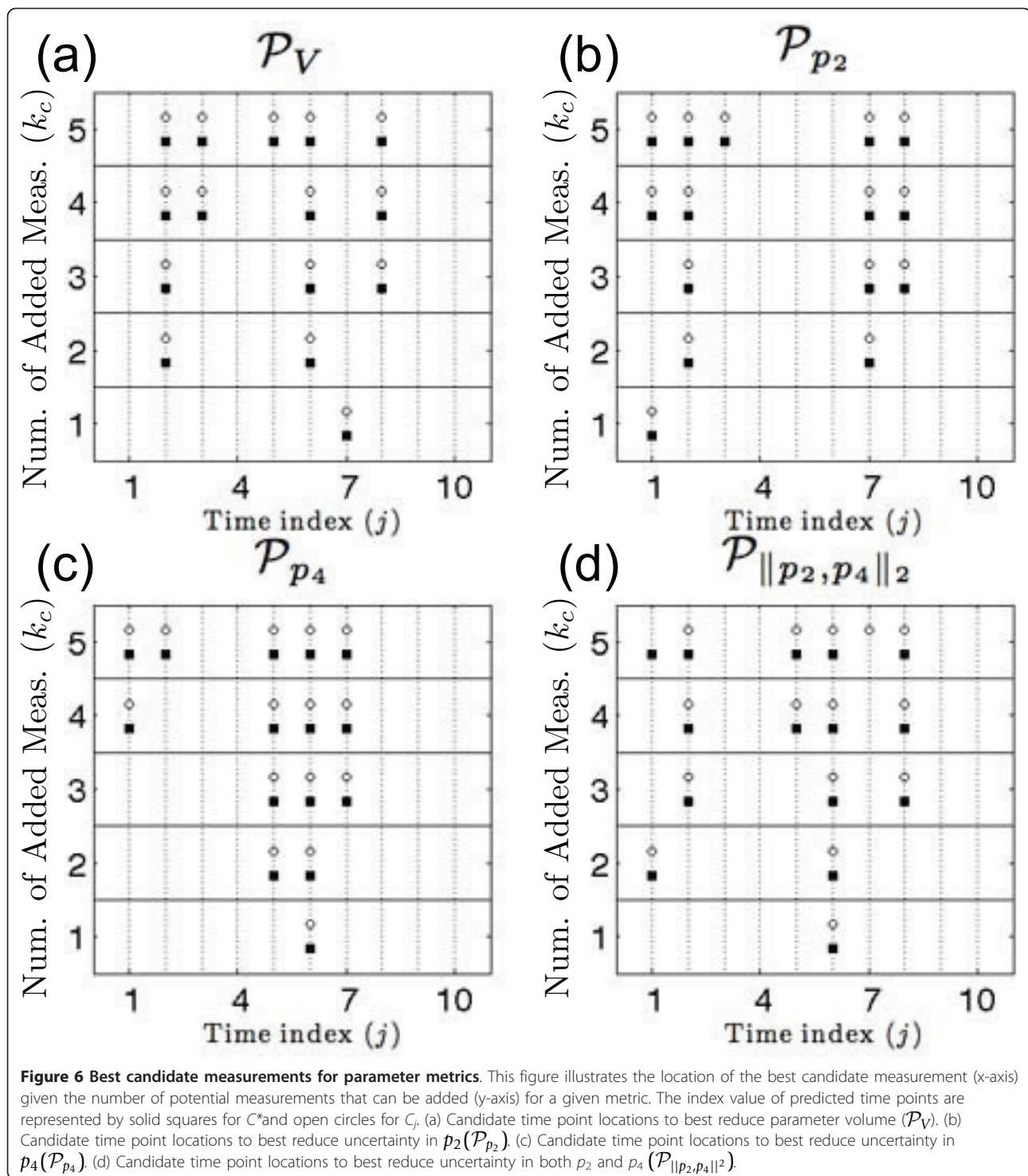
Figure 5 Combination of estimated state bounds. This figure shows the estimated state bounds assuming a candidate measurement at t_2 was added (x_2 , dashed black lines) and the estimated state bounds assuming a candidate measurement at t_6 was added (x_6 , dashed grey lines). The estimated state bounds for the combined candidate measurements, x_c , are the black lines, while the underlying true state values, x^* , are the solid grey lines.

reference and compared to the results obtained when using our estimated center points C_j .

Parameter information

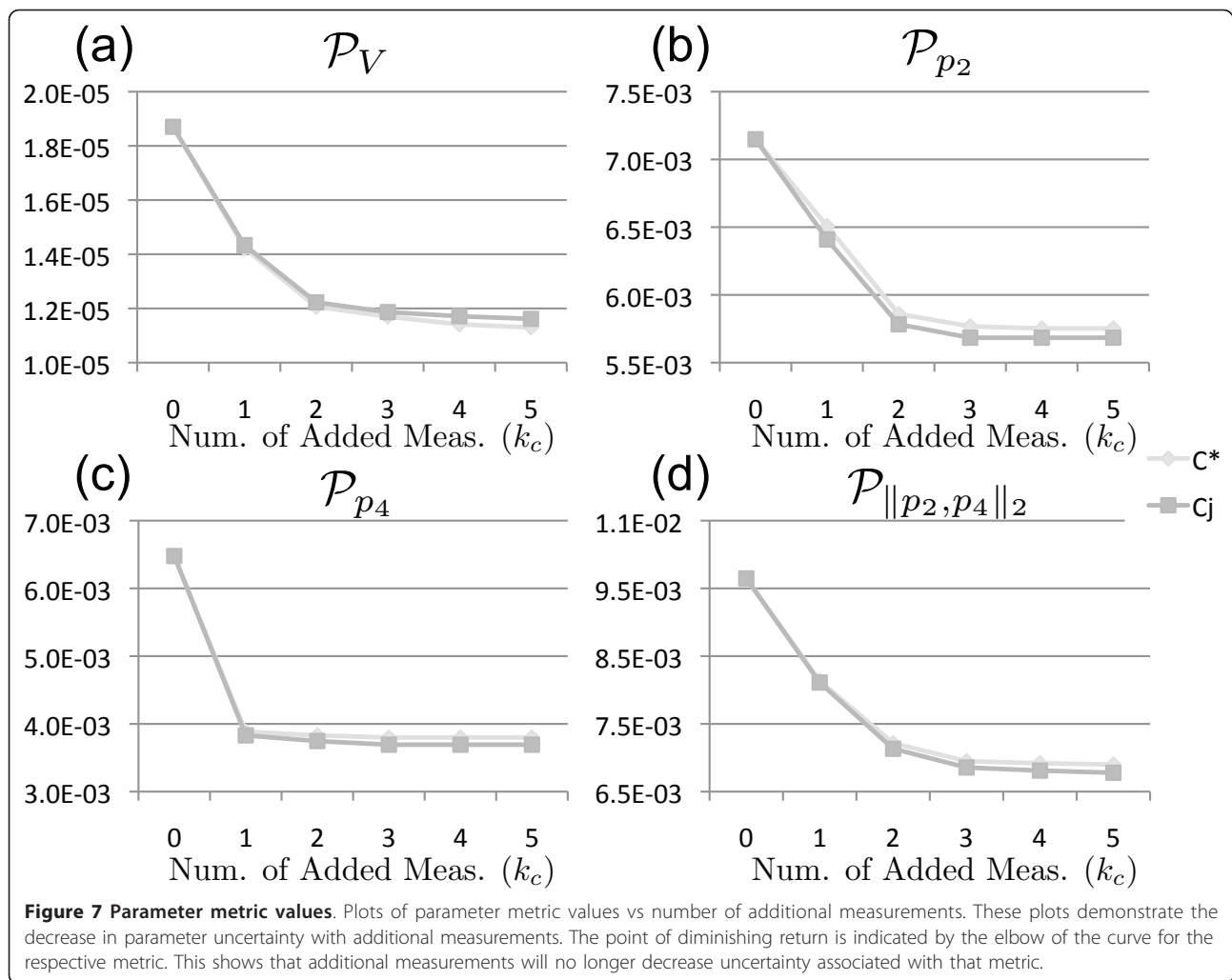
The prediction of the best time point locations, given the set of candidate measurements, for several parameter metrics are shown in Figure 6 when using center points C^* (solid squares) and C_j (open circles). This figure shows the best candidate measurement time point locations relative to the index of t_j for the parameter volume metric, \mathcal{P}_V , individual unknown parameter bounds, \mathcal{P}_{p_2} and \mathcal{P}_{p_4} , and combination of parameter

bounds, $\mathcal{P}_{||p_2, p_4||^2}$. Consider the design approach when there are only enough resources for a single additional measurement. Selecting a design to minimize the uncertainty of parameter p_2 (Figure 6b) would suggest placing a measurement at time $t_1 = 1.25$. However, to minimize the uncertainty of parameter p_4 a measurement at time $t_6 = 2.25$ would be more beneficial. If there are resources available for three additional measurements they would best be placed at times $t_2 = 1.5$, $t_6 = 2.75$, and $t_8 = 3.25$ to obtain additional information on both unknown parameters. We emphasize the established consistency between the best candidate time points



selected based on C^* and the best candidate time points selected based on our estimate C_j . The only inconsistent prediction between center points C^* and C_j occurs when applying the $\mathcal{P}_{\|p_2, p_4\|_2}$ metric for a combination of $k_c = 5$ time points, which results in a single time point difference.

The point at which additional measurements will not provide any additional information about the system can be predicted by observing the metric values for combinations of time points. This is especially beneficial for conserving resources that would otherwise be spent on experiments that yield no new information. The values



of the four parameter metrics are shown in Figure 7 as functions of the number of additional measurements. Using this information, an experimental designer could determine the desired number of additional measurements to collect without wasting resources. Consider selecting a set of measurements to reduce uncertainty for parameter p_4 . Estimating the impact of adding multiple measurements leads to the clear conclusion that a single additional measurement is all that is required. Similarly, reducing the uncertainty of the consistent parameter set volume may require 2 or 3 additional measurements. These metric value curves can be combined with cost functions to determine a design that efficiently utilizes experimental resources.

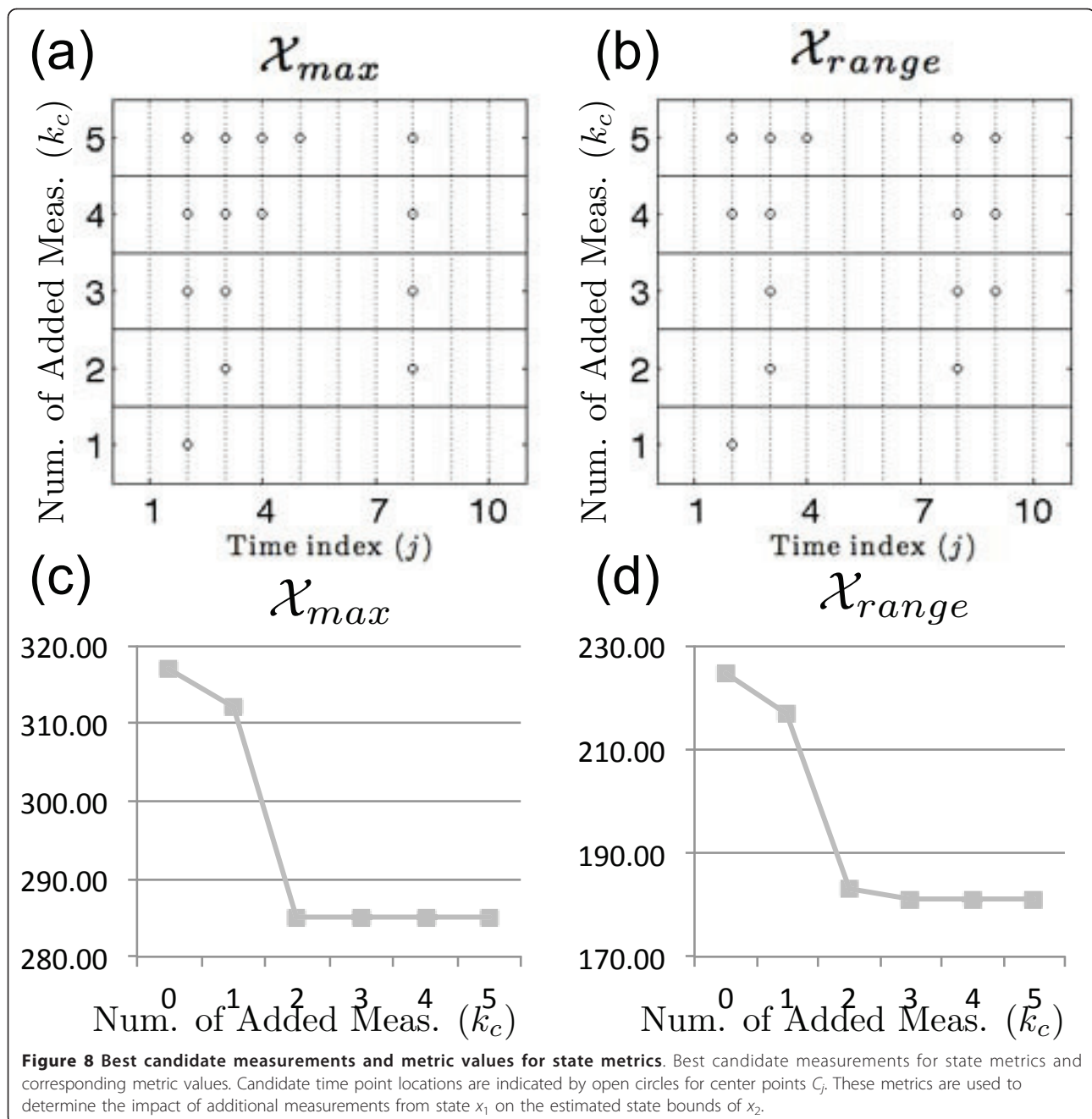
State information

Two metrics were applied to the unmeasured state, x_2 , to determine how its uncertainty is impacted when candidate measurements are applied to state x_1 using center points C_j . The first metric, \mathcal{X}_{max} , was used to select

candidate time points that would minimize the overall maximum value of x_2 . The second metric, \mathcal{X}_{range} , determines which candidate measurements will minimize the maximum uncertainty of x_2 over the simulation time $0 \leq t \leq 7$. The best time point locations and corresponding metric values are presented in Figure 8. Candidate measurement locations are fairly similar for the two metrics with \mathcal{X}_{max} slightly favoring candidate measurements located at earlier time points. A dramatic increase in information can be seen for both metrics when increasing from a single additional measurement to a combination of two measurements (Figure 8c-d). Little knowledge is gained when adding three or more measurements when compared to that gained from two additional measurements.

Comparison with FIM D-optimality

Scalar metrics of the Fisher Information Matrix (FIM) are often used to perform experimental design for many conventional problems [1,23,25]. We compared the



results of our set-based experimental design approach to results obtained using the D-optimality metric of the FIM. We did this to show how statistical assumptions that are often made to calculate the FIM could potentially impact the results when performing experimental design for biological processes. As stated previously, the number of measurements obtained for biological systems is very limited [4]. These data points are used to impose unwarranted statistics on the uncertainty, which are then used to calculate the FIM. Consider the scenario often encountered when quantifying biological

systems where resources are available for only four replicates of a given experiment, i.e. only four data points are generated for a given sample time t_i . The sets $\{74, 75, 80, 95\}_1$, $\{74, 80, 89, 95\}_2$ and $\{74, 89, 94, 95\}_3$ show three likely data sets containing four data points from experimental replicates for sample time t_i . All sets show data in the interval range 74 to 95. The small sample size of each set, however, implies that meaningful statistics of the uncertainty are difficult to obtain. In fact, each set has distinctively different means, with μ_1 , μ_2 , and μ_3 corresponding to 81, 84.5, and 88,

respectively. Given that the use of the FIM inherently assumes the use of Gaussian distributions [26], we use our results below to assess how these imposed Gaussian distributions, with their potentially different means, impact the decisions associated with experimental design.

We looked at three possible Gaussian distributions for each of the original measurement times, $t_i = \{2, 4, \text{ and } 6\}$, that could result from having small numbers of data samples (Figure 9a). Each distribution is characterized by its mean, $\mu_{t_i,s}$, and variance, $\sigma_{t_i,s}^2$. The variable t_i represents one of the original measurement time points and the variable s corresponds to the position of the distribution, i.e. $s = 1$ for shifted to the left, $s = c$ for shifted to the center, and $s = r$ for shifted to the right. All variances, $\sigma_{t_i,s}^2$, were calculated such that the distribution had a probability of 0.9 over the original interval uncertainty range. This ensures that each distribution, even though they have different means, has the same probability of producing population values over the uncertainty interval range.

We calculated the Maximum Likelihood (ML) estimate of the parameters [27] for the nine possible combinations of these distributions given the three initial measurement time points, $t_i = \{2, 4, \text{ and } 6\}$,

$$\hat{\theta}_{ML}^{\{s_2,s_4,s_6\}} = \min_{\theta} \sum_{t_i} \frac{1}{\sigma_{t_i,s}^2} (x_1(t_i, \theta) - \mu_{t_i,s})^2, \quad (7)$$

where S_{t_i} corresponds to the distribution type at time t_i . For example $\hat{\theta}_{ML}^{\{l_2,r_4,c_6\}}$ is the ML estimation resulting from using the left shifted distribution at time $t_1 = 2$, the right shifted distribution at time $t_2 = 4$, and the center distribution at time $t_3 = 6$. We computed the sensitivity matrix, S , using the method outlined in [28] by solving the ODE

$$\dot{S} = JS + A, \quad (8)$$

in combination with (6). Here, the $(i, j)^{\text{th}}$ element of these variables are $S_{i,j} = \partial x_i / \partial \theta_j$, $J_{i,j} = \partial f_i / \partial x_j$ and $A_{i,j} = \partial f_i / \partial \theta_j$. The FIM was then calculated as

$$FIM = \sum_{t_i \in I} \frac{1}{\sigma_{t_i,s}^2} \frac{\partial x_1(t_i)}{\partial \theta^T} \frac{\partial x_1(t_i)}{\partial \theta} \Bigg|_{\hat{\theta}_{ML}^{\{s_2,s_4,s_6\}}}, \quad (9)$$

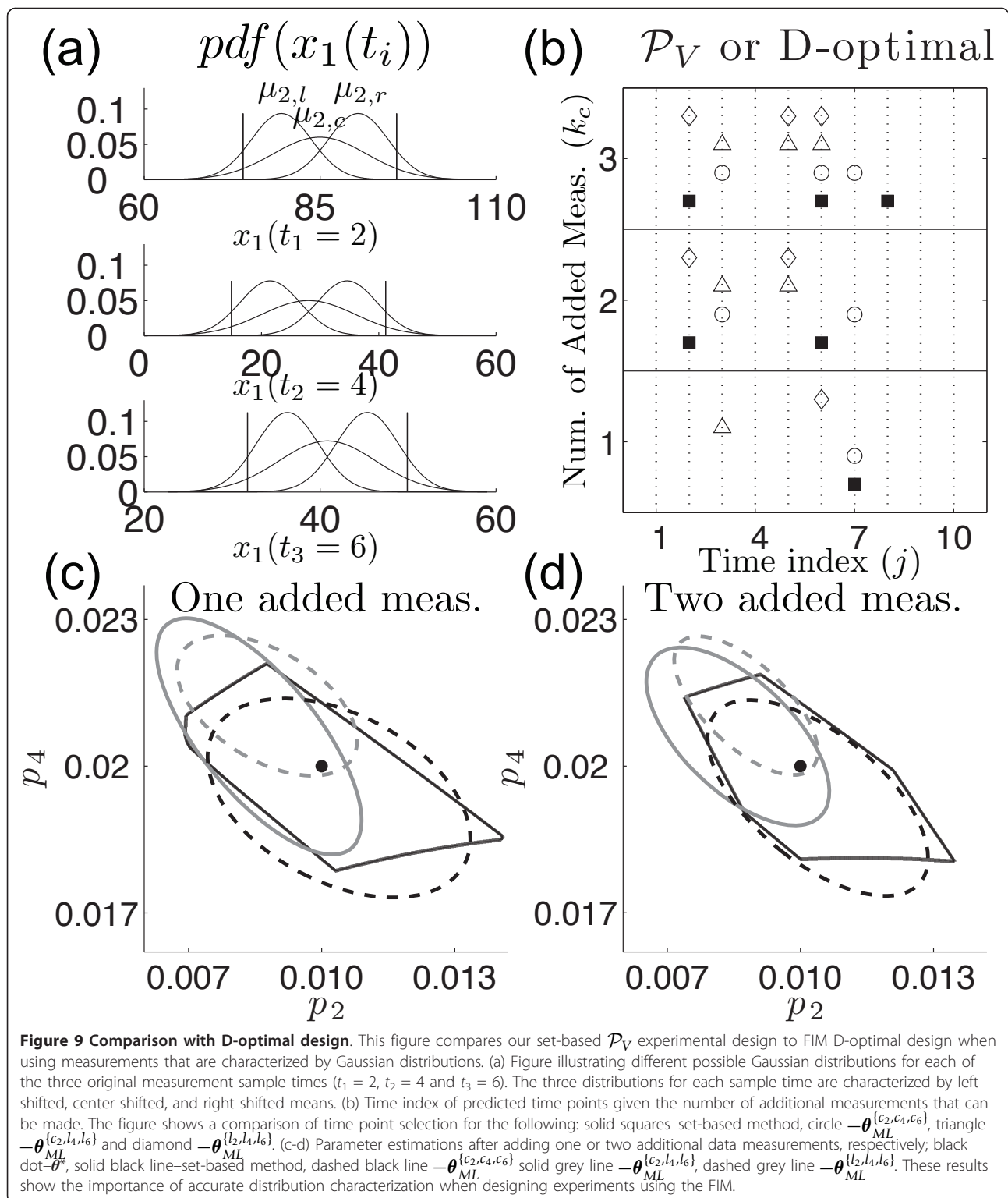
where I is the set of original measurement time points $\{2, 4 \text{ and } 6\}$ in addition to the subset of candidate time points, t_j , being evaluated, e.g. $\mathcal{J} = \{2, 2.75, 3.5, 4, 6\}$ where 2.75 and 3.5 would be the two candidate time points being evaluated. The variances at candidate time points were characterized in a way that was consistent with our set-based approach. The variance $\sigma_{t_j,s}^2$ for candidate

time point t_j was selected as the larger of the two variances of the adjacent initial measurements.

We computed D-optimal designs for the 9 distribution combinations and compared the selected candidate time points with our set-based method. The prediction of the best time point locations, given the set of candidate measurements, for our method (solid squares) and several D-optimal designs (circle $-\theta_{ML}^{\{c_2,c_4,c_6\}}$, triangle $-\theta_{ML}^{\{c_2,l_4,l_6\}}$ and diamond $-\theta_{ML}^{\{l_2,l_4,l_6\}}$) are shown in Figure 9b. The fluctuations in time point selection show that D-optimality is sensitive to our ability to correctly characterize the distributions of the initial data measurements, i.e. correctly characterizing the mean. Figure 9c and 9d show the corresponding parameter estimations for our method (solid black line) and the 95% confidence ellipsoids of D-optimal designs (dashed black line $-\theta_{ML}^{\{c_2,c_4,c_6\}}$, solid grey line $-\theta_{ML}^{\{c_2,l_4,l_6\}}$, dashed grey line $-\theta_{ML}^{\{l_2,l_4,l_6\}}$ after adding one and two measurements, respectively. The true parameter values are indicated with a point at (0.01,0.02). We are able to conclude based on these results that the selection of the time points for additional measurements, along with the assessment of the parameter uncertainty, changes depending on the characterization of the probability associated with the measurement uncertainty. Mischaracterization of the probability distribution is particularly possible when working with few data points, as is the case when modeling biological systems. This emphasizes the utility of our set-based experimental design approach. We also note that Figure 9c and 9d show that the resulting parameter uncertainty calculated using the FIM approach can result in an under or over estimation of the parameter range, depending on the characterization of the measurement uncertainty. This could be an important limitation in FIM experimental design approaches if one was interested in metrics related to absolute values of the parameter uncertainty (maximum value) instead of the relative change (minimum volume).

Conclusions

Developing accurate models is crucial for understanding, predicting and ultimately controlling biological processes. The limitation of costly resources and lengthy experiments associated with the study of biological systems promotes an experimental design approach for model development. Stochastic experimental design methods rely on correctly characterizing the distribution of uncertainty in the model, often requiring a large number of data measurements. This requirement is difficult to fulfill for many biological systems and alternative set-based experimental design approaches are more appropriate in these situations. In addition to the method used to characterize uncertainty, biological



interpretations of experimental design metrics are important because they provide a logical link between physical resources and mathematical constructs.

We have developed a novel experimental design framework using bounded-error methods and biologically relevant design metrics to select desirable time point

locations where additional measurements will be collected for the purpose of improving resource allocation for biological experiments. Our method propagates the uncertainty resulting from a small collection of data measurements, which may contain information for only a subset of the model states, through time to estimate parameter and state bounds for a given system model. We used these bounded-error results to estimate candidate measurement time points, center points and ranges. We proposed a method for combining candidate time points and present several biologically meaningful design metrics.

Measurement estimation is an important component of this method. We used a set-based approach to estimate measurements at time points where no information was available. We were able to estimate measurement bounds at candidate time points by combining information from the initial data measurement bounds with the estimated state bounds generated by the EMV algorithm. Our method resulted in a good estimate when compared to true measurements for the purpose of identifying where additional measurements should take place. The granularity of candidate time points can be made as fine as desirable at the cost of additional computation time. The computational expense to search all possible time points may make identifying globally optimal time point locations impractical using this method. However, the accuracy of when measurements are collected during biological experiments is often on the order of minutes, hours or days and locally optimal time points from an experimentally feasible set of time points is often sufficient.

The ability to estimate the effects of adding measurements at multiple time points is often desirable. A brute force method to explore all combinations of time points is computationally expensive. However, we found that the parameter estimation for a combination of time points can be directly obtained by intersecting the individual estimated parameter spaces. Estimated state bounds can then be determined using the intersected parameter space. The experimenter can determine when additional measurements will provide little or no additional information by exploring the effects of adding multiple measurements and will not needlessly spend limited resources on experiments that yield no additional information.

The framework presented here can be used to predict at what time additional measurements should be made to maximize information based on biologically relevant metrics and to determine the number at which additional measurements being to provide insignificant information. Problems of this sort are often faced by biologists when modeling biological processes. Selecting an appropriate metric is made more straightforward by associating it

with biologically relevant information. For example, the uncertainty of a parameter may be associated with specific characteristics of an engineered enzyme, while the limitations on the uncertainty of estimated state bounds can provide critical bounds on unmeasured component concentrations, allowing systems to maintain chemical and physiological phenotypes.

Acknowledgements

The authors acknowledge financial support from NCSU startup funds.

Authors' contributions

SM designed the study and prepared the manuscript. CW participated in the design and in revising the draft. Both authors read and approved the final manuscript.

Competing interests

The authors declare that they have no competing interests.

Received: 2 August 2011 Accepted: 21 March 2012

Published: 21 March 2012

References

1. Pronzato L: **Optimal experimental design and some related control problems.** *Automatica* 2008, **44**(2):303-325.
2. Markov S: **On the use of computer algebra systems and enclosure methods in the modelling and optimization of biotechnological processes.** *International Journal Bioautomation* 2005, **3**:1-9.
3. Balsa-Canto E, Alonso A, Banga J: **Computational procedures for optimal experimental design in biological systems.** *Systems Biology, IET* 2008, **2**(4):163-172.
4. Kreutz C, Timmer J: **Systems biology: experimental design.** *FEBS J* 2009, **276**(4):923-942.
5. Walter E, Piet-Lahanier H: **Estimation of parameter bounds from bounded-error data: a survey.** *Math Comput Simul* 1990, **32**(5-6):449-468.
6. Pronzato L, Walter E: **Experiment design for bounded-error models.** *Math Comput Simul* 1990, **32**(5-6):571-584.
7. Jaulin L: **Nonlinear bounded-error state estimation of continuous-time systems.** *Automatica* 2002, **38**(6):1079-1082.
8. Raissi T, Ramdani N, Candau Y: **Set membership state and parameter estimation for systems described by nonlinear differential equations.** *Automatica* 2004, **40**(10):1771-1777.
9. Pronzato L, Walter E: **Robust experiment design via maximin optimization.** *Math Biosci* 1988, **89**(2):161-176.
10. Pronzato L, Walter E: **Experiment design in a bounded-error context: comparison with D-optimality.** *Automatica* 1989, **25**(3):383-391.
11. Pronzato L, Walter E: **Minimum-volume ellipsoids containing compact sets: application to parameter bounding.** *Automatica* 1994, **30**(11):1731-1739.
12. Hasenauer J, Waldherr S, Wagner K, Allgower F: **Parameter identification, experimental design and model falsification for biological network models using semidefinite programming.** *Systems Biology, IET* 2010, **4**(2):119-130.
13. Chernousko F: **Ellipsoidal state estimation for dynamical systems.** *Nonlinear Analysis* 2005, **63**(5-7):872-879, Invited Talks from the Fourth World Congress of Nonlinear Analysts (WCNA 2004).
14. Combastel C: **A State Bounding Observer for Uncertain Non-linear Continuous-time Systems based on Zonotopes.** *Decision and Control, 2005 and 2005 European Control Conference. CDC-ECC'05. 44th IEEE Conference on* 2005, 7228-7234.
15. Lin Y, Stadtherr MA: **Guaranteed state and parameter estimation for nonlinear continuous-time systems with bounded-error measurements.** *Ind Eng Chem Res* 2007, **46**(22):7198-7207.
16. Marvel SW, de Luis Balaguer MA, Williams CM: **Parameter Estimation in Biological Systems Using Interval Methods with Parallel Processing.** *8th International Workshop on Computational Systems Biology* Zürich, Switzerland; 2011, 129-132.
17. Moore RE: *Interval Analysis* Englewood Cliffs, NJ: Prentice-Hall; 1966.

18. Nedialkov NS, Jackson KR, Corliss GF: **Validated solutions of initial value problems for ordinary differential equations.** *Appl Math Comput* 1999, **105**:21-68.
19. Rihm R: **Interval Methods for Initial Value Problems in ODEs.** *Topics in validated computations: Proceedings of the IMACS-GAMM international workshop on validated computations* 1994, 173-208.
20. Neumaier A: *Interval Methods for Systems of Equations* Cambridge, UK: Cambridge University Press; 1990.
21. Jaulin L, Walter E: **Set inversion via interval analysis for nonlinear bounded-error estimation.** *Automatica* 1993, **29**(4):1053-1064.
22. Gropp W, Lusk E, Skjellum A: *Using MPI: portable parallel programming with the message-passing interface* Cambridge, MA, USA: MIT Press; 1994.
23. Vanrolleghem P: **Bioprocess Model Identification.** In *Advanced Instrumentation, Data interpretation, and Control of Biotechnological Processes.* Edited by: Impe JV, Vanrolleghem P, Iserentant D. Dordrecht, The Netherlands: Kluwer Academic Publishers; 1998:251-318.
24. Voit E, Chou IC: **Parameter estimation in canonical biological systems models.** *International Journal of Systems and Synthetic Biology* 2010, **1**:1-19.
25. Rodriguez-Fernandez M, Mendes P, Banga JR: **A hybrid approach for efficient and robust parameter estimation in biochemical pathways.** *Biosystems* 2006, **83**(2-3):248-265.
26. Faller D, Klingmüller U, Timmer J: **Simulation methods for optimal experimental design in systems biology.** *Simulation* 2003, **79**(12):717-725.
27. Bro R, Sidiropoulos ND, Smilde AK: **Maximum likelihood fitting using ordinary least squares algorithms.** *J Chemom* 2002, **16**(8-10):387-400.
28. Dickinson RP, Gelinias RJ: **Sensitivity analysis of ordinary differential equation systems A direct method.** *J Comput Phys* 1976, **21**(2):123-143.

doi:10.1186/1752-0509-6-21

Cite this article as: Marvel and Williams: **Set membership experimental design for biological systems.** *BMC Systems Biology* 2012 **6**:21.

**Submit your next manuscript to BioMed Central
and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

