

Selection of Perturbation Experiments for Model Discrimination

Ivayla Vatcheva¹, Hidde de Jong² and Nicolaas J.I. Mars¹

¹ Department of Computer Science, University of Twente

P.O. Box 217, 7500 AE Enschede, the Netherlands

² Institut National de Recherche en Informatique et en Automatique (INRIA) Rhône-Alpes

655, avenue de l'Europe, 38330 Montbonnot Saint Martin, France

Email: {ivayla, mars}@cs.utwente.nl, Hidde.de-Jong@inrialpes.fr

Abstract

In many situations a system is described by several competing models. In order to distinguish among the proposed models, further information about the behavior of the system is required. One way to obtain such information is to perform suitably chosen perturbation experiments. This paper introduces a method for the selection of optimal perturbation experiments for discrimination among a set of competing dynamical models. The models are assumed to have the form of semi-quantitative differential equations. The method employs an optimization criterion based on the entropy measure of information.

Introduction

Scientists and engineers are frequently faced with situations in which a given system can be described by several *competing models*. The predictions of the models match available observations about the system behavior obtained in one or more experiments. When analyzing the synthesis rate of a product in a catalyzed chemical reaction as a function of the partial pressures of the input substances, one often arrives at several equations that all satisfy a set of measurements (Swaan 1992). For the mitotic clock of early embryos, a dozen of models predicting the observed periodic behavior of the concentrations of key proteins have been suggested (Obeyesekere, Tucker, & Zimmerman 1992).

In order to identify which of the proposed models best describes the real setting, new observations have to be made. These can be obtained by performing supplementary *perturbation experiments* on the system. In a perturbation experiment the structure of the system and/or the experimental conditions are changed. An experiment discriminates between the competing models, if the predictions of some of the candidates, which have been properly modified to reflect the experimental change, fit the newly obtained data whereas others show a lack of fit. The problem of experiment selection for model discrimination can then be defined as the problem of selecting a perturbation that gives rise to observations matching the predictions of as few of the proposed models as possible.

The imprecision of measurements in the experiments,

and the complexity of the system to be understood, do not always permit detailed quantitative analyses to be performed. Both the lack of accurate and reliable measurements, and the approximate models of real world systems, appeal to a qualitative or a semi-quantitative approach to the model discrimination problem. We assume that the models are given in the form of semi-quantitative differential equations. Predictions, in the form of intervals for the model variables, are derived by means of semi-quantitative simulation techniques. Measurements are considered to be intervals as well.

This paper presents a method for the systematic choice of perturbation experiments for model discrimination. Experiments are selected on the basis of an entropy criterion suggested by Box & Hill (1967), which measures the information increment provided by each of the experiments. The concept of entropy as a discrimination criterion has also been used in the work of Reilly (1970), and Fedorov (1972) in statistics, and in the work of de Kleer & Williams (1987) and Struss (1994b) in model-based diagnosis. A novel aspect of our work is that we extend this concept to the case of perturbation experiments and to situations in which experimental systems are described by nonlinear, dynamical models.

The in-principle applicability of our approach is illustrated on a set of competing models of an oscillatory, second-order system. We will consider six models of a mass-spring system and illustrate the choice of suitable perturbations to discriminate between the models. The principles involved in this example are applicable to the investigation of more complex and less understood oscillating systems.

The presentation starts with a description of the problem of model discrimination. A number of basic concepts are introduced and the relationship between models and experiments is given. The criterion for choosing a maximally-discriminating perturbation is described in the next section, and embedded in a simple algorithm for the discrimination of a set of competing models. Next, the application of the method is illustrated on the example. The last section discusses limitations and extensions of our method, in the context of related work in statistics and model-based diagnosis.

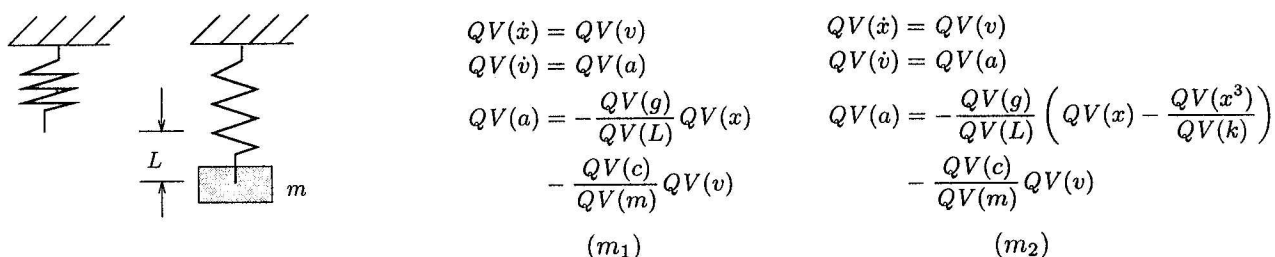


Figure 1: Damped mass-spring system and two QDEs m_1 and m_2 describing the system (for the notation, see Vatcheva & de Jong (1999)). The variables refer to the position x , velocity v , acceleration a , mass m , gravity constant g , initial spring elongation L and friction constant c . The following intervals complete the QDEs to SQDEs: $range(m) \in [2.95, 3.05]$, $range(g) = [9.83, 9.83]$, $range(L) = [5.8, 6.0]$, $range(c) = [0.3, 0.4]$. The initial values for the position and the velocity are $[0.9, 1.1]$ and $[0, 0]$, respectively. The constant k is specified by $range(k) = [6, 6]$.

Model discrimination by perturbation experiments

In this section the concepts of experimental system, perturbation experiment and model perturbation are introduced. Since the focus on the paper is on the method for model discrimination, we provide an intuitive explanation of these concepts rather than giving a well-founded formalization framework. Attempts to formalize discriminating tests can be found in (McIlraith 1994) and (Struss 1994a).

The systems we will be concerned with in this paper are (physical) systems controlled in experiments, also called *experimental systems*. An example of an experimental system is a cell culture allowed to grow under controlled environmental conditions, including nutrient supply and temperature. Control over an experimental system is achieved by creating and maintaining its structure and by regulating the experimental conditions under which the behavior of the system evolves.

Suppose a set M of models of a system being investigated in an experiment has been proposed. Let $p(m_i)$ be the a priori probability of model $m_i \in M$ to be the correct model of the system. The model probabilities can be derived from preliminary observations on the system behavior or theoretical considerations. If no prior knowledge about the relative plausibilities of the models exists, equal probabilities are assumed. We say that the models $m_i \in M$ are *competing*. M is assumed to be complete, that is, $\sum_{m_i \in M} p(m_i) = 1$. This may seem a strong assumption, but its practical consequences are limited as its violation can be tested.

In this paper we will model experimental systems by means of *semi-quantitative differential equations* (SQDEs), that is, qualitative differential equations (QDEs) enhanced with numerical information. The semi-quantitative information completing a QDE takes the form of numerical ranges added to landmarks and envelopes for monotonic function constraints (Berleant & Kuipers 1997). In this way, uncertainty about the exact values of parameters and the precise form of functional relations can be expressed. Fig. 1 shows two SQDEs describing a simple experimental system, a

damped mass-spring system. The models assume that the forces playing a role in the experiment are a spring force and a friction force, but they differ in the precise nature ascribed to the former. The initial ranges for the position x and velocity v are considered to be part of the model.

In order to distinguish between the models, additional information about the system is required. This information can be obtained by performing a suitably chosen *perturbation experiment*. In a perturbation experiment the system structure or the experimental conditions are modified. By allowing changes in the system structure we extend existing approaches to model discrimination which are only focused on changes in the system inputs, e.g. (Box & Hill 1967; Reilly 1970). The changes have to be reflected on the competing models in such a way that the operations on a model correspond with perturbations of the experimental system (Fig. 2).

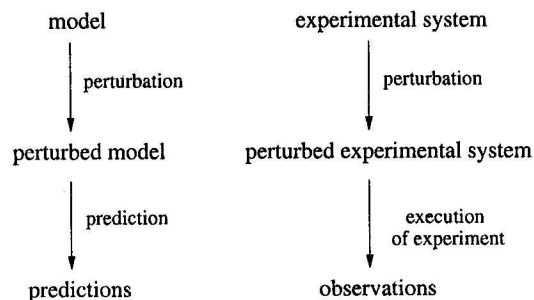


Figure 2: Correspondence between perturbation experiments and model perturbations.

In order to predict the consequences of a perturbation, we employ the semi-quantitative simulation techniques Q2 and Q3 (Kuipers 1994; Berleant & Kuipers 1997). Q2 and Q3 exploit the ranges of landmarks and the envelopes of monotonic function constraints in an SQDE to refine a qualitative behavior tree produced by QSIM. More specifically, they rule out qualitative behaviors or transform qualitative behaviors into *semi-quantitative behaviors* (SQBs) in which the qualitative values are annotated with numerical ranges. Fig. 3

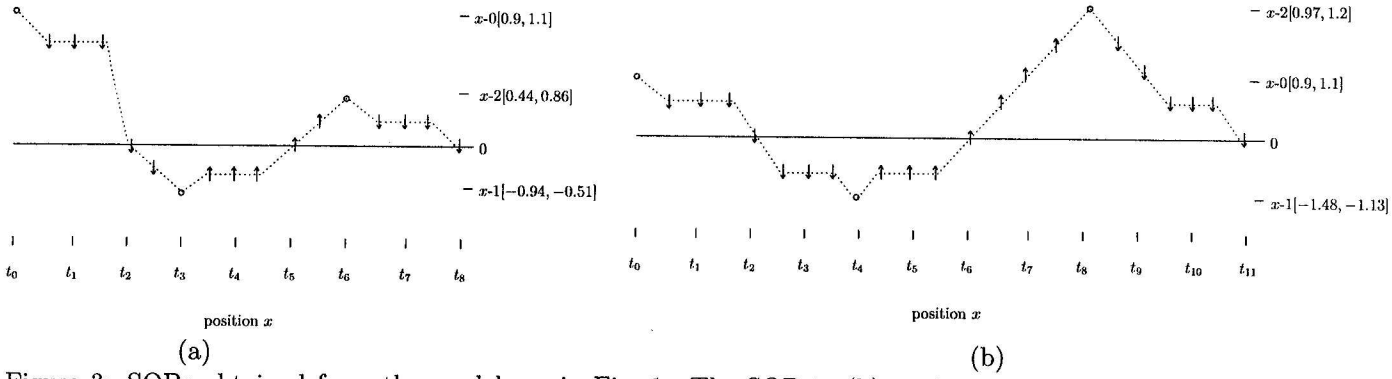


Figure 3: SQBs obtained from the model m_1 in Fig. 1. The SQB in (b) is obtained after a perturbation of m_1 , increasing the initial velocity from $[0, 0]$ to $[1.9, 2.1]$.

shows the oscillations predicted by model m_1 , before and after a perturbation that consists in releasing the object with nonzero initial velocity.

A semi-quantitative behavior is a prediction of the interval value of the variables at the *distinguished time-points*, the time-points at which some variable changes its qualitative value. For instance, the SQB in Fig. 3(a) shows that at t_3 , the time-point at which x reaches its maximum for the first time, the value of x lies in the interval $[-0.94, -0.51]$.

In addition to predictions of the value of a variable at a time-point, we might be interested in knowing the difference in value of a variable before and after a perturbation. Predictions of the relative interval value of variables can be obtained by subtracting the predicted interval values of the variable at corresponding distinguished time-points in the behavior before and after a perturbation, so-called *meaningful pairs of comparison* (de Jong & van Raalte 1999). As a consequence of the use of semi-quantitative information, these predictions may be weaker than necessary. We use the comparative analysis technique SQCA to obtain more precise predictions (Vatcheva & de Jong 1999). The information in Fig. 3 allows one to infer, by subtracting interval ranges, that the difference $\hat{x} - x$ at the pair of comparison $\{t_3, \hat{t}_4\}$ lies in the interval $[-0.97, -0.19]$. \hat{x} and \hat{t} refer to variables in the perturbed system. Application of SQCA refines this prediction by narrowing the relative value to $[-0.93, -0.19]$.

The amplitude in the behavior of the perturbed mass-spring system, or the difference in amplitude in the behaviors of the perturbed and unperturbed systems, are examples of *behavioral features* that help in discriminating competing models of an experimental system. Let \mathcal{P} be the sets of possible predictions of the interval values and relative interval values of the variables of an experimental system before and after a perturbation. A predicted behavioral feature is an interval value calculated from a set of predictions by means of an arithmetic function $f : \mathcal{P} \rightarrow \mathcal{I}(\mathbb{R})$.¹

¹ $\mathcal{I}(\mathbb{R})$ is the set of intervals with bounds in \mathbb{R} .

The function f may simply select a predicted value or relative value from the set of predictions \mathcal{P} , as in the case of a predicted amplitude. An example of a less trivial feature is the frequency of an oscillation, which can be calculated from the interval ranges of the distinguished time-points of two successive maxima. The concept of behavioral feature can be generalized to more complex features, in particular to qualitative features abstracted from the predicted behaviors of the systems. In the case of a mass-spring system, for instance, the mass could be increased to such an extent that the damped oscillation changes into an overdamped return to the rest state.

In order to be useful, predicted behavioral features need to correspond with *observed* behavioral features of the experimental system. That is, it must be possible to relate a predicted behavioral feature to some direct or indirect measurement of quantities of the system. As measurements will be assumed to have the form of *confidence intervals*, observed behavioral features are intervals.

The results of a perturbation experiment can be used to recompute the probabilities of the competing models. Models of which the predictions do not agree with the observations will have an a posteriori probability equal to 0. The model discrimination problem can now be intuitively stated as follows: find the perturbation experiment with values for the observed behavioral features that make a maximum number of models improbable. In the next section, we elaborate this intuition by means of an approach based on concepts from information theory.

Method for the selection of perturbation experiments

In order to maximally discriminate between a set of models, we will be interested in finding the perturbation yielding the highest increment in information (Box & Hill 1967). Consider a behavioral feature Y defined by some function f , mapping to intervals in a domain D . Consider a perturbation experiment $e \in E$, whose outcome yields a value $Y^e = [y^e - \epsilon/2, y^e + \epsilon/2]$ of the

behavioral feature, where y^e is the middle point of the interval Y^e and ϵ is the size of the confidence interval for Y . We can formulate the *information increment* of e as

$$\Delta H(e) = - \sum_{m_i \in M} p(m_i) \ln p(m_i) + \sum_{m_i \in M} p(m_i | Y^e) \ln p(m_i | Y^e), \quad (1)$$

where $p(m_i)$ and $p(m_i | Y^e)$ are the a priori and a posteriori probabilities of model m_i . The function ΔH reaches its maximum when the a posteriori probabilities of all competing models but one are 0. On the other hand, a minimal value is attained when the a posteriori probabilities are equal.

The $p(m_i | Y^e)$ s in (1) are not known, since they are determined by the outcome of the experiment. However, we can express the expected value of ΔH in terms of the probability distributions $g_i^{\{e,Y\}}$ of the behavioral feature Y . For brevity, g_i^e instead of $g_i^{\{e,Y\}}$ will be further used if no confusion about the behavioral feature being considered is possible. The value of Y predicted by model m_i under perturbation e is an interval $V_i^e \subseteq D$, with a probability distribution $g_i^e : D \rightarrow \mathbb{R}_{\geq 0}$ defined as follows

$$g_i^e(y) = \begin{cases} \frac{|[y - \frac{\epsilon}{2}, y + \frac{\epsilon}{2}] \cap V_i^e|}{|V_i^e|} & , [y - \frac{\epsilon}{2}, y + \frac{\epsilon}{2}] \cap V_i^e \neq \emptyset, \\ 0 & , [y - \frac{\epsilon}{2}, y + \frac{\epsilon}{2}] \cap V_i^e = \emptyset \end{cases}$$

where $|\cdot|$ denotes an interval length. $g_i^e(y)$ determines the probability of the empirically-determined value of Y to be $[y - \epsilon/2, y + \epsilon/2]$ if the model m_i is the correct model of the system. (2) can be replaced by the following equivalent expression, where the g_i^e s are defined as piece-wise linear functions:

$$g_i^e(y) = \begin{cases} \frac{y - V_i^e + \epsilon/2}{|V_i^e|} & , y \in [V_i^e - \epsilon/2, V_i^e + \epsilon/2], \\ \frac{\epsilon}{|V_i^e|} & , y \in [V_i^e + \epsilon/2, \bar{V}_i^e - \epsilon/2], \\ \frac{-y + \bar{V}_i^e + \epsilon/2}{|V_i^e|} & , y \in [\bar{V}_i^e - \epsilon/2, \bar{V}_i^e + \epsilon/2], \\ 0 & , y \notin [V_i^e - \epsilon/2, \bar{V}_i^e + \epsilon/2] \end{cases} \quad (3)$$

where V_i^e and \bar{V}_i^e denote the lower and the upper bound of V_i^e , respectively. Fig. 4(a) illustrates the function g_3^e for an experiment e consisting of replacing the object in the mass-spring system by a lighter object (e_4 in the next section). The behavioral feature considered is the interval value for the amplitude of the system and ϵ is taken to be 0.1.

Call the expected value of the information increment $\Delta J(e)$. By definition,

$$\Delta J(e) = \int_{y \in D} \Delta H(e) g^e(y) dy, \quad (4)$$

where

$$g^e(y) = \sum_{m_i \in M} p(m_i) g_i^e(y). \quad (5)$$

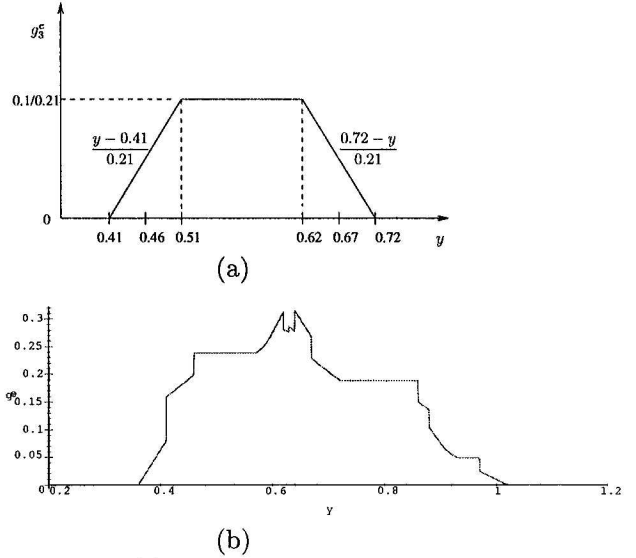


Figure 4: (a) A plot of the function g_3^e for the interval value of the amplitude of the mass-spring system in an experiment consisting of replacing the object by a lighter object (experiment e_4 in the next section). The prediction of model m_3 (see next section) perturbed according to this experiment is $V_3^e = [0.46, 0.67]$, and $\epsilon = 0.1$. (b) A plot of g^e defined for the same experiment and behavioral feature, and the six competing models given in Fig. 1 and Fig. 6.

Fig. 4(b) shows the plot of the function g^e for the experiment and the behavioral feature mentioned above, and the predictions of the six competing models given in Fig. 1 and Fig. 6.

By substituting the expression for $\Delta H(e)$ in (4) we get,

$$\Delta J(e) = \sum_{m_i \in M} p(m_i) \int_{y \in D} g_i^e(y) \{ \sum_{m_j \in M} p(m_j | Y) \ln p(m_j | Y) - \sum_{m_j \in M} p(m_j) \ln p(m_j) \} dy, \quad (6)$$

where $Y = [y - \epsilon/2, y + \epsilon/2]$ and

$$p(m_j | Y) = \frac{p(m_j) g_j^e(y)}{g^e(y)} \quad (7)$$

via the Bayes rule. Combination of (6) and (7) gives, after algebraic simplification,

$$\Delta J(e) = \sum_{m_i \in M} p(m_i) \int_{y \in D} g_i^e(y) \ln \frac{g_i^e(y)}{g^e(y)} dy. \quad (8)$$

If several behavioral features Y_1, \dots, Y_k are taken into account, the formula in (8) remains unchanged, except for replacing y by \mathbf{y} , D by $\mathbf{D} = D_1 \times \dots \times D_k$, the distributions $g_i^e(y)$ by joint probability distributions $g_i^{\{e, Y_1, \dots, Y_k\}}(\mathbf{y})$, and the integral by a multiple integral.

Intuitively, the criterion now tries to maximize the non-overlapping parts of the k -dimensional boxes in \mathbf{D} that are defined by the values for the behavioral features predicted by the m_i s.

Denote with $\Delta J(e, Y_1, Y_2)$ the expected increment of information of the experiment e when the behavioral features Y_1 and Y_2 are both taken into account, and $\Delta J(e, Y_2|Y_1)$ the expected increment of information of measuring Y_2 if Y_1 has been measured. The following properties of ΔJ are easily provable (Fedorov 1972):

1. $\Delta J(e) \geq 0$;
2. $\Delta J(e, Y_1) + \Delta J(e, Y_2|Y_1) = \Delta J(e, Y_1, Y_2)$;
3. $\Delta J(e, Y_2|Y_1) \leq \Delta J(e, Y_2)$ with equality iff Y_1 and Y_2 are independent, that is, $g_i^{\{e, Y_1, Y_2\}}(y_1, y_2) = g_i^{\{e, Y_1\}}(y_1)g_i^{\{e, Y_2\}}(y_2)$.

From the last two properties, it is readily seen that if the measurements in the experiment are independent, an experiment in which a set of quantities are measured is as informative as performing the same experiment a couple of times, each time measuring a single quantity.

The optimal next perturbation experiment to perform is the one for which (8) is maximized. Intuitively, the criterion favors experiments for which the corresponding model perturbation results in predicted intervals of the behavioral feature that overlap as little as possible. On average, less overlap of the intervals will increase the chance that a measurement of the feature discriminates between the models. This can be illustrated by means of the predictions of the relative amplitude by three alternative equiprobable models of the mass-spring system. Consider the case of a perturbation e_1 replacing the medium with almost a frictionless medium (setting c to 0), and a perturbation e_3 increasing the mass to [11.95, 12.05] (see Fig. 5). Notice that the expected information increment is higher for the second experiment, as the predicted intervals have less overlap.

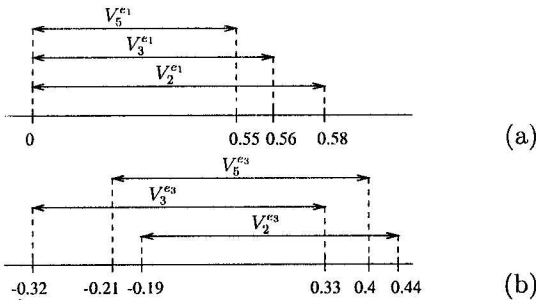


Figure 5: Two sets of predicted behavioral features (relative interval values for the amplitude). In (a) $V_2^{e_1}$, $V_3^{e_1}$ and $V_5^{e_1}$ are obtained from models m_2 , m_3 and m_5 , and their perturbations according to e_1 (see next section). In (b) $V_2^{e_3}$, $V_3^{e_3}$ and $V_5^{e_3}$ are obtained from the same models and their perturbations according to e_3 . In (a) $\Delta J(e_1)$ is 0.0176, and in (b) $\Delta J(e_3)$ is 0.1246. The models have been assumed equiprobable.

The soundness of the simulation algorithms referred to above guarantees that the models will never be falsely discriminated by a perturbation. If the measurement of the corresponding observed behavioral feature is correct, this implies that a model will never be rejected on false grounds. However, as a consequence of the incompleteness of the algorithms, competing models may fail to be discriminated while they should be.

On the basis of the selection criterion, a simple algorithm can be imagined to identify the model from M (if any) that best describes the real system by means of a minimal number of experiments. Let θ be a number between 0 and 1, determining the threshold above which we consider a model to be the best representation of the system. That is, m_i is assumed to best describe the system if $p(m_i) \geq \theta$. Let $p(m_j)$ be the a priori probabilities of the models and E is a set of pre-defined perturbation experiments.

```

set  $E_{discr}$  to {}
while  $\exists m_i \in M : p(m_i) \neq 0$  and  $\forall m_i \in M : p(m_i) \leq \theta$ 
and not  $E$  empty do
    determine  $e \in E$  for which  $\Delta J(e)$  is maximal
    perform experiment corresponding to  $e$ , determine  $Y^e$ 
    compute the a posteriori probabilities  $p(m_j|Y^e)$  of
        the models
    set  $p(m_j)$  to  $p(m_j|Y^e)$ 
    add  $e$  to  $E_{discr}$ 
    remove  $e$  from  $E$ 

```

The algorithm selects perturbation experiments till one of the following happens: a model has a sufficiently high probability, all models have zero probabilities or all possible experiments have been executed. If the algorithm terminates with $p(m_i) = 0$ for all models, obviously the assumption for completeness of M is violated.

Behavioral features are values of some properties of the system at a certain time-point, a pair of comparison, a sequence of time-points, or a sequence of pairs of comparisons. If many behavioral features are possible we can modify the algorithm above not only to select the best perturbation experiments but also to determine which set of behavioral features to be taken into account. This means that we can specify what needs to be measured and at which time-points in the experiment. The problem of selecting behavioral features is related to the problem of selecting best measurement points for the discrimination between a set of competing diagnoses in model-based diagnosis, e.g. (de Kleer & Williams 1987; Struss 1994b). If F is the set of all possible behavioral features, at each step we then determine not only the experiment $e \in E$ but also the set of features $Y \in F$ for which ΔJ is maximal. Then the tuple $\langle e, Y \rangle$ is added to E_{discr} and this tuple is not considered in subsequent iterations.

Example and evaluation

Consider the six models of a mass-spring system listed in Fig. 1 and Fig. 6 (Stoker 1992). The models differ

$$\begin{aligned}
(m_3) \quad QV(a) &= -\frac{QV(g)}{QV(L)} QV(x) - \frac{QV(c)}{QV(m)} QV(|v|) QV(v) \\
(m_4) \quad QV(a) &= -\frac{QV(g)}{QV(L)} \left(QV(x) - \frac{QV(x^3)}{QV(k)} \right) - \frac{QV(c)}{QV(m)} QV(|v|) QV(v) \\
(m_5) \quad QV(a) &= -\frac{QV(g)}{QV(L)} \left(QV(x) + \frac{QV(x^3)}{QV(k)} \right) - \frac{QV(c)}{QV(m)} QV(v) QV(a) \\
(m_6) \quad QV(a) &= -\frac{QV(g)}{QV(L)} \left(QV(x) + \frac{QV(x^3)}{QV(k)} \right) - \frac{QV(c)}{QV(m)} QV(|v|) QV(v)
\end{aligned}$$

Figure 6: Models m_3 - m_6 together with m_1 and m_2 in Fig. 1 form a set of competing models for the mass-spring system. Models m_1 and m_3 assume linear spring force. Models m_2 and m_4 assume soft spring forces (the stiffness of the spring decreases with the displacement), while the spring force in m_5 and m_6 is hard (the stiffness increases with the displacement). In m_1 , m_2 and m_5 the acceleration depends linearly on the velocity. The models m_3 , m_4 and m_6 assume quadratic dependency. The meaning of the variables and the constraints for $QV(\dot{x})$ and $QV(\dot{v})$ is the same as in Fig. 1.

	e_1	e_2	e_3	e_4	e_5		$\Delta J(e_i)$		$\Delta J(e_i)$
m_1	[0.9, 1.1]	[0.12, 0.39]	[0.8, 1.01]	[0.63, 0.97]	[1.13, 1.48]	e_1	0.0	e_1	0.0
m_2	[0.9, 1.1]	[0.11, 0.38]	[0.81, 1.01]	[0.62, 0.88]	[1.57, 2.02]	e_2	0.5315	e_2	0.3163
m_3	[0.9, 1.1]	[0.29, 0.52]	[0.88, 1.01]	[0.46, 0.67]	[1.24, 1.65]	e_3	0.0934	e_3	0.1239
m_4	[0.9, 1.1]	[0.34, 0.55]	[0.8, 1.0]	[0.41, 0.62]	[1.22, 1.64]	e_4	0.7499	e_5	0.1694
m_5	[0.9, 1.1]	[0.18, 0.43]	[0.81, 1.01]	[0.64, 0.86]	[1.34, 1.77]	e_5	0.6504		
m_6	[0.9, 1.1]	[0.23, 0.51]	[0.83, 1.02]	[0.41, 0.62]	[1.09, 1.47]				

(a)

(b)

(c)

Table 1: (a) Predictions for feature f_1 , the interval value for the amplitude, derived from the models for the perturbations e_1, \dots, e_5 . (b) The values of ΔJ computed for all perturbations, and (c) some values of ΔJ after application of e_4 (see text).

in the terms for the spring and the friction force. The experiment consists in stretching and then releasing the spring. Suppose that damped oscillations to the rest position have been observed and assume the following perturbation experiments can be performed:

- e_1 Replace the medium by an approximately frictionless medium ($c = 0$);
- e_2 Replace the medium by a more compact medium ($c = [2.85, 3.15]$);
- e_3 Test with a heavier object having mass $[11.95, 12.05]$;
- e_4 Test with a lighter object having mass $[0.7, 0.8]$;
- e_5 Release the object with initial velocity $[1.9, 2.2]$.

We consider four behavioral features:

- Y_1 is the interval value of the maximum distance from the rest position (the amplitude);
- Y_2 is the interval value of the frequency of the system;
- Y_3 is the relative interval value of the maximum amplitude for the perturbed and the original system; and
- Y_4 is the relative interval value of the frequency.

Values for Y_1 to Y_4 have been derived from the perturbed models by means of semi-quantitative simulation and comparative analysis. The predicted intervals for the amplitude are shown in Table 1(a). The first perturbation gives rise to identical predictions from all models. It is evident, even without looking at the value of $\Delta J(e_1)$, that the corresponding experiment will never

be able to distinguish between the models. The rest of the perturbations also do not give distinct intervals for this feature, but the predictions are not entirely overlapping. Hence, the measurements in the corresponding experiments may discriminate between at least some of the models.

Assume the amplitude of the system to be the only quantity measured in the experiments. Suppose the models have equal a priori probabilities $p(m_1) = \dots = p(m_6) = 1/6$ and $\theta = 0.75$, that is, a model is considered best if its probability is larger than 0.75. At the first step of the algorithm, e_4 is chosen since it maximizes ΔJ (see Table 1(b)). Assume the experiment is executed and a measurement $[0.4207, 0.5207]$ is obtained. The measurement is not consistent with the predictions derived from m_1 , m_2 , and m_5 for this perturbation and their a posteriori probabilities, therefore, become 0. The a posteriori probabilities of the other three models after the experiment are $p(m_3) = 0.2330$, $p(m_4) = 0.3835$ and $p(m_6) = 0.3835$. E_{discr} is set to $\{e_4\}$. In the next iteration, e_2 is selected (Table 1(c)). Assume the measurement $[0.3080, 0.4080]$ is obtained which gives rise to the posteriori probabilities $p(m_3) = 0.2794$, $p(m_4) = 0.3428$ and $p(m_6) = 0.3778$. e_2 is added to E_{discr} . Next, e_5 is chosen. A measurement $[1.1340, 1.2340]$ causes $p(m_3) = 0$ and the algorithm terminates, giving m_6 as the most appropriate model of the system with $p(m_6) = 0.8964$, and $E_{discr} = \{e_4, e_2, e_5\}$.

In order to evaluate the performance of the method we have adopted the following strategy. First, one of the models (m_6) was arbitrary selected. "Experimental" data was then produced by generating random intervals within the predictions of m_6 . The length of the random intervals was set equal to the size of the confidence interval of the behavioral feature ($\epsilon = 0.1$ in the case of the amplitude). Finally, the algorithm of the previous section was applied given these "measurements". This procedure was repeated 20 times and the results analyzed.

In only 15% of the cases the model that was used to generate the data was identified as the single remaining candidate. In the rest of the cases the algorithm terminated with two to three candidate models that could not be discriminated. On average, for the identification of the model 4 experiments were necessary. For comparison, when the size of the confidence interval was taken to be 0.01, in 40% of the cases m_6 was identified with average number of experiments 2.5. The results show, not surprisingly, that when the measurement error is smaller, better discrimination is achieved.

Now suppose all four features are considered, the other circumstances remaining the same. Assume further that the measurements of the amplitude and the period are independent. The joint probability distribution is then given by

$$g_i^{\{e, Y_1, Y_2, Y_3, Y_4\}}(y_1, y_2, y_3, y_4) = g_i^{\{e, Y_1, Y_3\}}(y_1, y_3) g_i^{\{e, Y_2, Y_4\}}(y_2, y_4),$$

where

$$g_i^{\{e, Y_1, Y_3\}}(y_1, y_3) = \begin{cases} g_i^{\{e, Y_1\}}(y_1) g_i^{\{e, Y_3\}}(y_3) & \text{if } y_3 = a - y_1, \\ 0 & \text{otherwise,} \end{cases}$$

with $A = [a - \epsilon_a/2, a + \epsilon_a/2]$ the amplitude of the unperturbed system, and

$$g_i^{\{e, Y_2, Y_4\}}(y_2, y_4) = \begin{cases} g_i^{\{e, Y_2\}}(y_2) g_i^{\{e, Y_4\}}(y_4) & \text{if } y_4 = t - y_2, \\ 0 & \text{otherwise,} \end{cases}$$

with $T = [t - \epsilon_t/2, a + \epsilon_t/2]$ the period of the unperturbed system. a and t have been taken 0.8 and 4.2, respectively - values that agree with the predictions of all models, and $\epsilon_a = \epsilon_t = 0.1$. In this case, e_5 maximizes ΔJ and it is selected as the best experiment (see the table below). Values of [1.134, 1.234] and [4.12, 4.22] for the amplitude and the period of the perturbed system, for instance, give rise to the posteriori probabilities $p(m_1) = \dots = p(m_5) = 0$ and $p(m_6) = 1.0$.

	e_1	e_2	e_3	e_4	e_5
$\Delta J(e_i)$	0.4049	1.1429	1.5548	1.5285	3.1811

The above evaluation procedure was again applied 20 times, now for the situation that all four features are taken into account. We found that the average number of experiments necessary to identify model m_6 was 1.1. In only two of the cases a second iteration in the algorithm was necessary. In all cases complete discrimination was achieved.

The example illustrates that when more behavioral features are considered, a higher efficiency may be achieved: measuring only the amplitude, we needed 4 experiments to discriminate between the models, while taking into account all four behavioral features a single experiment turned out to be sufficient.

Evaluation by means of random data was used to investigate the performance improvement of the algorithm for experiment selection with respect to random selection of perturbation experiments. Assume the amplitude of the system is the only quantity being measured ($\epsilon = 0.1$). After 20 times we again obtained that 4 experiments are necessary, on average, to identify the correct model. The reasons for the lack of any improvement of our method with respect to random selection are the large overlap between (some of) the predictions, the high measurement error assumed, and the low number of experiments provided. However, selecting the experiments in random order when all four features were considered, required 3.2 experiments on average to identify the correct model, whereas selecting the experiments by our method required only 1.1 experiments.

Discussion and conclusion

We have presented a method for the discrimination among competing models by selecting suitable perturbation experiments. The method chooses a maximally-discriminating experiment by means of a criterion based on the entropy measure of information. The application of this criterion was illustrated in an example concerning a set of competing models of a mass-spring system. The models in the example had the form of semi-quantitative differential equations.

Information theory has been used in model-based diagnosis (MBD) to distinguish among competing diagnoses of a faulty system (e.g. (de Kleer & Williams 1987); see (Narasimhan, Mosterman, & Biswas 1998) for other approaches). Like our method, these methods proceed by making new observations on the system. However, the work mentioned above is limited to determining the best measurement point within a given experiment, while we seek also the best experiment that would permit optimal discrimination. Struss (1994b) has extended the approach of (de Kleer & Williams 1987) by finding the best operating conditions that would give rise to the most discriminatory observations. Our work attempts to generalize this method by employing dynamical models and by extending the concept of discriminating test to discriminating perturbation experiment.

In statistics, the idea of employing the entropy measure as a discrimination criterion has been illustrated by Box & Hill (1967), Reilly (1970) and Fedorov (1972) for distinguishing between quantitative algebraic models. Kettunen, Sirvio & Varic (1988) have used the entropy to design observations discriminating among rival water quality models. Burke, Duever & Penlidis (1997;

1994) have applied the entropy to discriminate between copolymerization models. These examples are restricted to fully numerical models with precise point measurements. In this paper, we have shown how the criterion can be used when only imprecise, approximate observations and nonlinear models of the system are available.

The idea of planning perturbation experiments for model discrimination based on an entropy measure has also been proposed by (Karp, Stoughton, & Yeung 1999). They distinguish between models of a genetic regulation network by varying the expression level of involved genes or the influence of external stimuli. Their method, however, is limited to models in the form of Boolean networks without feedback and to binary perturbations. This article generalizes their approach by employing more advanced dynamical models and by extending the concept of perturbation experiments.

The work presented here can be extended into several directions. In practice, the number of possible perturbations will be infinite when the value of a quantity can be changed continuously. The problem of model discrimination as defined here should then be generalized. Instead of selecting a discrete perturbation that has been specified beforehand, a value for the quantity that maximizes (8) has to be chosen. An issue neglected thus far are the costs associated with experiments. In practice, the costs for performing an experiment may need to be balanced against its expected utility. In these cases, the problem can be reformulated as the selection of an experiment that maximizes $\Delta J(e)/h(cost(e))$, where h is a function depending on the intended application: one may be interested in effective experiments without caring about expenses, or prefer less costly tests.

Further research will concentrate on the extension of the method along the lines mentioned above, its comparison with other model discrimination techniques (e.g. (Atkinson & Fedorov 1975), (Hsiang & Reilly 1971)), and its application to real-world systems. Currently, we are applying the approach to a model discrimination problem in biology: the regulation of the cell cycle in early embryos (Obeyesekere, Tucker, & Zimmerman 1992). This system is described by second-order models and exhibits periodic behavior similar to the oscillations of the mass-spring example considered here.

References

Atkinson, A., and Fedorov, V. 1975. Optimal design: experiments discriminating between several models. *Biometrika* 62:289–303.

Berleant, D., and Kuipers, B. 1997. Qualitative and quantitative simulation: Bridging the gap. *Artificial Intelligence* 95:215–256.

Box, G., and Hill, W. 1967. Discrimination among mechanistic models. *Technometrics* 9(1):57–71.

Burke, A.; Duever, T.; and Penlidis, A. 1994. Model discrimination via designed experiments: discriminat-

ing between the terminal and penultimate models based on triad fraction data. *Macromolecular Theory Simulation* 3:1005–1031.

Burke, A.; Duever, T.; and Penlidis, A. 1997. Choosing the right model: case studies on the use of statistical model discrimination experiments. *The Canadian Journal of Chemical Engineering* 75:422–436.

de Jong, H., and van Raalte, F. 1999. Comparative environment construction: A technique for the comparative analysis of dynamical systems. *Artificial Intelligence* 115:145–214.

de Kleer, J., and Williams, B. 1987. Diagnosing multiple faults. *Artificial Intelligence* 32:97–130.

Fedorov, V. 1972. *Theory of Optimal Experiments*. Academic Press.

Hsiang, T., and Reilly, P. 1971. A practical method for discriminating among mechanistic models. *The Canadian Journal of Chemical Engineering* 49:865–871.

Karp, R.; Stoughton, R.; and Yeung, K. 1999. Algorithms for coosing differential gene expression experiments. In *RECOMB'99*, 208–217.

Kettunen, J.; Sirvio, H.; and Varis, O. 1988. Design of observations for discrimination among rival water quality models. In Dodge, Y.; Fedorov, V.; and Wynn, H., eds., *Optimal Design and Analysis of Experiments*. Elsevier Science Publishers. 257–268.

Kuipers, B. 1994. *Qualitative Reasoning: Modeling and Simulation with Incomplete Knowledge*. MIT Press.

McIlraith, S. 1994. Towards a theory of diagnosis, testing and repair. In *Working Notes DX-94*.

Narasimhan, S.; Mosterman, P.; and Biswas, G. 1998. A systematic analysis of measurement selection algorithms for fault isolation in dynamic systems. In *Working Notes DX-98*.

Obeyesekere, M.; Tucker, S.; and Zimmerman, S. 1992. Mathematical models of the cellular concentrations of cyclin and mpf. *Biochemical and Biophysical Research Communications* 184(2):782–789.

Reilly, P. 1970. Statistical methods in model discrimination. *The Canadian Journal of Chemical Engineering* 48:168–173.

Stoker, J. 1992. *Nonlinear vibrations in mechanical and electrical systems*. Wiley.

Struss, P. 1994a. Models abstraction for testing of physical systems. In *Proceedings AAAI-94*.

Struss, P. 1994b. Testing for discrimination of diagnoses. In *Working notes DX-94*.

Swaan, H. 1992. *The oxidative dehydrogenation of ethane using promoted lithium doped magnesium oxide catalysts*. Ph.D. Dissertation, University of Twente.

Vatcheva, I., and de Jong, H. 1999. Semi-quantitative comparative analysis. In *Proceedings IJCAI'99*, 1034–1040.