

Integrating Knowledge-Driven and Data-Driven Approaches to Modeling

Ljupčo Todorovski and Sašo Džeroski*

Abstract

In this paper, we present a framework for modeling dynamic systems that integrates the knowledge-based theoretical approach to modeling with the data-driven empirical modeling. The framework allows for integration of modeling knowledge specific to the domain of interest in the process of model induction from measured data. The knowledge is organized around the central notion of basic processes in the domain and it includes models thereof as well as guidelines for combining models of individual processes into a model of the entire observed system. The presented framework is applied to three tasks of modeling dynamic environmental systems from noisy measurement data in the domains of population and hydro dynamics. In all applications, the models induced with the framework can be used both to accurately predict and explain the behavior of the observed dynamic systems.

1 Introduction

Scientists and engineers build mathematical models in order to analyze and better understand the behavior of real world systems. Establishing an acceptable model of an observed system is a challenging task that occupies a major portion of the mathematical modeler's work. It involves observations and measurements of the system behavior under various conditions, selecting a set of variables that are important for modeling, and formulating the model itself. The first milestone in the process of modeling a real-world system is the choice of the modeling formalism. Ordinary differential equations (ODEs) are one of the most widely accepted formalisms for modeling of dynamic systems, i.e., systems that change their state over time [7]. In this paper, we deal with inducing models of dynamics systems based on ODEs from observed behavior of the selected system variables.

*Jožef Stefan Institute, Department of Knowledge Technologies, Jamova 39, SI-1000 Ljubljana, Slovenia

There are two main aspects of formulating a model of a real-world system. First, an appropriate model (equations) structure has to be established (the structure identification problem). Second, acceptably accurate values for the parameters are to be determined (the parameter estimation problem). Methods for system identification [13] focus mainly on solving the parameter estimation problem and make one of the following two assumptions. They either assume that the structure of the model is provided by a human expert, or assume that the structure is chosen from some general well-known class of model structures, such as linear equations, polynomials, or neural networks.

Formulating a model based on the assumption that the structure identification problem is solved by a human expert is known as the theoretical or knowledge-driven approach. Following this approach, the expert first identifies the processes that govern the behavior of the observed system. Then, using domain-specific knowledge about the identified processes, the expert writes down a proper structure of the model equations. In contrast to the theoretical approach, the empirical approach adopts a data-driven trial-and-error paradigm. The expert first chooses a general class of structures (such as linear or polynomial) that he/she believes to be adequate, fits its constant parameters, and checks how well the simulation of the model matches the observed data. If the match is not close enough, the procedure is repeated until an adequate model is found. A very limited portion (if any) of existing domain-specific knowledge is used in the modeling process.

In this paper, we aim at integrating the theoretical and empirical approaches to modeling. We present a modeling framework that integrates domain-specific knowledge in the process of data-driven induction of models based on ODEs. The knowledge is used to constrain the space of candidate models considered in the induction process. Before we can use the knowledge, we have to encode it — in the first part of the paper, we present a formalism for encoding knowledge in a form of taxonomy of basic processes and their models. The encoded knowledge includes also guidelines for composing models of the whole system from the models of the basic processes that govern its behavior. We show that this kind of knowledge, available in many textbooks on ecological modelling, can be efficiently transformed to grammars that specify the space of candidate models and can be used to guide the induction process. Our framework then uses the grammar-based equation discovery method Lagrange [17] to heuristically search through the space of candidate

models, match them against data, and find the one that fits the data best. We illustrate the generality and usefulness of the framework by applying it to three tasks of modeling environmental dynamic systems from noisy measurement data.

The paper is organized as follows. We present the formalism for encoding domain-specific modeling knowledge in Section 2. Section 3 presents the automated modeling framework, which is then applied to three tasks of inducing models of environmental dynamic systems. We present the results of these applications in Section 4. Section 5 puts our work in the context of related research. Finally, Section 6 concludes the paper with a summary and directions for further research.

2 Encoding modeling knowledge

Models of dynamic systems are often stated in terms of *basic processes* that govern the behavior of the observed system. Each basic process influences the change of one or more variables of the observed system, while the model of a basic process specifies the equations used to model its influence. Our formalism for encoding domain-specific knowledge targets knowledge about what are the basic processes in the domain of interest as well as what models of their influence on the system behavior are typically used by domain experts. The encoded knowledge also includes guidelines for combining the models of the individual basic processes into a single model of the entire observed system.

We will illustrate the use of the formalism on the example domain of population dynamics [14]. Population dynamics studies the structure and dynamics of populations, where each population is a group of individuals of the same species sharing a common environment. We consider models of the dynamic change of population concentrations (or densities) that take the form of ordinary differential equations.

2.1 Taxonomy of basic process classes

The knowledge about basic processes in the domain of interest is encoded in the taxonomy of process classes. Figure 1 presents an example taxonomy of process classes in the population dynamics domain. At the highest level of the taxonomy, we distinguish between basic processes that involve a single species or an inorganic nutrient and processes that repre-

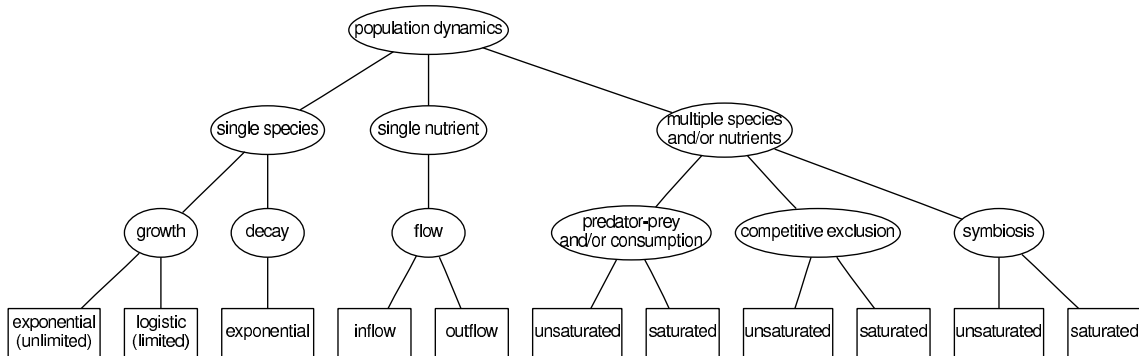


Figure 1: A taxonomy of classes of basic processes in the population dynamics domain.

sent interactions between two or more species and/or nutrients. Down the taxonomy tree, the process classes become more specific. For example, there are two kinds of processes that involve a single species: growth of a population or its decay. The leaf nodes represent particular modeling alternatives typically used by domain experts to model individual processes, such as exponential or logistic growth.

Table 1 presents the formalization of several process classes from the population dynamics. Note that we present the complete formalization of the process classes from the Figure 1 in Appendix A. The first process class *Growth* represents processes of species growth in the absence of any interaction with other species and inorganic nutrients. It has two sub-classes, specifying two models of growth. The *Exponential_growth* class specifies unlimited exponential growth of the observed population. This kind of growth is inappropriate in many real-world cases, since the environment typically has a limited carrying capacity for the given species. In such cases, the alternative *Logistic_growth* model is more appropriate.

The definition of each process class consist of three parts, which concerns variables, conditions on the variables and the equation templates used to model processes of the class.

First, we specify which types of variables are involved in processes from the given class. For example, each process in the *Growth* and *Decay* process classes involves a single variable p of type *population*. The processes in the *Feeds_on* class involve one population variable p and a set of variables cs of type concentration, which can be either a population or an

Table 1: Formalization of some of the population dynamics process classes from the taxonomy presented in Figure 1.

process class Growth(Population p)
process class Exponential_growth is Growth
expression $\text{const}(\text{growth_rate}, 0, 1, \text{Inf}) * p$
process class Logistic_growth is Growth
expression $\text{const}(\text{growth_rate}, 0, 1, \text{Inf}) * p * (1 - p / \text{const}(\text{capacity}, 0, 1, \text{Inf}))$
process class Decay(Population p)
process class Exponential_decay is Decay
expression $\text{const}(\text{decay_rate}, 0, 1, \text{Inf}) * p$
process class Feeds_on(Population p, set of Concentration cs)
condition $p \notin cs$
expression $p * \prod_{c \in cs} c$

inorganic nutrient. The declarations of variable types are inherited through the taxonomy of process classes: the *Exponential_growth* class inherits from the parent class *Growth* the fact that growth processes involve a single variable of type *population*.

The second part of the process class definition specifies constraints on the variables involved in a process of the given class. The condition $p \notin cs$ in the *Feeds_on* process class specifies that a population cannot feed (predate) on itself.

The third part of the process class definition specifies the model (i.e., equation) structure used to model the influences of processes of the given class. The model structure is based on variables involved in the process and generic constant parameters. The values of the generic constant parameters are not known and are to be fitted against measurement data. The term *const(name, lower_bound, initial, upper_bound)* is used to specify a generic constant parameter: its name, as well as lower bound, default, and upper bound values. For example, the *Exponential_growth* model involves a single nonnegative (note that a lower bound of 0 as well as infinite upper bound are specified) constant parameter that represents the growth rate with the default value of 1. The default value of the constant parameter is used as its initial value by the iterative non-linear method for parameter fitting, which matches the model structure against measured data and find the values of the model parameters that fit the data best.

Table 2: Schemes for combining the models of basic population dynamics processes into a model of the entire system.

combining scheme Population_dynamics(Inorganic i) $di/dt = - \sum_{\text{food}, f \in \text{food}} \text{const}(-, 0, 1, \text{Inf}) * \text{Feeds_on}(p, \text{food})$
combining scheme Population_dynamics(Population p) $dp/dt = \text{Growth}(p) - \text{Decay}(p) + \sum_{\text{food}} \text{const}(-, 0, 1, \text{Inf}) * \text{Feeds_on}(p, \text{food})$ $- \sum_{\text{pred}, \text{food}, p \in \text{food}} \text{const}(-, 0, 1, \text{Inf}) * \text{Feeds_on}(\text{pred}, \text{food})$

2.2 Combining scheme

The knowledge representation formalism also encodes the scheme that is used to combine the models of individual basic processes into a model of the entire system. Table 2 presents two combining schemes for building population dynamics models.

The first combining scheme specifies how to build the equation for the rate of change of an inorganic nutrient concentration i , which is represented by the time derivative di/dt of i . The change of i is negatively influenced by all the *Feeds_on* interactions in which i is being consumed by an arbitrary population p . The term *Feeds_on*(p, i) denotes the model of the *Feeds_on* process influence. The \sum aggregation function is used to sum up the influences of all such processes.

The second combining scheme specifies how to combine individual process models into a model of the rate of change of a population concentration p . The first line specifies that the time derivative of p increases with the population growth *Growth*(p) and decreases with its decay *Decay*(p). *Feeds_on* processes that involve p as a predator (or consumer) positively influence the rate of change of p , while the *Feeds_on* processes where p is involved as a prey negatively influence its rate of change. Again, influences of these processes are summed up, as shown in Table 2.

3 The modeling framework

The knowledge we encoded so far is independent of the particular modeling task at hand — it allows for modeling of an arbitrary system in the population dynamics domain. Before using the knowledge for modeling of a particular population dynamics system, we have

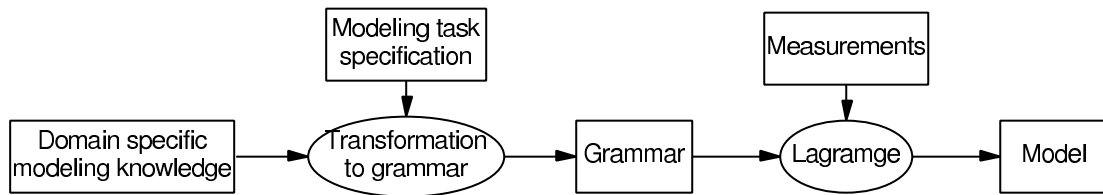


Figure 2: An automated modeling framework based on the integration of domain-specific modeling knowledge in the process of equation discovery.

Table 3: A task specification used for modeling a simple aquatic ecosystem consisting of two consumption interactions between three populations of inorganic nutrient, phytoplankton, and zooplankton.

system variable	Inorganic nut
system variable	Population phyto, zoo
processes	Growth(phyto), Decay(zoo)
processes	Feeds_on(phyto, {nut}), Feeds_on(zoo, {phyto})

to specify the modeling task. The specification includes system variables and their types along with the processes (and their process classes) that are expected to influence the system behavior. Given the modeling task specification, the encoded knowledge can be then used to generate a grammar that specifies the space of candidate model structures for the observed system. The grammar-based equation discovery method LAGRANGE is then used to search through the space of candidate models and find the one that fits the measured data best. Figure 2 gives a schematic representation of this modeling framework.

Table 3 gives an example of a modeling task specification for a simple ecosystem with three system variables *nut*, *phyto*, and *zoo* representing the concentrations of an inorganic nutrient and two populations, one of phytoplankton and one of zooplankton, respectively. The first two (growth and decay) processes specify that the populations of phytoplankton/zooplankton tend to increase/decrease in the absence of any interactions with the environment and other species. The following two *Feeds_on* processes specify that phytoplankton consumes the inorganic nutrient and that zooplankton feeds on phytoplankton.

The transformation of the modeling task specification to a grammar proceeds in a top-down manner. It starts with the starting symbol and assigns productions to it, then

proceeds with other nonterminal symbols. Nonterminal symbols in the grammar denote process classes, while the alternative productions for each nonterminal symbol specify possible expressions for modeling the corresponding process class. The starting symbol of the grammar combines the expressions for individual processes into candidate models of the whole system, according to the combining scheme.

Table 4: The grammar specifying the candidate models for modeling the simple aquatic ecosystem presented in Table 3.

Start \rightarrow $\text{time_deriv}(\text{nut}) = - \text{const}[_ : 0 : 1 :] * \text{Feeds_on_phyto_nut};$ $\text{time_deriv}(\text{phyto}) = \text{Growth_phyto} - 0 + \text{const}[_ : 0 : 1 :] * \text{Feeds_on_phyto_nut}$ $\quad - \text{const}[_ : 0 : 1 :] * \text{Feeds_on_zoo_phyto};$ $\text{time_deriv}(\text{zoo}) = 0 - \text{Decay_zoo} + \text{const}[_ : 0 : 1 :] * \text{Feeds_on_zoo_phyto}$
$\text{Feeds_on_phyto_nut} \rightarrow \text{Unsaturated_feeds_on_phyto_nut}$ $\text{Feeds_on_phyto_nut} \rightarrow \text{Saturated_feeds_on_phyto_nut}$ $\text{Unsaturated_feeds_on_phyto_nut} \rightarrow \text{phyto} * \text{nut}$ $\text{Saturated_phyto_nut} \rightarrow \text{phyto} * \text{nut} / (\text{nut} + \text{const}[\text{saturation_rate} : 0 : 1 :])$
$\text{Growth_phyto} \rightarrow \text{Exponential_growth_phyto}$ $\text{Exponential_growth_phyto} \rightarrow \text{const}[\text{growth_rate} : 0 : 1 :] * \text{phyto}$ $\text{Logistic_growth_phyto} \rightarrow \text{const}[\text{growth_rate} : 0 : 1 :] * \text{phyto} * (1 - \text{const}[\text{capacity} : 0 : 1 :])$
$\text{Feeds_on_zoo_phyto} \rightarrow \text{Unsaturated_feeds_on_zoo_phyto}$ $\text{Feeds_on_zoo_phyto} \rightarrow \text{Saturated_feeds_on_zoo_phyto}$ $\text{Unsaturated_feeds_on_zoo_phyto} \rightarrow \text{zoo} * \text{phyto}$ $\text{Saturated_zoo_phyto} \rightarrow \text{zoo} * \text{phyto} / (\text{phyto} + \text{const}[\text{saturation_rate} : 0 : 1 :])$
$\text{Decay_zoo} \rightarrow \text{Exponential_decay_zoo}$ $\text{Exponential_decay_zoo} \rightarrow \text{const}[\text{decay_rate} : 0 : 1 :] * \text{zoo}$

For example, consider the aquatic ecosystem example from Table 3. The start symbol uses the combining schemes from Table 2 to compose a model of the whole ecosystem (see the first production in the grammar from Table 4). The right hand side of the production for the starting symbol builds a system of three differential equations, one for each variable of the observed ecosystem. The second equation (for the time derivative of *phyto*) is built by summing up the growth process for the *phyto* population, i.e., $\text{Growth}(\text{phyto})$, the *phyto* decay process (0, since none of them are specified), the consumption process where *phyto*

is involved as consumer, i.e., $Feeds_on(phyto, \{nut\})$, and the predator-prey process where *phyto* has the role of prey, i.e., $Feeds_on(zoo, \{phyto\})$. The productions for the other nonterminal symbols are built according to the specification of the taxonomy of process classes from Table 1.

The grammar from Table 4 generates a set of candidate model structures for modeling the simple aquatic ecosystem. LAGRAMGE can heuristically search through this space and find the model that fits the data best. The LAGRAMGE algorithm employs an iterative non-linear optimization method [3] to fit the constant parameters of the candidate model structures. Heuristic function based on the discrepancy between simulated and observed values of the system variables, guides the search. Further details about LAGRAMGE can be found in [17, 16].

4 Modeling environmental dynamic systems

We applied our framework to three tasks of modeling dynamic systems from real-world data. Two tasks from the population dynamics domain have been already addressed by other equation discovery methods. The third task is from the domain of hydrodynamics. The evaluation criteria used in these experiments are predictive accuracy (as measured by root mean squared error, RMSE) and comprehensibility of the discovered models as evaluated by domain experts.

4.1 Modeling algal growth in the Lagoon of Venice

The Lagoon of Venice measures 550 km², but is very shallow, with an average depth of less than 1 m. It is heavily influenced by anthropogenic inflow of nutrients — 7 [mio kg/year] of nitrogen and 1.4 [mio kg/year] of phosphorus [1]. These (mainly nitrogen) loads are above the Lagoon’s admissible trophic limit and generate its dystrophic behavior, which is characterized by excessive growth of algae, mainly *Ulva rigida*. Four sets of measured data were available [4] for modeling the growth of algae in the Lagoon. The data were sampled weekly for slightly more than one year at four different locations in the Lagoon. Location 0 was sampled in 1985/86, locations 1, 2, and 3 in 1990/91. The sampled quantities are nitrogen in ammonia NH_3 , nitrogen in nitrate NO_3 , phosphorus in orthophosphate PO_4

Table 5: A task specification used for modeling the growth of algae biomass in the Lagoon of Venice.

variable	Inorganic temp, DO, NH3, NO3, PO4
system variable	Population biomass
process	Growth(biomass) biomass_growth
process	Decay(biomass) biomass_decay
process	Feeds_on(biomass, *) biomass_grazing

(all in $[\mu\text{g}/\text{l}]$), dissolved oxygen DO (in percentage of saturation), temperature T ([degrees C]), and algal biomass B (dry weight in $[\text{g}/\text{m}^2]$).

Previous experiments with automated modeling of algal growth in the Lagoon of Venice with equation discovery [9] used the GOLDHORN method [10]. Since GOLDHORN could not find an accurate model based on the set of measured variables, two additional variables were calculated and added to the set of variables. These are the growth and mortality rates, which are known quantities in ecological modeling and were calculated according to the simplified version of an existing model of algal growth in the lagoon proposed by [4]. From the extended set of variables and data measured at Location 0, GOLDHORN discovered a difference equation for predicting biomass that, due to the large measurement errors (estimated at the level of 20-50%), does not fit the data perfectly, but it still predicts most of the peaks and crashes of the biomass concentration correctly [9]. Although the equation model involves the mortality rate, as calculated by domain experts, the model itself is still a black-box model that does not reveal the limiting factors for the biomass growth in the lagoon.

The task specification of modeling algal growth in the Lagoon of Venice, given in Table 5, lists the types of the observed variables and the processes that are important for the growth of algae (and their biomass) in the lagoon. Note that we are only interested in a model of the biomass growth, so biomass is the only system variable. The specification of the `biomass_grazing` process leaves the nutrient parameter of the `Feeds_on` process class unspecified (denoted using the symbol `*`). Since ecologists did not know the limiting factors for biomass growth, they let LAGRANGE search for the model that would reveal them.

The experiments with LAGRAMGE were performed using a grammar automatically built from the task specification in Table 5 and the library of modeling knowledge outlined in Section 2 and completely specified in Appendix A. The grammar generates 6248 candidate model structures, among which the following model matches the measured data on the Location 0 best:

$$\begin{aligned} d\text{biomass}/dt = & 6.17 \cdot 10^{-5} \cdot \text{biomass} \cdot \left(1 - \frac{\text{biomass}}{1.80}\right) \\ & + 3.01 \cdot 10^{-4} \cdot \text{biomass} \cdot \text{DO} \cdot \frac{\text{NO3}}{\text{NO3} + 6.28} - 0.0319 \cdot \text{biomass}. \end{aligned}$$

The model for Location 0 tells us that the limiting factors for the biomass growth of the *Ulva rigida* in the lagoon are dissolved oxygen (DO) and nitrogen in nitrate (NO3).

Furthermore, LAGRAMGE discovered another model from the Location 2 data:

$$\begin{aligned} d\text{biomass}/dt = & 4.79 \cdot 10^{-5} \cdot \text{biomass} \cdot \left(1 - \frac{\text{biomass}}{0.844}\right) \\ & + 0.406 \cdot \text{biomass} \cdot (1 - e^{-0.216 \cdot \text{temp}}) \cdot (1 - e^{-0.413 \cdot \text{DO}}) \cdot \frac{\text{NH3}}{\text{NH3} + 10} \\ & - 0.0343 \cdot \text{biomass}. \end{aligned}$$

This model tells us that the limiting factors for the biomass growth are temperature (`temp`), dissolved oxygen (DO), and nitrogen in ammonia (NH3). Although the two models are not identical, they both identify dissolved oxygen and nitrogen based nutrients to be limiting factors for biomass growth. The differences between the two models may be due to the fact that the measurements were taken during two different time periods.

In the experiments with the data measured at the other two locations (1 and 3), LAGRAMGE did not find an accurate model of the biomass growth. Note that these results still compare favorably with the results obtained by GOLDHORN, which discovered an acceptable model for Location 0 only.

Figure 3 compares the measured and simulated values of the biomass change over time for both models. We ran long-term simulations of the models from the initial value of the biomass without restarting the simulation process at each measurement point. For values of all other variables needed during the simulation, we used the measurement at the nearest time point in the past. As in the GOLDHORN experiments, due to the high measurement errors (of the order 20-50%), the models discovered by LAGRAMGE do not

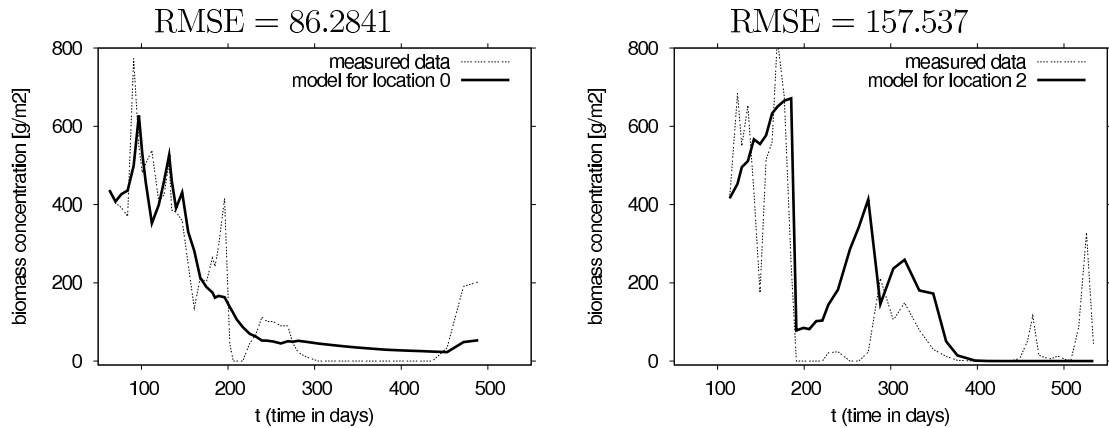


Figure 3: Simulations of the two models of the biomass growth in the Lagoon of Venice, discovered by LAGRAMGE, compared to the measured biomass concentration (left-hand side: Location 0, right-hand side: Location 2).

fit the measured data perfectly. However, they correctly predict most of the peaks and crashes of the biomass. These events are more important to ecologists than the degree of fit. Note an important advantage of these models over the one discovered by GOLDHORN. While the GOLDHORN model is black-box, the models discovered by LAGRAMGE identify the most important limiting factors for biomass growth in the Lagoon of Venice.

4.2 Modeling phytoplankton growth in Lake Glumsø

Lake Glumsø [8] is situated in a sub-glacial valley in Denmark. It is shallow with average depth of about 2 [m] and its surface area is 266,000 [m²]. For several years, prior to the measurements discussed here, it was receiving mechanically-biologically treated waste water from a community with 3,000 inhabitants and a surrounding area which was mainly agricultural with almost no industry. The high nitrogen and phosphorus concentration in the treated waste water had caused hypereutrophication. The lake contained no submerged vegetation, probably due to the low transparency of the water and the oxygen deficit at the bottom.

Domain experts considered the concentrations of phytoplankton (`phyto`), zooplankton (`zoo`), soluble nitrogen (`nitro`), soluble phosphorus (`phosp`), and the water temperature (`temp`) relevant for modeling the phytoplankton growth. These variables were measured at 14 distinct time points over a period of two months. This amount of measured data itself

Table 6: A specification of the modeling the phytoplankton growth in Lake Glumsø task.

variable	Inorganic temp, nitro, phosp
variable	Population zoo
system variable	Population phyto
process	Decay(phyto) phyto_decay
process	Feeds_on(phyto, *) phyto_grazing
process	Feeds_on(zoo, phyto) zoo_grazing

was insufficient for modeling, so additional processing was applied to obtain a suitable data set [18]. First, dotted graphs of the measurements were plotted and given to three human experts to draw a curve that, in their own opinion, described the dynamic behavior of the observed variable between the measured points. A properly plotted expert curve can be regarded as an additional source of reliable data. Curves drawn by the human experts were then smoothed with Bezier splines. Finally, three data sets were obtained by sampling the splines derived from each of the three human experts' approximations at regular time intervals with time step $h = 0.03215$ day. The data set provided by the first expert was used for the experiments.

The task of modeling phytoplankton growth in Lake Glumsø, as specified in Table 6, lists the types of the observed variables and the processes that are important for the growth. Note that the specification of the `phyto_grazing` process leaves the nutrient parameter of the `Feeds_on` class unspecified (denoted using the symbol `*`). This is because the experts did not know the limiting factors for the biomass growth, and therefore we let LAGRAMGE search for the model that would identify them.

The experiments with LAGRAMGE were performed using a grammar automatically built from the task specification and the library of modeling knowledge outlined in Section 2. The grammar generates 496 candidate model structures. Among these, the model with the minimal value of the LAGRAMGE heuristic function on the measured data was:

$$\dot{\text{phyto}} = 0.553 \cdot \text{temp} \cdot \frac{\text{phosp}}{0.0264 + \text{phosp}} - 4.35 \cdot \text{phyto} - 8.67 \cdot \text{phyto} \cdot \text{zoo}.$$

The structure of the discovered equation tells us that phosphorus is a limiting factor for phytoplankton growth in the lake and that the growth is temperature dependent.

Note that the same model was already discovered by LAGRAMGE [18] using a hand-crafted grammar based on the human expert knowledge about modeling population dynamics. However, there is an important difference between the experiment performed here and the one performed with the previous version of LAGRAMGE. In the previous study, we were not able to specify bounds on the values of the constant parameters, so the output of LAGRAMGE was manually post-processed in order to filter out the equations with invalid values of the constant parameters (e.g., negative growth or saturation rate). In the experiment presented here, there is no need for this additional step, since the knowledge about the valid values of the constant parameters was encoded within the library of domain-specific knowledge.

4.3 Modeling the water level variation in Ringkøbing fjord

In the last series of experiments, we illustrate that the proposed formalism allows for partial model specification. In such a case, a human expert specifies only some parts of the model structure and leaves others unspecified or partly specified. Our framework can be then used to determine both the structure and the parameters of the unspecified parts.

An example of such a task is modeling water level variation in Ringkøbing fjord, a shallow estuary located at the Danish west coast, where it experiences mainly easterly and westerly winds.¹ Wind forcing causes large short term variation of the water level (h) measured at the gate between the estuary and the North Sea. Domain experts specified the following partial model for the temporal variation of the water level in the estuary:

$$\dot{h} = \frac{f(a)}{A}(h_{sea} - h + h_0) + \frac{Q_f}{A} + g(W_{vel}, W_{dir}).$$

The water level response to the wind forcing, dependent on both wind speed (variable W_{vel} , measured in [m/s]) and direction (W_{dir} , measured in degrees), is modeled by an unknown function g . Apart from wind forcing, the water level is dominated by the fresh water supply (Q_f , measured in [m³/s]). When the gate is closed, fresh water is accumulated in the estuary causing a water level rise of Q_f/A , where A is the surface area of the estuary

¹The task was used as an exercise within a post-graduate course on modeling dynamic systems organized in 2000. Since the Web page of the course is no longer available, we cannot provide a proper reference to the original task specification. Note also that we could not consult domain experts and therefore could not obtain expert comments on the induced models.

Table 7: Formalization of the partially specified model of the water level variation in the Ringkøbing fjord and specification of the modeling task.

function class Salt_water_drive(Opening a, Level h_sea, Level h, Surface A)
expression (h_sea - h + const[h_0:-5:0.1:5])
function class Fresh_water_flow(Flow Q_f, Surface A)
expression Q_f / A
combining scheme Water_level_change(Level h)
time_deriv(h) = (F(a) / A) * Salt_water_drive(a, h_sea, h, A) + Fresh_water_flow(Q_f, A) + G(W_vel, W_dir)

variable Opening a
variable Level h_sea
system variable Level h
variable Surface A
variable Flow Q_f
variable Velocity W_vel
variable Direction W_dir
Function F(a) flow_friction
Function G(W_vel, W_dir) wind_forcing

measured in squared meters. During periods when the gate is open, the stored fresh water is emptied in the North Sea. The gate is also opened in order to maintain sufficient water level in the estuary, in which case the water rise is driven by the difference between the water level in the open sea (variable h_{sea} , measured in meters), the water level in the estuary (h , measured in meters), and the constant parameter (h_0). The flow is restricted by the friction of the flow, modeled by an unknown function f of the number of gate parts being open (a). Namely, the gate consists of 14 parts and allows for opening some parts and closing others. The value of A is not directly observed, but a function that calculates A on the basis of h is provided, so A can be also treated as observed variable.

In order to apply our framework to the task of model completion, we first encode the partial specification within our formalism. The formalization of the partial model specification from Table 7 follows the partial model formula proposed by the domain experts. The

Table 8: Formalization of the modeling alternatives for the unspecified parts of the model of the water level variation in the Ringkøbing fjord.

function class F(Opening a)
function class F_0 is F
expression const[$[-5000:0.1:5000]$]
function class F_1 is F
expression Polynomial($\{a\}$, const[$[-5:0.1:5]$], 5)
function class G(Velocity W_vel, Direction W_dir)
function class G_0 is G
expression const[$[-5000:0.1:5000]$]
function class G_1 is G
expression Polynomial($\{W_vel, W_dir\}$, const[$[-5:0.1:5]$], 5)
function class G_2 is G
expression Polynomial($\{W_vel, \sin(W_dir), \cos(W_dir)\}$, const[$[-5:0.1:5]$], 5)

formula is decomposed into two building blocks following the explanation of the partial model specification².

In the second step, we formalize the modeling alternatives for each of the unspecified parts of the model, i.e., the f and g functions. In the experiments, we use simple constant and polynomial models due to the lack of additional domain knowledge. The modeling alternatives used in the experiments are presented in Table 8. The first modeling alternatives **F_0** and **G_0** for f and g are the simplest possible models, i.e., constants within the interval $[-5000, 5000]$ with the initial values of 0.1. The next two alternatives (**F_1** and **G_1**) are polynomials of the appropriate system variables with constant parameters within the interval $[-5, 5]$ (and initial values of 0.1). The maximal degree of the polynomials is five. Finally, we used one additional modeling alternative for the g function (**G_2**) that replaces the wind direction value (that represents angle) with the sine and cosine transformation thereof in the polynomial.

The data about the observed variables is collected by hourly measurements of all the observed variables within the period from 1st of January to 10th of December 1999. We

²In order to specify different alternatives for modeling flow friction F and wind forcing G , we introduce a taxonomy of function classes. The definitions of function classes are the same as definitions of process classes. We distinguish between functions and processes because the former do not represent processes from the domain; in this case, we use them to specify alternative models for the individual building blocks.

Table 9: The root mean squared errors (RMSE, estimated on both training data and using 10-fold cross-validation) of the four water level variation models induced by LAGRAMGE with (three first rows) and without (last row) using the partial model specification provided by the domain experts. The last column gives number of candidate model structures (#CMS) considered during the search.

task specification	training RMSE	cross-validated RMSE	#CMS
F_0 + G_0	0.0848	0.106	1
F_1 + G_1	0.0655	0.0931	378
F_1 + G_2	0.0585	0.0903	2184
polynomial	0.0556	2.389	2801

used the task specification presented in Tables 7 and 8 to induce a model from the measurements with LAGRAMGE. We performed three experiments. In the first, we used the F_0 and G_0 modeling alternatives, in the second we used the F_1 and G_1 modeling alternatives, and in the third we used the F_1 and G_2 modeling alternatives. In order to evaluate the benefit of using partial model specification, we performed one additional experiment in which no knowledge about model structure was used. In this last experiment, we used a polynomial model of the water level change in the fjord. We used 10-fold cross-validation for estimating the RMSE of the induced models.

Table 9 summarizes the results of the experiments. The best cross-validated performance is gained using the partial model specification provided by the experts in combination with the F_1 and G_2 modeling alternatives for the unspecified parts of the structure. LAGRAMGE proposed the following models for f and g :

$$\begin{aligned}
 f(a) &= 5 + 5 \cdot a + 5 \cdot a^2 + 5 \cdot a^3 - 1.01 \cdot a^4 \\
 g(W_{\text{vel}}, W_{\text{dir}}) &= -0.00137 - 0.0106 \cdot \cos W_{\text{dir}} + 0.218 \cdot \cos W_{\text{dir}} \cdot \sin W_{\text{dir}} \\
 &\quad + 0.0106 \cdot W_{\text{vel}} \cdot \cos W_{\text{dir}} \cdot \sin W_{\text{dir}} - 0.0128 \cdot W_{\text{vel}}^2 \cdot \cos W_{\text{dir}} \cdot \sin W_{\text{dir}} \\
 &\quad - 0.000428 \cdot W_{\text{vel}}^3 \cdot \cos W_{\text{dir}} \cdot \sin W_{\text{dir}}.
 \end{aligned}$$

The graph on the left-hand side of Figure 4 shows the simulation of this model compared to the measured water level in the Ringkøbing fjord. We ran a long-term simulation of the model from the initial value of the water level without restarting the simulation process

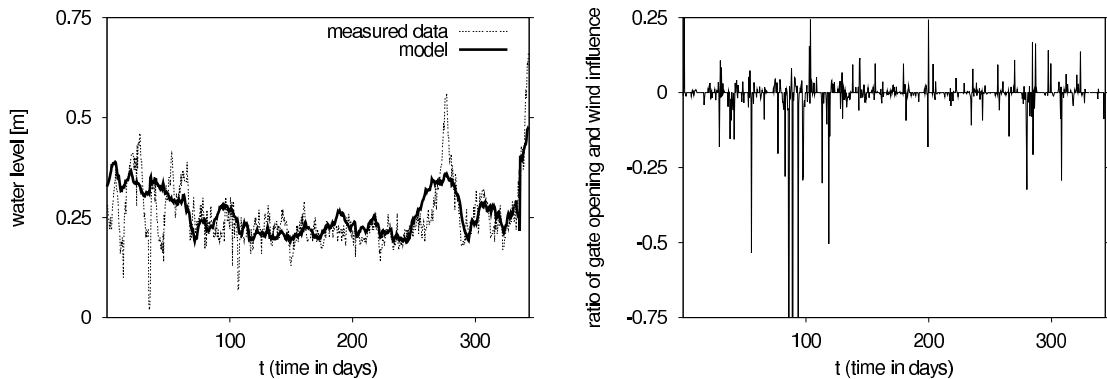


Figure 4: Simulation of the water level variation model induced by LAGRAMGE compared to the measured water level (left-hand side) and the ratio of the gate opening and the wind influences on the water level change in the Ringkøbing fjord as modeled by LAGRAMGE.

Table 10: The RMSE and correlation coefficient (r) for the short-term (one hour and one day) prediction of the water level in the Ringkøbing fjord compared to the RMSE and r of the simulation over the whole observation period.

prediction/simulation period	RMSE	r
one hour	0.0168	0.976
one day	0.0425	0.845
whole observation period	0.0585	0.659

at any measurement point. For the values of all other system variables needed during the simulation, we used the measurement at the nearest time point in the past.

Note that the model follows the general pattern of water level variation. The long-term simulation of the model, however, fails to precisely capture the short-term (hour) changes of the water level in the fjord. To test the short-term prediction power of the model, we performed two additional simulations, which we restarted with the true measured water level values at every hour and at every day (24 hours). Table 10 presents the results of this analysis. They show that the model is suitable for short-term prediction of the water level in the Ringkøbing fjord.

Since the model induced by LAGRAMGE follows the partial structure specification provided by the human experts, further analysis of the model can be performed. For example, we can compare the influence of the gate opening (modeled by $f(a)(h_{\text{sea}} - h + h_0)/A$) with

the effect of the wind (modeled by $g(W_{\text{vel}}, W_{\text{dir}})$). The graph on the right-hand side of Figure 4 shows the ratio of the gate opening and the wind influences on the water level change in the Ringkøbing fjord. The low magnitude of the ratio shows that the influence of the wind prevails over the influence of the gate opening most of the time. The only exceptions occur in the period from 80 to 100 days from the beginning of the measurement, that is, the end of March and beginning of April 1999.³

The polynomial model of the water level variation that ignores the partial specification of the model performs best on the training data. However, the model’s small RMSE is due to the overfitting of the training data, since the cross-validated RMSE of this model (2.389) is much larger than the cross-validated RMSE of the models that follow the partial structure specification.

In sum, the Ringkøbing fjord experiments show the capability of our framework to address modeling tasks in which human experts can partially specify the model structure and leave some of its parts unspecified.

5 Related work

The presented modeling framework follows the paradigm of compositional modeling [6], an automated approach to building qualitative models from observations in presence of domain-specific modeling knowledge. In the compositional modeling, knowledge is organized as a library of model fragments. Given a modeling task specification (or scenario), compositional modeling methods compose a set of appropriate fragments into a model that is suitable for modeling the observed system. The obtained model is evaluated by qualitative simulation [11]. The compositional modeling approach is mainly applied to the tasks of building qualitative models. For example, authors of [5] apply this approach to the task of inducing qualitative models of chemical reaction pathways from noisy measurement data.

Although the concepts introduced within the QR area are also relevant for automated building of quantitative models of real-world systems, this idea has not been widely explored. A notable exception is the PRET reasoning system for automated modeling of

³Note again that, unfortunately, we could not obtain expert comments on these results.

dynamic systems [2], which employs two kinds of knowledge. The first is domain-specific knowledge in the form of “conservation rules”, such as Kirchoff’s law in the domain of electrical circuits, which specifies that the sum of input and output currents at any observed point in the circuit is zero. Similarly, the force balance rule in the mechanics specifies that the sum of forces at any observed coordinate of the mechanical system is zero. These rules are more general than domain knowledge about model fragments used in compositional modeling approaches, and constrain the space of possible models much less. This kind of knowledge more appropriate for modeling in engineering domains such as electrical circuits and mechanics. Knowledge about basic processes in the domain of interest (that is more appropriate for building explanatory models in the biological or chemical domains) can not be fitted within PRET’s framework.

The approach presented here builds up on previous work on inducing process-based models presented in [12]. Note however, that the taxonomy of process classes introduced here allows for better organization of knowledge. Also, here we made the combining scheme explicit, as opposed to the implicit combining scheme used there — models of individual processes are always additively combined. Finally, note that methods for inducing equation-based models from data have been already used to induce models of environmental systems, see, e.g., [9, 18]. They have been also successfully applied to tasks of revising existing models of environmental system, see, e.g., [21, 15, 19].

6 Summary and further work

In this paper, we presented a framework for automated modeling of dynamic systems based on equation discovery. The framework integrates the theoretical knowledge-driven and the empirical data-driven approaches to modeling. The framework provides a formalism for encoding and integrating domain-specific knowledge in the process of model induction. The knowledge is organized in a taxonomy of process classes, each representing an important class of processes in the observed domain. This high-level knowledge representation can be automatically transformed to the operational form of grammars that specify the space of candidate models of the observed system. The equation discovery method LAGRAMGE can be then used to search through the space of candidate models and find the one that fits the measured data best.

The results of the application of the modeling framework to two real-world modeling tasks show that our framework is capable of inducing comprehensible dynamic systems' models from real-world measurement data. Our framework performs better than existing equation discovery methods on the tasks of modeling algae growth in Lagoon of Venice in terms of performance, flexibility, and comprehensibility of the discovered models. The experiment on modeling water level variation in Ringkøbing fjord illustrates the capability of our framework to address modeling tasks, in which a human expert partially specifies the model structure and leaves other parts unspecified — our framework can be then applied to induce the unspecified parts of the model from data.

Although the scope of this paper is modeling dynamic change of the observed system through time, the approach can be extended to incorporate knowledge about spatial processes also. In the population dynamics domain, these processes would reflect the spatial diffusion of populations through the environment. The extended approach would allow modeling of spatio-temporal dynamic systems with partial differential equations [20].

References

- [1] G. Bendoricchio, G. Coffaro, and C. DeMarchi. A trophic model for *Ulva Rigida* in the Lagoon of Venice. *Ecological Modelling*, 75/76:485–496, 1994.
- [2] E. Bradley, M. Easley, and R. Stolle. Reasoning about nonlinear system identification. *Artificial Intelligence*, 133:139–188, 2001.
- [3] D. S. Bunch, D. M. Gay, and R. E. Welsch. Algorithm 717: Subroutines for maximum likelihood and quasi-likelihood estimation of parameters in nonlinear regression models. *ACM Transactions on Mathematical Software*, 19:109–130, 1993.
- [4] G. Coffaro, G. Carrer, and G. Bendoricchio. Model for *Ulva Rigida* growth in the Lagoon of Venice. Technical report, University of Padova, Padova, Italy, 1993. UNESCO MURST Project: Venice Lagoon Ecosystem.
- [5] G. M. Coghill, S. M. Garrett, and R. D. King. Learning qualitative models in the presence of noise. In *Proceedings of the Sixteenth International Workshop on Qualitative Reasoning*, 2002.

- [6] B. Falkenheiner and K. D. Forbus. Compositional modeling: Finding the right model for the job. *Artificial Intelligence*, 51:95–143, 1991.
- [7] N. Gershenfeld. *The Nature of Mathematical Modeling*. Cambridge University Press, Cambridge, UK, 1999.
- [8] S. E. Jørgensen, L. Kamp-Nielsen, T. Chirstenen, J. Windolf-Nielsen, and B. Westergaard. Validation of a prognosis based upon a eutrophication model. *Ecological Modelling*, 32:165–182, 1986.
- [9] B. Kompare and S. Džeroski. Getting more out of data: automated modelling of algal growth with machine learning. In *Proceedings of the International symposium on coastal ocean space utilisation*, pages 209–220, Yokohama, Japan, 1995.
- [10] V. Križman. *Avtomatsko odkrivanje strukture modelov dinamičnih sistemov*. PhD thesis, Faculty of computer and information science, University of Ljubljana, Ljubljana, Slovenia, 1998. In Slovene.
- [11] B. Kuipers. *Qualitative reasoning: modeling and simulation with incomplete knowledge*. MIT Press, Cambridge, MA, 1994.
- [12] P. Langley, J. Sanchez, L. Todorovski, and S. Džeroski. Inducing process models from continuous data. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 347–354, San Mateo, CA, 2002. Morgan Kaufmann.
- [13] L. Ljung. Modelling of industrial systems. In *Proceedings of Seventh International Symposium on Methodologies for Intelligent Systems*, pages 338–349, Berlin, 1993. Springer.
- [14] J. D. Murray. *Mathematical Biology*. Springer, Berlin, 1993. Second, Corrected Edition.
- [15] K. Saito, P. Langley, T. Grenager, C. Potter, A. Torregrosa, and S. A. Klooster. The computational revision of quantitative scientific models. In *Proceedings of the Fourth International Conference on Discovery Science*, pages 336–349, Berlin, 2001. Springer.

- [16] L. Todorovski. *Using domain knowledge for automated modeling of dynamic systems with equation discovery*. PhD thesis, Faculty of computer and information science, University of Ljubljana, Ljubljana, Slovenia, 2003. Available for download at <http://www-ai.ijs.si/~ljupco/ed/>.
- [17] L. Todorovski and S. Džeroski. Declarative bias in equation discovery. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 376–384, San Mateo, CA, 1997. Morgan Kaufmann.
- [18] L. Todorovski, S. Džeroski, and B. Kompare. Modelling and prediction of phytoplankton growth with equation discovery. *Ecological Modelling*, 113:71–81, 1998.
- [19] L. Todorovski, S. Džeroski, P. Langley, and C. S. Potter. Using equation discovery to revise an earth ecosystem model of the carbon netproduction. *Ecological Modelling*, 170:141–154, 2003.
- [20] L. Todorovski, S. Džeroski, A. Srinivasan, J. Whiteley, and D. Gavaghan. Discovering the structure of partial differential equations from example behavior. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pages 991–998, San Mateo, CA, 2000. Morgan Kaufmann.
- [21] P. A. Whigham and F. Recknagel. Predicting chlorophyll-a in freshwater lakes by hybridising process-based models and genetic algorithms. *Ecological Modelling*, 146(1-3):243–252, 2001.

A Complete library of modeling knowledge for population dynamics

```

type Concentration is real
type Concentrations is set(Concentration)

type Population is Concentration
type Populations is set(Population)

type Inorganic is Concentration

```

```

function class Saturation(Concentration c)

function class No_saturation() is Saturation
  expression c

function class Saturation_type_1() is Saturation
  expression c / (c + const(saturation_rate,0,1,Inf))

function class Saturation_type_2() is Saturation
  expression c * c / (c * c + const(saturation_rate,0,1,Inf))

function class Saturation_type_3() is Saturation
  expression 1 - exp(-const(saturation_rate,0,1,Inf) * c)

process class Growth(Population p)

process class Exponential_growth() is Growth
  expression const(growth_rate,0,1,Inf) * p

process class Logistic_growth() is Growth
  expression const(growth_rate,0,1,Inf) * p * (1 - p / const(capac,0,1,Inf))
process class Decay(Population p)

process class Exponential_decay() is Decay
  expression const(decay_rate,0,1,Inf) * p

process class Flow(Concentration c)

process class Constant_inflow() is Flow
  expression const(inflow_rate,0,1,Inf)

process class Constant_outflow() is Flow
  expression -const(outflow_rate,0,1,Inf)

process class Feeds_on(Population p, Concentrations cs)
  condition p not in cs
  expression p * product({c}, c in cs, Saturation(c))

process class Interaction(Populations ps)
  condition cardinality(ps) at least 2

process class Competitive_exclusion() is Interaction
  expression -product({p}, p in ps, Saturation(p))

process class Symbiosis() is Interaction
  expression product({p}, p in ps, Saturation(p))

```



```
combining scheme Population_dynamics(Inorganic i)
time_deriv(i) = + sum({}, true, Flow(i))
               - sum({p, food}, i in food, const(_,0,1,Inf) * Feeds_on(p, food))
```

```
combining scheme Population_dynamics(Population p)
time_deriv(p) = + sum({}, true, Growth(p))
               + sum({}, true, Flow(p))
               + sum({food}, true, const(_,0,1,Inf) * Feeds_on(p, food))
               - sum({}, true, Decay(p))
               - sum({p1, food}, p in food, const(_,0,1,Inf) * Feeds_on(p1, food))
               + sum({ps}, p in ps, const(_,0,1,Inf) * Interaction(ps))
```