

knowledge-based framework for modeling dynamic environmental systems

sašo džeroski (joint work with ljupčo todorovski)
jožef stefan institute, ljubljana, slovenia

4th workshop on environmental applications of machine learning
4th european conference on ecological modelling

bled, september 2004

motivation: automated modeling

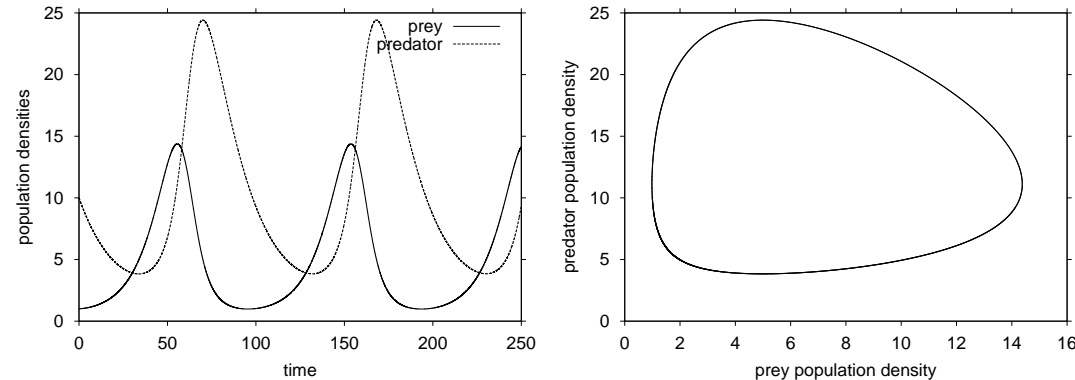
- computational support for building equation-based models
- existing modeling (system identification) methods:
 - assume that the model structure is known
 - or assume linear/NN structure (if the model is unknown)
 - task: determine the values of the model parameters
- in contrast, our (equation discovery based) method:
 - does not assume (a single) prescribed model structure
 - task: determine both structure and parameters of the model

motivation: the use of domain-specific knowledge

- two approaches to modeling of real-world systems:
 1. theoretical (knowledge-driven)
 - domain expert derives a proper model structure based on domain-specific modeling knowledge
 2. empirical (data-driven)
 - try different models to fit observed data (trial-and-error)
- in 1., a lot of domain knowledge and few data are needed;
 - terms in equations typically correspond to processes (equations are meaningful to domain scientists)
- in 2., many data of good quality are needed;
 - no domain knowledge can be used
 - terms in equations are rarely meaningful to the domain scientist
- to integrate 1. and 2., that is the question!

discovering dynamics: problem definition and example

- GIVEN an example behavior of a dynamic system:



- FIND the system dynamics equations:

$$\dot{N} = a \cdot N - b \cdot NP$$

$$\dot{P} = c \cdot NP - d \cdot P$$

- N is prey population density (hares)
- P is predator population density (lynx)

discovering dynamics: a declarative bias approach

- space of possible equations (language bias):
 - user defined (declarative)
 - based on the domain-specific knowledge
- declarative bias formalism in lagrange: context-free grammar
 - prescribes the form of the expression on the right-hand side
 - generates legal expressions in the C programming language
 - result of the derived expressions: (double)
- discovered equations are of the form $\dot{v}_d = E$,
where E is an expression derived using the given grammar

common language biases for equations

- universal grammar (arithmetical expressions):

$$E \rightarrow E + F \mid E - F \mid F$$

$$F \rightarrow F * T \mid F / T \mid T$$

$$T \rightarrow \textit{const} \mid v \mid (E)$$

- multivariate polynomials:

$$E \rightarrow \textit{const} \mid \textit{const} \cdot F \mid E + \textit{const} \cdot F$$

$$F \rightarrow v \mid v \cdot F$$

monod example (subexpressions): context free grammar

- context free grammar used for environmental dynamic systems:

$$E \rightarrow \text{const} \mid \text{const} \cdot F \mid E + \text{const} \cdot F$$

$$F \rightarrow v \mid Y \mid v \cdot Y$$

$$Y \rightarrow \text{monod}(\text{const}, v)$$

$$Y : \frac{v}{v + \text{const}}$$

- user defined function monod:

```
double monod(double c, double v) {  
    return(v / (v + c));  
}
```

lagramge: search in the space of structures

- the grammar defines the space of possible equation structures: the structures that can be derived with at most N applications of rules from the grammar
- the grammar also
 - orders the search space
 - starting from the simplest equation structure ...
 - allows the use of different search strategies
- two search strategies implemented in lagramge:
 - beam search
 - exhaustive search

lagrange: constant parameters fitting

- nonlinear optimization method:
 - downhill simplex
 - minimize the difference between given and simulated data
 - use integration instead of differentiation:

$$\dot{V} = F(V) \longrightarrow V(t_i) = \int_{t_0}^{t_i} F(V) dt$$

- two different heuristic functions:
 - SSE = sum of squared errors
 - MDL = SSE + equation length penalty

lagrange: summary

- advantages of using declarative bias in discovering dynamics:
 - can use different search strategies (e.g., beam search)
 - can discover complex models
 - can use background knowledge from the study domain
 - and thus compensate for incomplete (and noisy) data
- Q: where do the grammars come from?

grammar source 1: domain-specific modeling knowledge

- knowledge organized around central notion of process
 - identify basic processes in the domain
 - what models are used for individual processes?
 - how are they combined into models of the entire system?
- similar to compositional modeling [Falkenheiner & Forbus, 1994]
 - model fragments + rules for combining them
 - model fragments = models of individual processes
 - building model = search for appropriate combination of fragments
 - support for building QUALITATIVE models only

example: population dynamics modeling knowledge

- collected from textbooks in modeling of biological systems
- population: group of individuals of the same species inhabiting the same area
- population dynamics: change of the population density
- basic population dynamics processes:
 - population growth, population decay
 - interaction between populations

example: simple population dynamics model

- Lotka-Volterra population dynamics model:

$$\begin{aligned}\dot{N} &= a \cdot N - b \cdot NP \\ \dot{P} &= c \cdot NP - d \cdot P\end{aligned}$$

- general scheme:

$$\begin{aligned}\dot{N} &= \text{growth}(N) - b \cdot \text{feeds_on}(P, N) \\ \dot{P} &= c \cdot \text{feeds_on}(P, N) - \text{decay}(P)\end{aligned}$$

- assumes:

- unlimited prey growth – $\text{growth}(N) = aN$
- unlimited predator decay – $\text{decay}(P) = dP$
- unlimited predator capacity – $\text{feeds_on}(P, N) = PN$

relaxing the assumptions: models of population growth

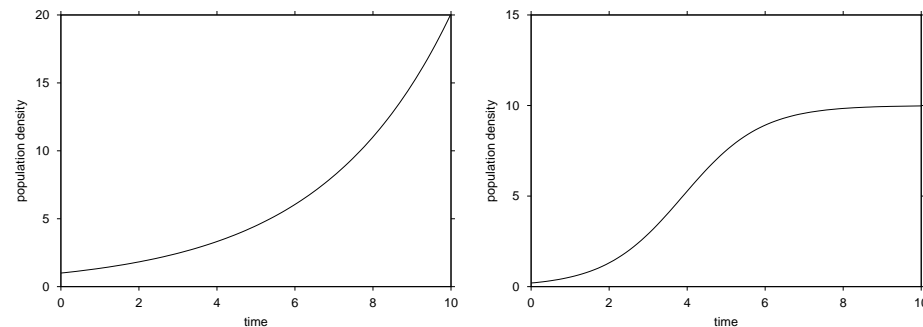
- exponential (unlimited) growth:

$$\text{growth}(N) = aN$$

- logistic (limited) growth:

$$\text{growth}(N) = aN(1 - N/K)$$

K is carrying capacity of the environment



models of predator-prey interaction

- unsaturated (unlimited) predation capacity:

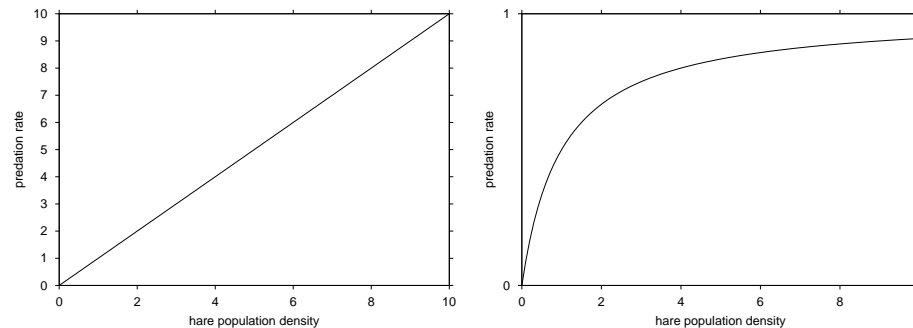
$$\text{feeds_on}(P, N) = aPN$$

- saturated predation capacity:

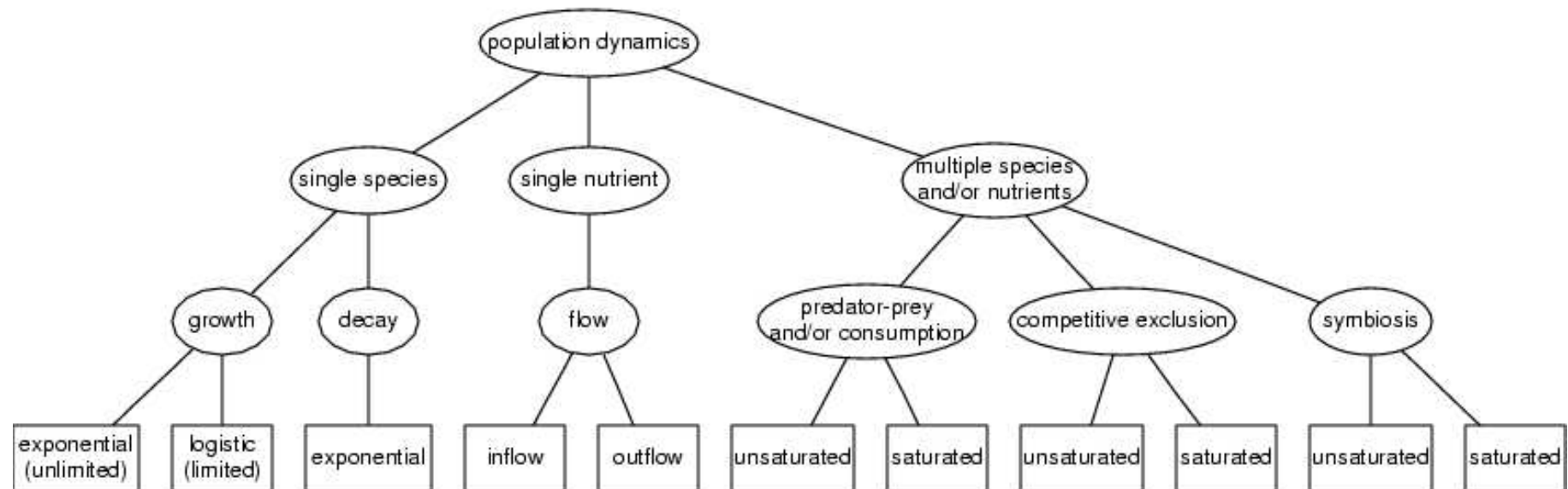
$$\text{feeds_on}(P, N) = P \cdot \text{saturation}(N)$$

$$\text{saturation}(N) = AN/(N + B)$$

A is the limit of the predation saturation, B is the saturation rate



domain-specific knowledge (1): taxonomy of process classes



domain-specific knowledge (2): models of the individual processes

- alternative sub-model templates for modeling individual processes:

```
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(capacity,0,1,Inf))
```

```
process class Feeds_on(Population p, set of Concentration cs)
  condition p  $\notin$  cs
  expression p *  $\prod_{c \in cs}$  Saturation(c)
```

domain-specific knowledge (3): combining scheme

- combining models of individual processes into a model of the entire system:

combining scheme Population_dynamics(Inorganic i)

$$\begin{aligned} \frac{d}{dt}i &= + \text{Flow}(i) \\ &\quad - \sum_{\text{food}, i \in \text{food}} \text{const}(_, 0, 1, \text{Inf}) * \text{Feeds_on}(p, \text{food}) \end{aligned}$$

combining scheme Population_dynamics(Population p)

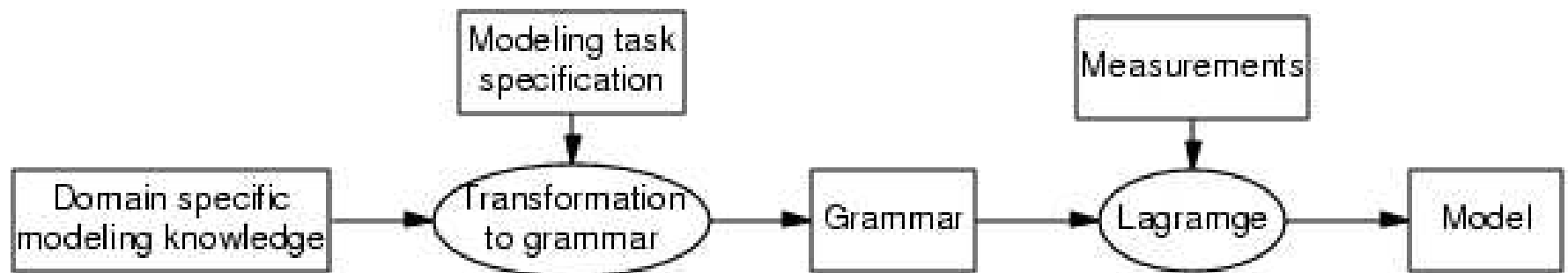
$$\begin{aligned} \frac{d}{dt}p &= + \text{Growth}(p) - \text{Decay}(p) \\ &\quad + \sum_{\text{food}} \text{const}(_, 0, 1, \text{Inf}) * \text{Feeds_on}(p, \text{food}) \\ &\quad - \sum_{\text{pred}, \text{food}, p \in \text{food}} \text{const}(_, 0, 1, \text{Inf}) * \text{Feeds_on}(\text{pred}, \text{food}) \end{aligned}$$

modeling task specification: Lotka-Volterra

```
variable Population hare  
variable Population lynx  
  
process Growth(hare)  
process Decay(lynx)  
process Feeds_on(hare, lynx)
```

- two system variables: hare (prey) and lynx (predator)
- three processes: growth, decay, and predator-prey interaction

integrating knowledge in the process of model induction



1. from task specification and knowledge to grammar
2. using the grammar for equation discovery with lagrange

domain-specific knowledge transformed into grammar

- context-dependent grammar for equation discovery:

```
general_lotka_volterra ->
```

```
    time_deriv(hare) = Growth(hare) - const[0:1:] * Feeds_on(lynx,hare)
```

```
    time_deriv(lynx) = const[0:1:] * Feeds_on(lynx,hare) - Decay(lynx)
```

```
Growth(hare) -> const[0:1:] * hare
```

```
Growth(hare) -> const[0:1:] * hare * (1 - hare / const[0:1:])
```

```
Decay(lynx) -> const[0:1:] * lynx
```

```
Feeds_on(lynx,hare) -> lynx * Saturation(hare)
```

```
Saturation(hare) -> hare
```

```
Saturation(hare) -> hare / (hare + const[0:1:])
```

- note the:
 - context-dependent constraint
 - bounds on the constant parameters

Lagoon of Venice: task specification

variable Inorganic temp, DO, NH₃, NO₃, PO₄

variable Population biomass

process Growth(biomass) biomass_growth

process Decay(biomass) biomass_decay

process Feeds_on(biomass, *) biomass_grazing

- six observed variables (only biomass modeled)
- two fixed processes
- one process template (unknown limiting factors)

Lagoon of Venice: results

- relatively high model errors (rmse – 86.2841; 157.537)
 - due to high measurement errors (order 20% – 50%)
 - predict most of the biomass concentration peaks and crashes
- comprehensible white-box models induced:
 - ecology expert can easily understand them
 - they reveal the limiting factors for algae growth
 - dissolved oxygen, nitrogen-based nutrients, and temperature

Lagoon of Venice: results

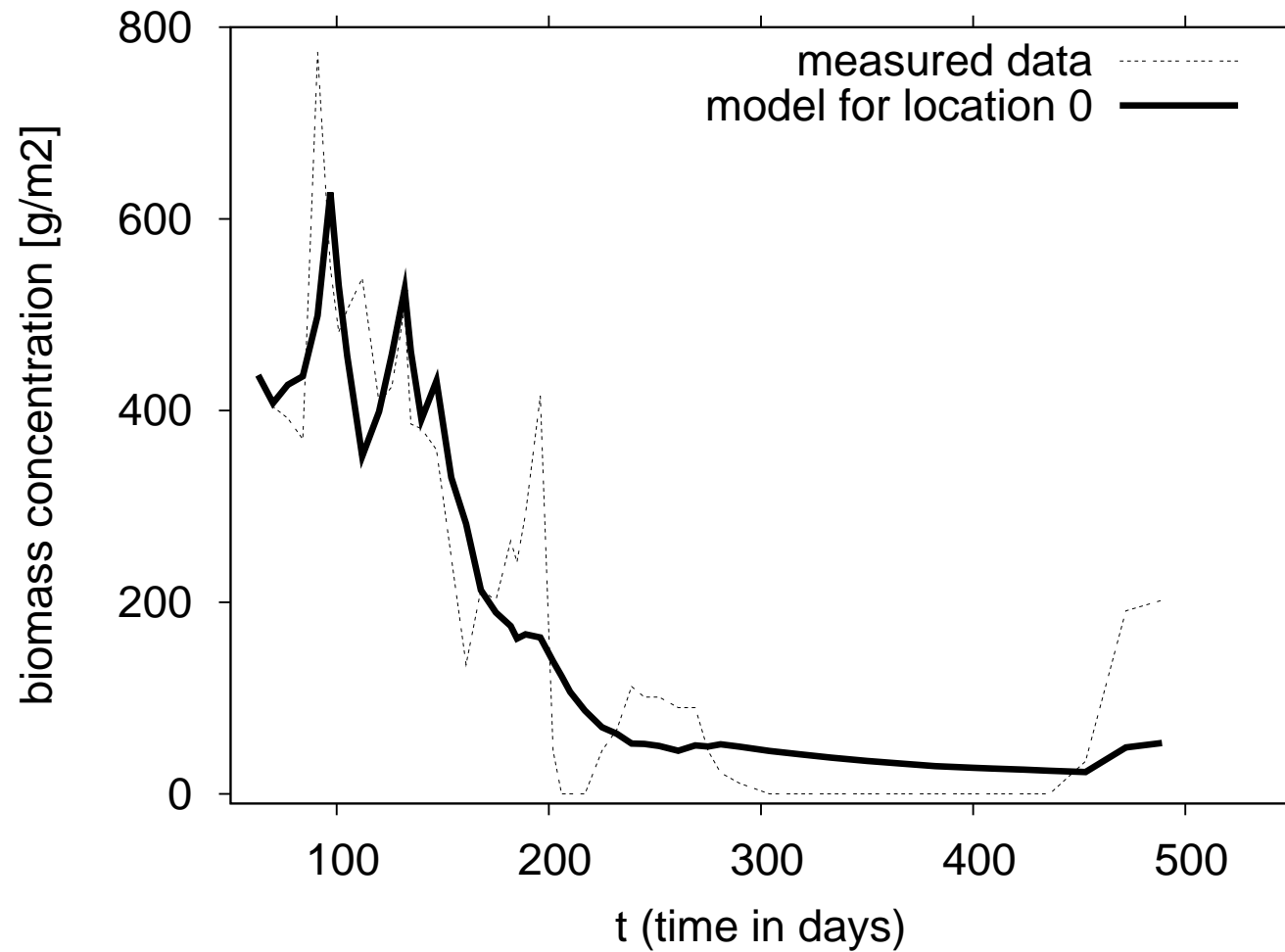
location 0:

$$\begin{aligned} \dot{\text{biomass}} = & 6.17 \cdot 10^{-5} \cdot \text{biomass} \cdot \left(1 - \frac{\text{biomass}}{1.80}\right) \\ & + 3.01 \cdot 10^{-4} \cdot \text{biomass} \cdot \text{D0} \cdot \frac{\text{N03}}{\text{N03} + 6.28} - 0.0319 \cdot \text{biomass}. \end{aligned}$$

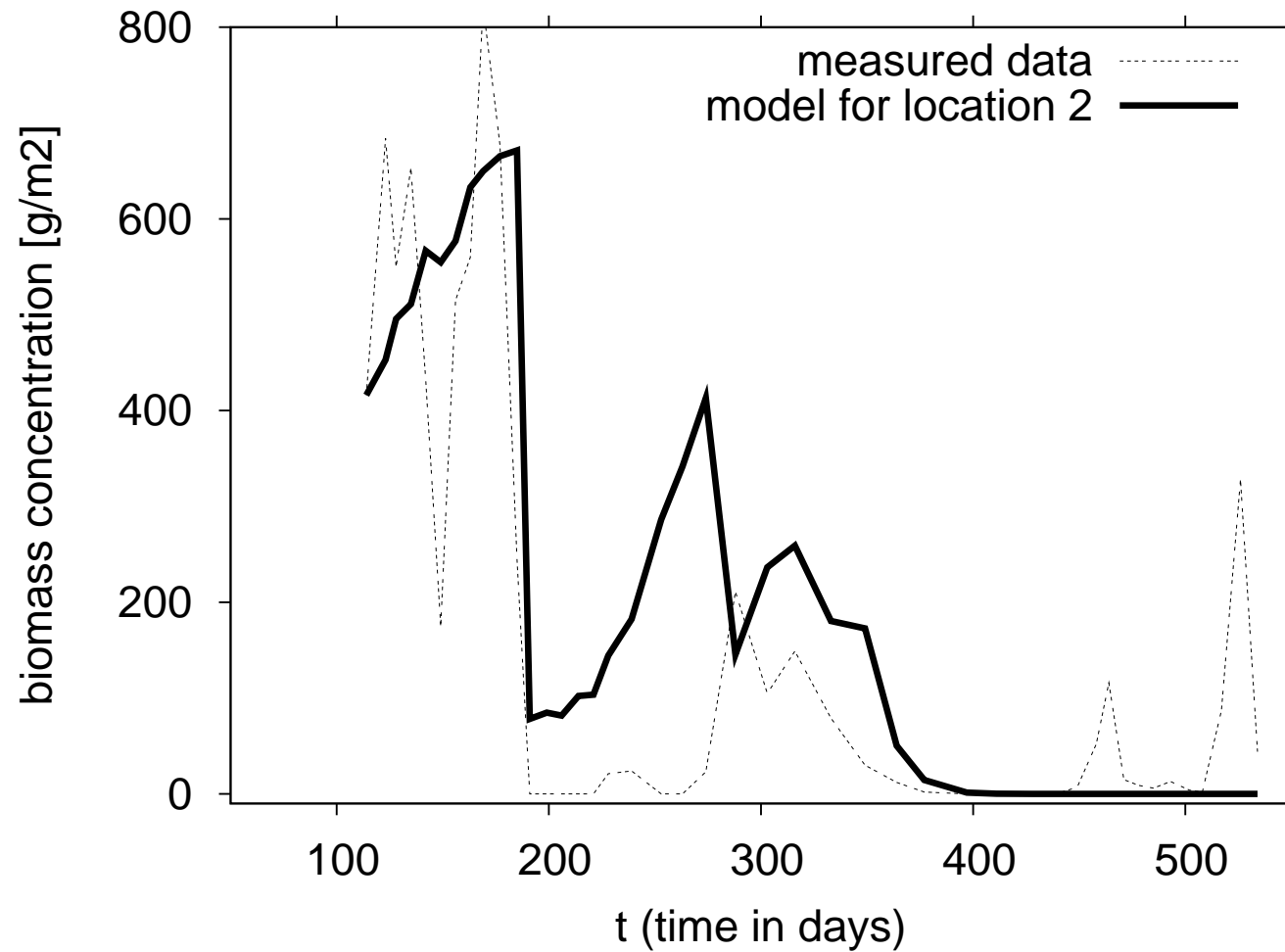
location 4:

$$\begin{aligned} \dot{\text{biomass}} = & 4.79 \cdot 10^{-5} \cdot \text{biomass} \cdot \left(1 - \frac{\text{biomass}}{0.844}\right) \\ & + 0.406 \cdot \text{biomass} \cdot (1 - e^{-0.216 \cdot \text{temp}}) \cdot (1 - e^{-0.413 \cdot \text{D0}}) \cdot \frac{\text{NH3}}{\text{NH3} + 10} \\ & - 0.0343 \cdot \text{biomass}. \end{aligned}$$

Lagoon of Venice: simulation of the induced model



Lagoon of Venice: simulation of the induced model



Lake Glumsoe

task specification

variable Inorganic temp, nitro, phosp

variable Population phyto, zoo

process Decay(phyto) phyto_decay

process Feeds_on(phyto, *) phyto_grazing

process Feeds_on(zoo, phyto) zoo_grazing

discovered model

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

grammar source 2: Ringkøbing fjord

- expert specified only part of the model structure:

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

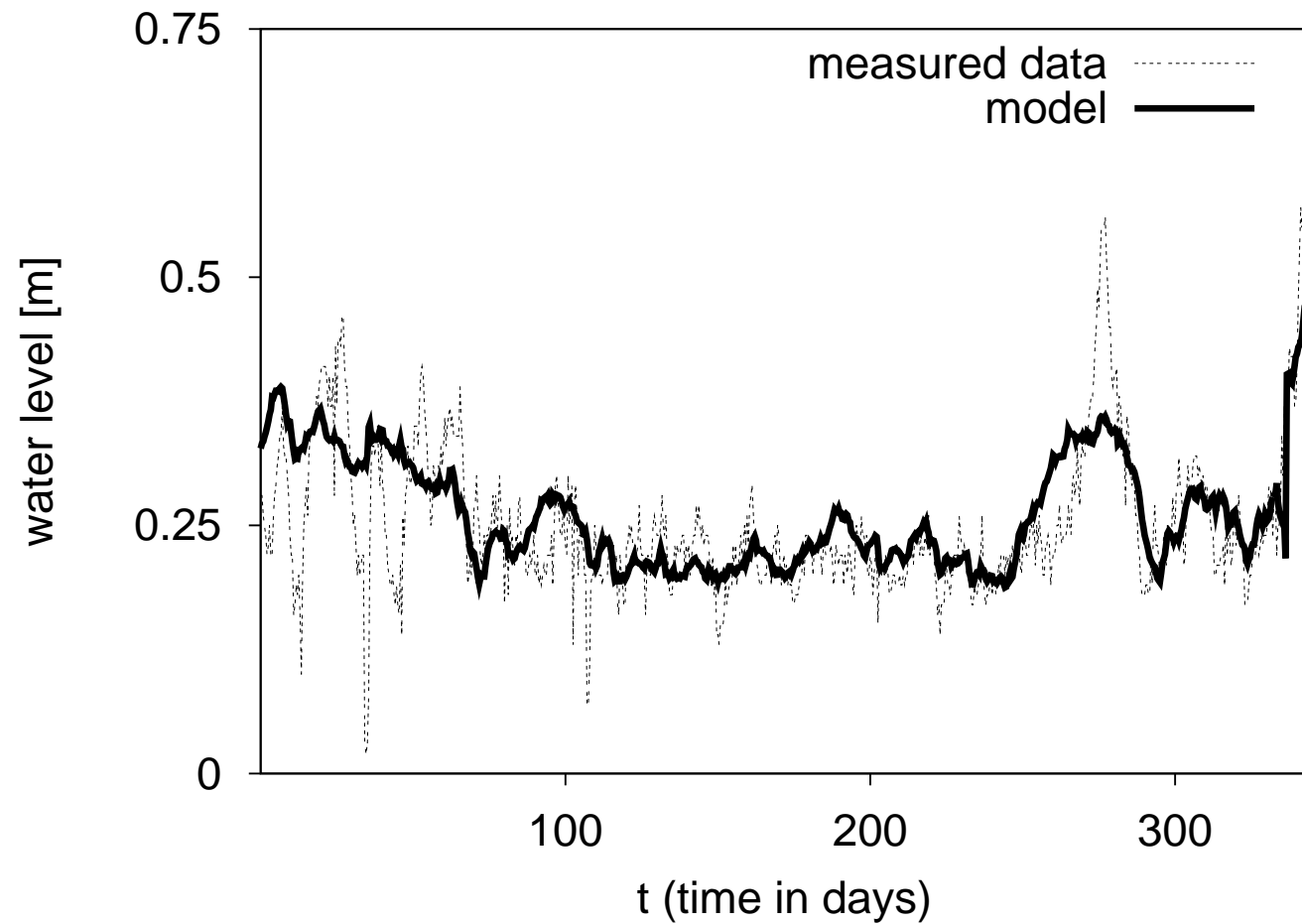
- two parts of the structure left unspecified:
 - gate opening f that depends on number of open gate parts
 - wind forcing g that depends on wind direction and speed
- experiments:
 - f – constant or polynomial
 - g – constant, polynomial, and trigonometric functions
 - maximal polynomial degree of 5

Ringkøbing fjord: results

task specification	training rmse	CV rmse	#CMS
constants	0.0848	0.106	1
polynomials	0.0655	0.0931	378
polynomials + sin/cos	0.0585	0.0903	2184
no partial structure (black-box)	0.0556	2.389	2801

- best model captures the long-term dynamics of the water level
- can also predict short-term changes (one hour or day)
- allows for comparison of wind and gate opening influence
- black-box polynomials overfit the training data

Ringkøbing fjord: simulation of the induced model



grammar source 3: existing models

- GIVEN:
 - an existing (imperfect) model M_0
 - a set of new observations/measurements
- FIND a revised model M_R that
 - minimizes the discrepancy between observed and measured values of the system variables
 - is as similar as possible to M_0

outline of the approach

- take an existing model
- construct a grammar that derives the model
- identify the unreliable parts of the model (expert)
- add alternative grammar rules for these parts (expert)
- use lagrange on the observations and grammar
- preference for models similar to the original one (minimality of change principle)

revising a part of a global vegetation model: CASA-NPPc initial

$$NPPc = \max(0, E \cdot IPAR)$$

$$E = 0.389 \cdot T1 \cdot T2 \cdot W$$

$$T1 = 0.8 + 0.02 \cdot topt - 0.0005 \cdot topt^2$$

$$T2 = 1.1814 / ((1 + \exp(0.2 \cdot (TDIFF - 10))) \cdot (1 + \exp(0.3 \cdot (-TDIFF - 10))))$$

$$TDIFF = topt - tempc$$

$$W = 0.5 + 0.5 \cdot eet / PET$$

$$PET = 1.6 \cdot (10 \cdot \max(tempc, 0) / ahi)^A \cdot pet_tw_m$$

$$A = 0.000000675 \cdot ahi^3 - 0.0000771 \cdot ahi^2 + 0.01792 \cdot ahi + 0.49239$$

$$IPAR = FPAR_FAS \cdot monthly_solar \cdot SOL_CONV \cdot 0.5$$

$$FPAR_FAS = \min((SR_FAS - 1.08) / srdiff, 0.95)$$

$$SR_FAS = (1 + fas_ndvi / 1000) / (1 - fas_ndvi / 1000)$$

$$SOL_CONV = 0.0864 \cdot days_per_month$$

proposed alternatives

- experts identified four "weak" parts of the model:
 - equations for E , $T1$, $T2$, and SR_FAS
- two alternatives for $E = 0.0389 \cdot T1 \cdot T2 \cdot W$:
 1. $E = \text{const} \cdot T1 \cdot T2 \cdot W$
 2. $E = \text{const} \cdot T1^{\text{const}} \cdot T2^{\text{const}} \cdot W^{\text{const}}$
- two alternatives for $T1 = 0.8 + 0.02 \cdot topt - 0.0005 \cdot topt^2$:
 1. $T1 = \text{const} + \text{const} \cdot topt + \text{const} \cdot topt^2$
 2. $T1 \rightarrow \text{const} \mid \text{const} + (T1) * topt$

the revised CASA-NPPc model

$$NPPc = \max(0, E \cdot IPAR)$$

$$E = 0.402 \cdot T1^{0.624} \cdot T2^{0.215} \cdot W^0$$

$$T1 = 0.680 + 0.270 \cdot topt - 0 \cdot topt^2$$

$$T2 = 1.1814 / ((1 + \exp(0.2 \cdot (TDIFF - 10))) \cdot (1 + \exp(0.3 \cdot (-TDIFF - 10))))$$

$$TDIFF = topt - tempc$$

$$TDIFF = topt - tempc$$

$$W = 0.5 + 0.5 \cdot eet / PET$$

$$PET = 1.6 \cdot (10 \cdot \max(tempc, 0) / ahi)^A \cdot pet_tw_m$$

$$A = 0.000000675 \cdot ahi^3 - 0.0000771 \cdot ahi^2 + 0.01792 \cdot ahi + 0.49239$$

$$IPAR = FPAR_FAS \cdot monthly_solar \cdot SOL_CONV \cdot 0.5$$

$$FPAR_FAS = \min((SR_FAS - 1.08) / srdiff, 0.95)$$

$$SR_FAS = (1 + fas_ndvi / 750) / (1 - fas_ndvi / 750)$$

$$SOL_CONV = 0.0864 \cdot days_per_month$$

- relative accuracy improvement 9%, revised model simpler

a brief summary of the talk

- integration of different aspects of knowledge:

- taxonomy of basic processes
- partial specification of the model structure
- existing models
- ...

in equation discovery through grammars (lagramge)

- grammars can be

- obtained by transforming the domain-specific knowledge
- in terms of textbook modeling knowledge
- or provided by human expert

further work

- establishing libraries of
 - modeling (process-based) knowledge in different domains
 - existing models
- application/evaluation on other real-world tasks
 - better parameters fitting procedure
 - modeling population dynamics in lake Bled
- inducing domain-specific knowledge (grammars):
 - from data
 - from existing models
 - from textbooks and articles on modeling