# On validation of individual-based models

Jochen Wittmann

wittmann@informatik.uni-hamburg.de
Universität Hamburg, Fachbereich Informatik, Arbeitsbereich Technische Informatiksysteme
Vogt-Kölln-Straße 30, 22527 Hamburg, Germany

**KEYWORDS**

model validation, individual-based models, experimental design

**ABSTRACT**

This paper analyses typical experimental set-ups for individual-based models on a not-aggregated level of model description in comparison to conventionally aggregated models.

It postulates that for real-world-applications additional assumptions become necessary which concern to the type and the parameters of the data transformation between the aggregated and the non-aggregated level.

The structure of the problem is analysed and typical scenarios for model usage and validation are listed. General methodological deliberations for each of these scenarios are made which offer a guideline for correct experimental design in order to validate the corresponding models.

## 1. Introduction

The object oriented modeling paradigma has established during the last years and leads especially for the application areas biology, sociology to its specialization in the form of so called „individual-based" models. This paper will not go into further discussions on the definitions of "individual-based" in contrary to "individual-oriented" or even "agent-based" modelling. A comprehensive summary concerning this topic can be found in (Ortmann 1999). However, the paper will analyse the validation step during a simulation study if an individual-based approach has been chosen.

With regard to the main application areas of individual-based models, which mainly are applied in domains without exact physically derived model descriptions, this important phase in a simulation study attracts special attention.

It is typical for individual-based models to model the reality, the objects under observation, and their behaviour in a very natural way by close analogy between the real world objects and the objects – or individuals – used on model description level. Therefore, this modelling

paradigm leads to a class of models which satisfy the criteria of adequate and easy understandable model structure on a very high level.

On the other hand, these models often are associated with the disadvantages caused by their demands concerning processor time and memory. This problem is a direct consequence of the non-aggregated model description and seems to be the price the user has to pay for comprehensibility and transparency on model specification level. This problem, too, shall not be discussed here.

This paper will focus on a further problem field which seems to be neglected in the main discussion of individual-based models: the problem of model validation. The pretended simpleness in model description often implies the need for a highly sophisticated analyse of the model and its results in the phase after the runs, in validation and interpretation of the results. In this situation, this paper analyses typical dilemmas and tries to give hints for a proper determination of the range of validity for individual-based models.

## 2. Simulation on local and on global level

To understand the problems concerning validation, we start with a view on the general design of a modelling and simulation study based on the individual-based paradigm. Figure 1 depicts the course of the argumentation in comparison to the use of a "conventional", i.e. non-individual-based model.
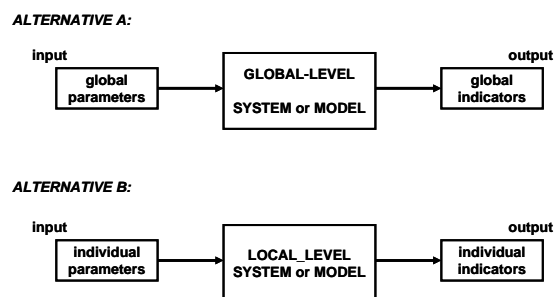


Fig. 1: Comparison between individual-based and non-individual-based modelling studies

Both alternatives work according to the same basic scheme: Alternative A shows the situation for a conventional model on global, which means here accumulated, level in model specification. The modeller and experimentator is interested in the effects of a change in a global parameter. This parameter is set for the simulation and after the run an other parameter on global level, a global indicator variable is observed.

Example: global input parameter is the reproduction rate of a population, the model is a common differential equation model for the population dynamics, and the model result is the population for a future point in time. Input, output, and model equations work on highly aggregated data for the population, which mirror the situation on individual level in statistical sense.

On the other hand, alternative B describes the system dynamics on the individual level. Example: For the population dynamics, a possible input parameter would be the mean number of children a woman gets during her life, one would have to model the interactions of the individuals and would be able to derive an individual curriculum vitae for each of the individuals. At the end, the actual number of children each individual has got would be the observation parameter on this level.

Both alternatives are proper implementations of the same basic modelling and simulation approach. The experiment deals with the objects input variable, the model itself, and the output. Accordingly the three basic tasks are identified: system identification (input and output given), forecast (input and system model given), and control (system model and output given). Differences between the alternatives A and B can only be found on the level of model description: In the first case, the complete model is specified using the population number as a cumulated value. The second case specifies the behaviour of the individuals and produces the population number as a dependent variable of the set of interacting individuals. Naturally, both model approaches have to be parameterised and validated on their specific level of model description. In consequence, even the results can only be interpreted and exploited on the level of specification the model offers.

# 3. Data transformation between the levels

There is one observation which appears from the simple description of the experimental set-up described so far: During the simulation run an individual-based model produces the curriculum vitae of the set of individuals under observation. If the experimenter is interested in more general model quantities, a recalculation and evaluation of those raw data will be necessary. (In our very simple example this recalculation step is realized by a simple summation of the individuals living at a

certain point of time and could be realised as a dependent model quality as well.) This argumentation implies a change of modelling level for data evaluation and interpretation (i.e. from level A to level B) concerning the two alternative scenarios introduced in figure 1.

Similar and much more complicated transformations from one level to the other can be necessary in a number of simulation experiments which deal with individual-based models.

In general, the change of levels is usefully applied if missing information on the one level is replaced by or can be derived from well known information on the other level. Such a level change can be done on the input-side as well as on the side of the outputs.

So far there are no problems in the experimental set-up and the situation can be recapitulated graphically by figure 2. The difficulties, however, arise when the model has to be validated and the situation escalates if there is a lack of comprehensive system data.
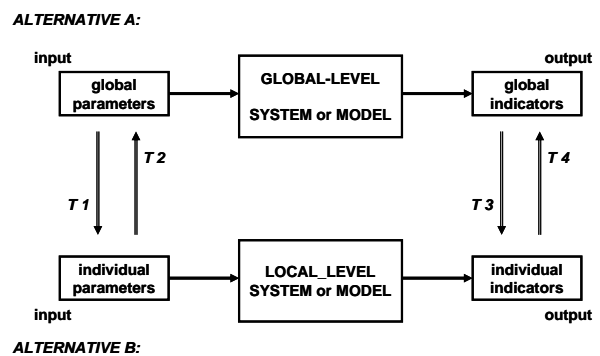


Fig.2: possible transformations between the levels of model description during experimentation

Example: For the very simple population dynamics example the transformations T1 to T4 introduced by the figure shall be exemplified:

1. A known mean life expectancy is transformed into determined ages for set of identical model individuals.

2. In a statistical sample size and weight of persons are measured, the mean values are used as parameters for a model on global level.

3. A certain mean value for energy consumption of a region has to be allocated to individual energy consumption values for each individual living in the region.

4. The total population number is summarised by counting the model individuals at a certain point of time.

Usually, the transformations from the individual level to the global level are evident and easily to execute. In this direction, there exist data on detail

level, which have to be aggregated to a more general, often statistical parameter value on the global level.

Transformations in the other direction are not possible without at least two further assumptions:

1. the type of distribution of the parameter transformed (e.g. uniform, normal, ...)

2. parameters of the distribution, such as mean value, variance, ...

But even the very simple transformation of type 4 (individual level to global level) can be more than a simple summation and has to be considered with carefulness. An example: The individually collected voices during an election could be weighted. Therefore an additional set of weight-parameters has to be specified for the model and the corresponding aggregation function has to be calculated for a correctly executed level change.

## 4. The problem

The argumentation so far explains the theoretical design of simulation experiments on the both levels introduced. However, in praxis and especially in the praxis of the application domains which like to use individual based models of cause of their easy and structure adequate model description facilities, the missing data forces to a more sophisticated, combined experiment design crossing the levels. Therefore transformations become necessary which imply additional parameters.

The systematic problem of these parameters is that their values cannot be acquired separately. If it would be possible to do so, the transformation and the level change would not have been necessary.

On the other hand, proper parameter identification needs measurements on both levels to identify the transformation parameters first and to calculate their values afterwards. This is an inherent contradiction of the experimental design. It is caused by the situation of system data and will not be dissolved by additional data acquisition in the real system.

For the modelling and simulation study follows: A separate validation of the assumptions concerning transformation parameters and their values is not possible. They have to be an additional task within the global model validation process.

To formulate constructively: The model experiments have to be designed in a manner that

1. the model results are independent of these transformation parameters, or

2. there is a proper distinction between the influence and effects of the transformations and their parameters and the effects of a change in the model parameters which in fact are under observation to achieve the experiments objectives.

In both cases the validation implies additional restrictions for the experimental design. The experiments have to assure that a statistical distinction between the effects of the transformations and those of the intended classical investigation according to the tasks identification, forecast, and control becomes possible.

Naturally, this problem escalates because even in the model there are variations in parameters to test, which are caused by uncertainties concerning model parameter values and even model structure. Figure 3 concludes these possibilities in argumentation for the different alternatives in experimental design.

It is obvious that the additional parameters make the study much more complex and the intended direct causality between the experimental parameters and their effects becomes more and more difficult to extract.
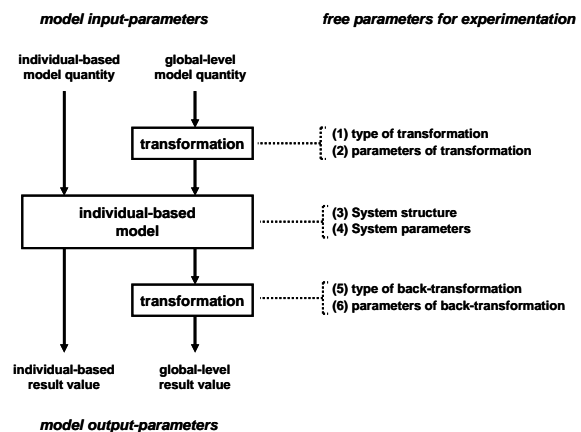


Fig.3: data-flow and free experimental parameters

## 5. Possible experimental designs for validation

So far, the need for sophisticated statistical methods for validation has been elaborated. Furthermore it is obvious that it will not be possible to validate the additional parameters separately, because there are no (or at least: not enough) system data.

In this situation, four possible and typical experimental designs shall be analysed with regard on a feasible model validation. The objective is to demonstrate the general argumentation and to explain the logical consequences of the initially chosen experimental design.

## 5.1 Individual behaviour under observation

The most obvious motivation for building individual-oriented models is to investigate in the behaviour of the individuals itself. This is represented by alternative B from figure 1. The experimental design is without any modification as it is usual in modelling and simulation because all operations take place on individual level. For the validation system data and model data have to be compared and the range of validity has to be determined from these deliberations. Concerning the structure of system and model equal assertions can be made, and the free parameters (numbers (1), (2), (5) and (6)) from figure 3 are not relevant in this case.

However, one should pay attention to the format of simulation results: To be accurate, only the life data for individuals are observed on the individual level. There is no aggregation of the data at all. Any aggregation would be interpreted as a change to the global level and would imply the necessity of a transformation of type T4 with the corresponding parameters and difficulties.

These deliberations lead to the next experimental scenario:

## 5.2.1 Structural adequate models for global processes

The motivation for this design variant comes from model description methodology: There exists the presumption that a model code as well as a program code is easier to understand and more efficiently to maintain if its structure mirrors the real world structure of the modelled system. With this background, the individual-based model description seems to offer the optimal level of comprehensibility because this model specification paradigm propagates to be nearly completely adequate.

For the validation context, one interesting observation must be made in connection with this approach: Even though the interests of the experimentation lie on the global level, model description and simulation work with the non-aggregated level. Therefore the model holds a scale in detail, which is not necessary for the level of results the experimentation intends. If the information on the detail level can be provided, this approach is very self-explaining and the advantages of the evident model structure overweight the demands in run-time those models usually need.

If there is a lack of information concerning parameters on the individual level, there are lots of additional hypothesis concerning type and parameter values of the transformations to calculate and validate, a task which has to be solved by data collected on the aggregated level solely. Thus, a serious validation for this kind of models succeeds only with great efforts in statistical determination of the missing parameters. In praxis the modeller will have to weight whether the adequate model structure will be worth these investments in statistical procedures. These deliberations show that the evaluation of this experimental design scenario has to be made for each application separately. The balance between investments and effort as described above should be considered very carefully.

## 5.3 Measurements are not possible on the desired level of model description

This scenario is very similar to the preceding one; however, in this case the experimenter has no choice between the alternatives in level because a missing access to the data on the one level forces him/her to substitute the missing information by investigations on the other one.

To be able to parameterise, validate, and work with the model at all, at least one of the transformations has to be specified and parameterised. Here the efforts are the prize for capacity to act not only the prize for an adequate, a nice model structure.

The limitations concerning accuracy and validity of the model have to be accepted. The experimental design has to be very sophisticated but the way of additional transformations is the only one, which provides access to a region of knowledge otherwise completely inaccessible.

## 5.4 Investigations on emergent behaviour

Highly interesting is an application field for individual based models not yet mentioned in this paper so far: the so called "emergent behaviour". In short, this means a behaviour of a group or mathematically spoken a set of identical individuals which is observed when these individuals interact, communicate, and cooperate but which is no specified explicitly within the behaviour specification of the single individual (e.g. the organisation of the ants, swarms, ..).

It is evident that the use of individual based models is inevitable in this case. Here, the experiment focuses on one of our transformations: The purpose of the model is to describe individual behaviour locally, let the individuals interact, and to observe behaviour of the group of individuals which has not been specified explicitly on the local level. The change of level is the trick: input on local, measurement of output on global level.

A further analysis touches the assumption that has been the base for all the deliberations before: the existence of well-defined rules for aggregation.

This assumption is challenged by the assumption of emergent behaviour.

There is no transformation specification in the form of rules or functions! In contrary, the observations on global level are generated by the behaviour specification on local level exclusively.

So far the theory. In real world applications the investigations on emergent behaviour naturally are superposed by the problems in getting proper system data. Therefore, very often level transformations are necessary to avoid data lacks. These transformations have to be parameterised and validated as described before. To prove real evident behaviour properly it is inevitable to separate the transformation and its effects from the observations and investigations made to prove the emergent behaviour.

If the parameters of the transformation are not known, complex additional experiments are necessary to determine their effects first, and let the argumentation turn to the phenomena of emergent processes only if there are no more doubts concerning "technical" transformation parameters. Especially for validation these interacting effects have to be differentiated and isolated to make real causalities between local behaviour specification and global level parameters evident.

## 6. Concluding example

The well known predator-prey model shall serve as a very simple example to illustrate the problems and the argumentation for the different experimental set-ups.

Alternative A implements the model by the well-known set of two differential equations for the two populations. Alternative B specifies the same situation in an individual-based manner. The question has to be discussed, how information on the one level can be completed by data on the other level and how far the two levels provide support for validation.

First the (well known) suppositions for the differential equation model explicitly in advance:

1. The equations are valid only for large population numbers N.

2. The parameter values are based on equal distribution of the individuals on the field. (e.g. for the meeting probability)

To demonstrate the dilemma comparing individual-based and global model to each other, the following deliberations will be enough:

1. If the individual-based model is run with low population number N, there is a direct contradiction to the assumption 1 for the global model.

2. If the individual-based model is run with large population number N, there will be a contradiction with the assumption 2: If there are lots of individuals, the distribution over the area under observation will not be equal. Normally, there are groups of hunting predators with no prey in between them in one block and in an other region other groups of prey with no predators in between.

The consequence for the experimenter is now: Is the group building process just a mistake in model description or should it be interpreted as emergent behaviour? Often the answer of this question draws upon the data produced by the model on the other level. As explained, such an argumentation breaks the assumptions. There is no other way out than to specify the transformations between the levels, determine their parameters and validate the hypothesis on this statistically detailed level.

Concerning the validation of models by a second model of the same system but on an other level of detail the conflict is obvious as well: The change of model specification level does not replace detailed validation based on additional experiments with the model and normally even with the real world system.

## 7. Resume

The paper tries to give a structure to discuss the validation of individual based models by mentioning the separate data transformation steps within the global and the local modelling level and between the levels themselves.

It emphasises that each transformation has additional parameters for its own that normally have to be determined by additional statistical experiments.

A comparison of results gained by models on the different levels may be interesting, however, its statistical value for validation and interpretation of possibly appearing effects is negligible.

The proposed scheme does not provide an algorithm to solve the problems in using individual-based models but it tries to make the typical structures of argumentation using such models transparent and tries to give a guideline for the discussion of critical aspects and common problems using these types of models.

## References

Ortmann, Jörg 1999. „Ein allgemeiner individuenorientierter Ansatz zur Modellierung von Populationsdynamiken in Ökosystemen unter Einbeziehung der Mikro- und Makroebene"; Dissertation am Fachbereich Informatik, Universität Rostock, 1999