# Chapter 1

# Introduction

*God and data have one thing in common.*
*They both know the truth, but remain silent about it.*
*E.D.*

The advancements in computational power allow to model, simulate and optimize real–world applications and problems to a much a larger extent than in the previous decades. Instead of concentrating on isolated effects, one is now able to run through complex scenarios and predict their future development. This trend of computer driven process design and product optimization intensifies as one can easily incorporate or eliminate new effects as well as apply and compare different computational methods. Even though this development supports the idea of an automated modelling and discrimination toolbox, experiences in several emerging fields of application have shown that existing setups and procedures may fail to deliver the very type of reliable statements that are needed for further inferences.

In computational chemistry and biotechnology, for example, homogeneous systems like the mass action law allow to establish mathematical models of the chemical processes in terms of ordinary differential equations (ODEs). The equation's right side is a function of reactants and parameters which determines the dynamical behavior of the system. There, effects can be added, substituted and eliminated.

When constructing new models, the modeler is usually confronted with a set of known analytically given effects that are worth considering. Some users refer to this setting as an effect–modelling–toolbox, which is typical for the stage of model assembling after having conducted data analysis.

The ease of assembling new models by the procedure described above, suggests to intensify computer driven model generation and validation. However, an automation is only sensible if the individual effects and sub–processes of the model to be built are known and documented and if modelling can be conceived more or less as a reconstruction.

However, in the field of biokinetics for example, the experience with comparable processes is very limited and in many cases not available as quantitative knowledge. Instead, it is necessary to identify reasonable models and adapt the respective parameters on the basis of a rather limited number of experiments in a permanently shortening cycle time. At this very stage of assembling a model, the knowledge about such *intermediate models* is as low as the required details of model, compared to the in-depth and trustworthy models at the end of the modelling process.

The already mentioned simplicity of constructing new models results in a huge number of candidate models. As a consequence one – reluctantly or willingly – encounters model uncertainty in terms of a systematic model–data–deviation, which cannot be explained statistically. This is typical and unavoidable at this stage of modelling, namely prior to having conducted some type of model validation.

The way towards a validated and trustworthy model mainly requires the work with model ideas, model alternatives of comparable quality, rough checks of the reasonability of the modelling approach in question, parameter estimation for these intermediate models, or decisions about further experiments putting doubts on certain models. The methods used there should differ from the calibration ones, which rely on some established trust and credibleness and apply at the end of the modelling process.

In view of what has been said above, especially of the surfacing model uncertainty, the crucial questions posing are (1) how to judge deviations between model and experimental data and (2) how to discriminate between competitive models if a systematic model–data–deviation is expected. Compared to the well–established and frequently employed methods, where trust in the model in question has already been established, one has to decide on a different goodness–of–fit interpretation and on an adapted model discrimination strategy: *the model–data–overlap approach.*

The model–data–overlap, which is introduced in this thesis paper, analyzes and

validates the ability of the model to take on a range of values by incorporating model variability in the discrimination target functional and can be applied for intermediate models. This new entity, in the following referred to as *model variability* is associated, to parameter sensitivity, resulting in a distributed parameter interpretation.

The model–data–overlap therefore is a tool to analyze and access the model–parameter–structure. It is a suitable approach when exploring this very structure with respect to the measured data, when gaining insight knowledge and when discriminating intermediate models. As a consequence, the model–data–overlap closes the gap between the existing data analysis at the beginning and the calibration methods at the end of the modelling process(see figure 1.1).
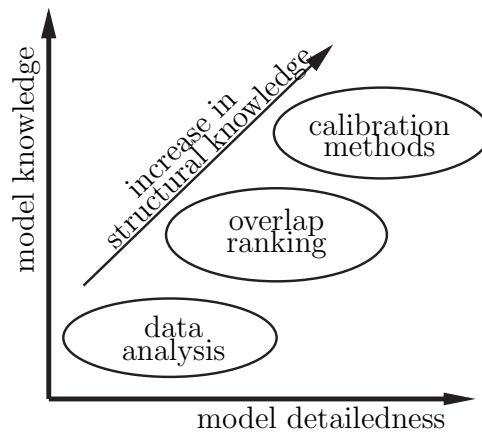


Figure 1.1: The overlap approach closes that gap between data analysis and calibration methods by analyzing the model–parameter–structure.

Generally spreaking, there are several strategic approaches to validate and discriminate models. According to JACOBS and GRAINGER in [127] as well as MYUNG and PITT in [180], one could for example investigate and check

(D1) *Falsifiability*: whether there exist potential observations that are incompatible with the model,

(D2) *Explanatory adequacy*: whether the theoretical account of the model helps to make sense of observed data but also established findings,

(D3) *Interpretability*: whether the components of the model, especially its parameters, are understandable and are linked to known processes,

(D4) *Faithfulness*: whether the model's ability to capture the underlying regularities comes from the theoretical principles the model purports to implement, not from the incidental choices made in its computational instantiation,

(D5) *Goodness–of–fit*: whether the model fits the observation data sufficiently well,

(D6) *Complexity and simplicity*: whether the model's description of the observed data is achieved in the simplest possible manner, or

(D7) *Generalizability*: whether the model provides a good prediction of future observations.

By assessing the model's ability to take on experimental data, the model–data–overlap follows the discrimination strategy (D1), (D3) as well as (D5). By that it shows again, that the overlap is a tool to analyze and to interpret the surfacing model–data–deviation for intermediate candidate models.

The different discrimination strategies (D1)–(D7) cannot be applied at once, since they differently validate model–data–deviations. Presently, there is no and is not going to be a general master strategy. At some point, the experimenter or modeler has to make some sort of assumption and interpretation. For example, one may want to access model complexity or to assume that the underlaying model is true. In the case of the model–data–overlap to be presented, the user has to accept that the parameters have to be interpreted as distributions in order to cope with structural model uncertainty.

Not only seems model discrimination to be an unsolved strategic problem, the same does apply for the implementation and algorithmic aspects. As each approach and each application class raises individual challenges, the discrimination strategy realization problems are equally important.

It is therefore not astonishing that the development of the model–overlap–concept took place within the inter– and intra–disciplinary research project: "Experimentally Controlled Discrimination of models, parameter estimation and overlap optimization."[1] The project was carried out in 2002/03 and was split up into three parts:

– The model–overlap–concept. A new approach to parameter estimation, model validation, selection and discrimination (Lorenz)

---

[1]Project number: 03SCM1B2 within the BMBF–framework (BMBF = Federal Ministry of Education and Research, Germany): "Mathematics in "New mathematical methods in industry and business services".

– Numerical aspects and challenges at implementing the model–overlap–concept in PREDICI KINETICS and PRESTO[2] and example (TELGMANN in [227])

– Dimension reduction of complex systems (DIEDERICS in [77])

The results of the project so far have been published in two publications: "Adaptive approach for nonlinear sensitivity analysis of reaction kinetics" [120] and "Discrimination of dynamical system models for biological and chemical processes" [160]. The second paper deals with the introduction of the model–data–overlap in general. The first one, focuses on a special implementation aspect of the overlap. It shows, how the existing TRAIL-algorithm (c.f. [119]) is adapted to a Fokker-Planck-setting in order propagate the model variability for the class of dynamical systems.

**Outline.** This thesis paper starts with a brief introduction of the model–data–overlap in chapter 2. Before giving arguments and motivations for its construction in chapter 4, existing concepts for model discrimination and parameter estimation are reviewed in chapter 3. The chapters 5 and 6 show the implementation and performance of the model–overlap–concept for selected scenarios, the first of the two chapters explains the algorithmic framework, the second one shows the numerical results.

---

[2]PREDICI and PRESTO KINETICS are trademarks by CiT–Dr.M. Wulkow Computing in Technology GmbH, Rastede, Germany, www.cit-wulkow.de

[3]Since October 2004 Professor and head at the Institute for Chemical Processes Engineering, University Stuttgart

strong support especially during the past two years, when I was already holding on a purely non–scientific position and was finishing this thesis paper exclusively in my spare time and vacation.

Berlin, December 2nd, 2005                                                  Sönke Lorenz