# Reaction rate ambiguities for perturbed spectroscopic data: Theory and implementation

Henning Schröder[a,c], Cyril Ruckebusch[b], Alexander Brächer[e], Mathias Sawall[a], Denise Meinhardt[a], Christoph Kubis[c], Sara Mostafapour[a,d], Armin Börner[c], Robert Franke[e,f], Klaus Neymeyr[a,c]

[a]*Universität Rostock, Institut für Mathematik, Ulmenstraße 69, 18057 Rostock, Germany*
[b]*U Lille, CNRS, LASIRE, Laboratoire de spectroscopie pour les interactions, la réactivité et l'environnement, F-59000, Lille, France*
[c]*Leibniz-Institut für Katalyse, Albert-Einstein-Straße 29a, 18059 Rostock, Germany*
[d]*Department of Chemistry, Shiraz University, 71454 Shiraz, Iran*
[e]*Evonik Performance Materials GmbH, Paul-Baumann Straße 1, 45772 Marl, Germany*
[f]*Lehrstuhl für Theoretische Chemie, Ruhr-Universität Bochum, 44780 Bochum, Germany*

## Abstract

The analysis of reaction systems and their kinetic modeling is important for both exploratory research and process design. Multivariate curve resolution (MCR) methods are state-of-the-art tools for the analysis of spectral series, but are also affected by an unavoidable solution ambiguity that impacts the obtained concentration profiles, spectra and model parameters. These uncertainties depend on the underlying model and the magnitude of the measurement perturbations.

We present a general theoretical approach together with a computational method for the analysis of the solution ambiguity underlying arbitrary kinetic models. The main idea is to determine all those model parameters for which the corresponding pure component factorizations satisfy all given constraints within small error tolerances. This makes it possible to determine bands of concentration profiles and spectra that reflect the underlying ambiguity and circumscribes the potential reliability of MCR solutions. False conclusions on the uniqueness of a solution can be prevented. The procedure can be applied as a post-processing step to MCR methods as MCR-ALS, ReactLab or others. The Matlab program code is freely accessible and includes not only the proposed ambiguity analysis but also an MCR hard-modeling approach. Application studies are presented for two experimental data sets, namely for UV/Vis spectra on the relaxation of a photoexcited state of benzophenone and for Raman spectra on an aldehyde formation process.

*Keywords:* spectral recovery, kinetic modeling, parameter identifiability, ambiguity of rate constants, factor analysis, nonnegative matrix factorization

## 1. Introduction

The identification and quantification of unknown species is a major task in the analysis of chemical reaction systems. The required data are (time) series of spectra. Any series of $m$ spectra with $n$ data channels each can be stored in a nonnegative matrix $D \in \mathbb{R}^{m \times n}$. For a chemical system with $s$ components the matrix $D$ has the rank $s$ in the ideal, noise-free case. The Lambert-Beer law justifies the existence of a nonnegative matrix factorization

$$D = CS^T \tag{1}$$

for which the columns of $C \in \mathbb{R}^{m \times s}$ can be assigned to the concentration profiles and the columns of $S \in \mathbb{R}^{n \times s}$ to the spectra of the $s$ pure components. It is well known that $C$ and $S$ are not uniquely determined only from $D$ [1, 2, 3]. Multivariate curve resolution (MCR) methods aim to resolve a chemical interpretable factorization within the set of all possible factorizations. Further constraints on the factors $C$ and $S$ can support the factor recovery process [1, 4, 5]. In this work we focus on the consistency of the factor $C$ with an optimally parameterized kinetic model. Such a problem is called a *constrained factorization problem*. Even under the constraint of an underlying kinetic model, a unique factorization cannot be presumed. This fact is well known especially for first-order reaction systems [6]. The

remaining factor ambiguity can be represented in the space of the parameters of the kinetic model, namely the reaction rate constants. This leads to the *set of feasible parameters* $\mathcal{K}^+$ [7]. Each element of this set is a vector of reaction rate constants. For each of them the associated model evaluation on the given time grid yields a factor $C$, which appears in a nonnegative matrix factorization $D = CS^T$ of the spectral data $D$. Typically, this construction is very restrictive as it implies the existence of a proper nonnegative matrix factor $S$. Originally, the set $\mathcal{K}^+$ of vectors of reaction rate constants is defined for noise-free and non-perturbed spectral data matrices $D$. However, the calculation of $\mathcal{K}^+$ also yields meaningful results for experimental data if the influence of perturbations and noise is small [8, 9].

In this paper we work with small error tolerances for the nonnegativity constraint on $S$ as well as for the model fit. Our aim is to define a generalization of the set $\mathcal{K}^+$ for perturbed data $D$. The relation to the idealized case is illustrated and a guideline for useful settings of the tolerance values is given. A novel parallelizable adaptive algorithm is introduced that can be used to compute $\mathcal{K}^+$ as well as its generalization. The program code of all implementations can be found on the FACPACK homepage[1] including possibly updated versions. The theoretical results are applied to two spectroscopic data sets.

### 1.1. Organization of the paper

Section 2 contains a short introduction to rate constant ambiguities for the noise-free case. A generalization for perturbed and noisy data sets is introduced in Section 3. A novel algorithm for the approximation of these parameter ambiguities of kinetic models is presented in Section 4. The results are applied to a UV/Vis and a Raman spectroscopic data set in Section 5. The final Section 6 contains a summary and a conclusion.

## 2. Introduction to parameter ambiguities

The coupling of MCR methods with kinetic models is a well-established approach in order to reduce solution ambiguities [10, 11]. In particular, a second or higher order kinetic model usually results in unique factors $C$ and $S$. For first-order models, on the other hand, often a solution set exists that contains significantly different factorizations. This has been observed for the parameterization of kinetic models for single wavelength measurements [12] and for the more general case of MCR problems [13]. The first theoretical analyses were carried out by Vajda and Rabits [6, 14]. A detailed mathematical investigation for arbitrary first-order models is given in [7]. Those parts in [7], that are important for the further understanding, are reviewed below.

### 2.1. Factor representation by model parameters

MCR methods can be coupled with kinetic modeling by evaluating the difference between the factor $C$, for example computed by (1), and the prediction $C^{\text{ode}}(k) \in \mathbb{R}^{m \times s}$ as computed by an evaluation of the underlying kinetic model. This difference is to be minimized along with further constraint functions, which for instance penalize negative entries of $C$ and $S$. Here $C^{\text{ode}}(k)$ is obtained by solving an initial value problem on a given time grid for the parameters of the given model. These parameters are represented with a parameter vector $k \in \mathbb{R}^q$.

Typically, the components of $k$ are the rate constants of the model. They are determined along with the matrices $C$ and $S$. Throughout the paper we denote a (general) solution of the constrained factorization problem by $C^*$, $S^*$ and $k^*$. Typically, the solution meets the following optimality conditions (or other comparable conditions)

$$\frac{\|D - C^*(S^*)^T\|_F}{\|D\|_F} \to \min, \quad \frac{\|C^* - C^{\text{ode}}(k^*)\|_F}{\|C^*\|_F} \to \min, \quad \frac{\|\min(C^*, 0)\|_F}{\|C^*\|_F} \to \min \quad \text{and} \quad \frac{\|\min(S^*, 0)\|_F}{\|S^*\|_F} \to \min.$$

Therein the Frobenius norm $\|\cdot\|_F$ is the square root of the sum of squared matrix elements. A solution can be computed by means of various chemometric tools of which MCR-ALS [15] and ReactLab [16] appear to be the most popular ones. In this paper the hard-modeling approach presented in [7, 8] is used whose Matlab program code can be found on the FACPACK homepage. Its numerical algorithm is based on a truncated singular value decomposition $D = U\Sigma V^T$

---

[1] www.math.uni-rostock.de/facpack/

with orthogonal matrices $U \in \mathbb{R}^{m \times s}$ and $V \in \mathbb{R}^{n \times s}$ as well as a diagonal matrix $\Sigma \in \mathbb{R}^{s \times s}$. The factorization (1) can always be represented in the form

$$D = \underbrace{U\Sigma T^{-1}}_{C} \underbrace{TV^T}_{S^T}$$

by means of a regular matrix $T \in \mathbb{R}^{s \times s}$. If we additionally assume the consistency with a kinetic model, then it is possible to represent $T$ as a function of the parameter vector $k$ underlying the present concentration profiles. For a given $k$ the matrix $T = T(k) = ((U\Sigma)^+ C^{\text{ode}}(k))^+$ is a solution of the least-squares problem

$$\min_{T \in \mathbb{R}^{s \times s}} \|U\Sigma T - C^{\text{ode}}(k)\|_F .$$

In words, $T(k)$ minimizes the difference between the concentration factor and the parameterized model in a least-squares sense. Based on this, the factors $C(k) = U\Sigma(T(k))^+$ and $S(k) = T(k)V^T$ can be defined in a low-dimensional way. Therein the superscript + denotes the pseudoinverse of a matrix. The solution $C^* = C(k^*)$ and $S^* = S(k^*)$ of the kinetically constrained factorization problem makes the computation numerically effective.

### 2.2. The sets of D-consistent parameters and feasible parameters

The previous section shows that a solution $C^*(S^*)^T$ of the constrained factorization problem is entirely determined by its vector of parameters $k^*$. For the case of ambiguous solutions it appears to be useful to consider the *set of D-consistent parameters*
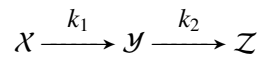
$$\mathcal{K} = \left\{ k \in \mathbb{R}^q : D = C^{\text{ode}}(k)S^T \text{ for a (not necessarily nonnegative) matrix } S \in \mathbb{R}^{n \times s} \right\} \tag{2}$$

and the *set of feasible parameters*

$$\mathcal{K}^+ = \left\{ k \in \mathbb{R}^q : D = C^{\text{ode}}(k)S^T \text{ for a nonnegative matrix } S \in \mathbb{R}^{n \times s} \right\}. \tag{3}$$

The set $\mathcal{K}^+$ contains only those $k$ respecting the nonnegativity constraints for $C$ and $S$. By contrast, the definition of $\mathcal{K}$ neglects the nonnegativity constraint for $S$. It has an auxiliary role and can be interpreted as an intermediate step towards the definition of its subset $\mathcal{K}^+$ in theoretical considerations and numerical computations.

The well-known slow-fast ambiguity [17] illustrates that these set definitions are meaningful. This ambiguity refers to the circumstance that the constrained factorization problem for the kinetic model

$$\mathcal{X} \xrightarrow{k_1} \mathcal{Y} \xrightarrow{k_2} \mathcal{Z}$$

may have two solutions. Let $k^* = (k_1, k_2)$ be such a solution. The second solution $k^{**} = (k_2, k_1)$ with swapped entries exists if the factor $S(k^{**})$ is nonnegative. Then the set $\mathcal{K}$ is given by $\{k^*, k^{**}\}$ and the set $\mathcal{K}^+$ is either equal to $\mathcal{K}$ or contains only $k^*$.

## 3. Rate constant ambiguities for perturbed and/or noisy data

The necessity for considering noise or perturbations is explained by a simple introductory example. A reaction system is assumed that is based on the reversible model

$$\mathcal{X} \underset{k_2}{\overset{k_1}{\rightleftharpoons}} \mathcal{Y} \tag{4}$$

with initial concentrations $(\mathcal{X}_0, \mathcal{Y}_0) = (1, 0)$. For the computation of the spectral mixture data $D$, see Figure 1, a concentration factor $C = C^{\text{ode}}(k^*)$ with $k^* = (2, 1)^T$ is used. Each column of the factor $S$ contains the evaluation
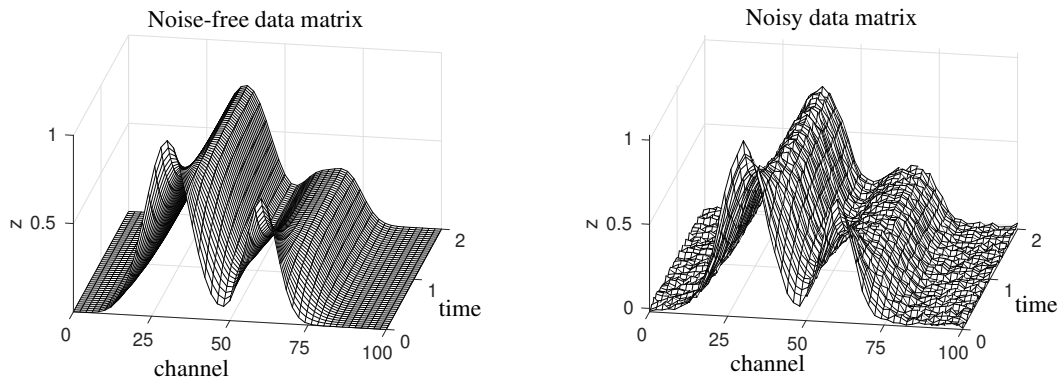
Figure 1: Two series of simulated spectra without (left) and with (right) noise are shown. The pure spectra of the noise-free case are modeled by two Gaussians each. The kinetic model $X \rightleftarrows Y$ with reaction rate constants $k_1 = 2$ and $k_2 = 1$ is used to generate the underlying concentration factor. For the noisy case 2% homoscedastic and 2% heteroscedastic noise have been added.

of two overlapping Gaussian curves. Next, the solution defined by $k^*$ (respectively $C^* = C(k^*)$ and $S^* = S(k^*)$) is considered to be the "true solution" and is to be determined. Furthermore with $\kappa = k_1^* + k_2^* = 3$ the sets

$$\mathcal{K} = \left\{ \begin{pmatrix} k_1 \\ \kappa - k_1 \end{pmatrix} : 0 < k_1 \leq \kappa \right\} \quad \text{and} \quad \mathcal{K}^+ = \left\{ \begin{pmatrix} k_1 \\ \kappa - k_1 \end{pmatrix} : 1.9984 < k_1 \leq \kappa \right\} \tag{5}$$

are known to describe the solution ambiguity. See Appendix A for details.

First, the idealized case of the noise-free matrix $D$ is analyzed. The kinetic hard-modeling approach [7] is used to solve the constrained factorization problem. It utilizes a minimization of a cost function that depends only on the parameter vector $k$. The results for 30 randomly chosen initial settings of $k$ are shown in the top row of Figure 2. The resulting optimized parameter vectors are shown in the left plot by gray crosses along with $\mathcal{K}$, $\mathcal{K}^+$ and $k^*$. They are located within the set $\mathcal{K}^+$. The center and right plots show the associated factors $C$ and $S$. As expected, it is not possible to identify $k^*$ with this method (or by any other MCR approach that uses only kinetic hard-modeling as additional constraint).

In the second case, we add 2% of homoscedastic and heteroscedastic noise to the model data. Similarly to the noise-free case, the corresponding MCR problem is solved by applying the hard-model approach. The solver for the underlying minimization problem is initialized with the same 30 vectors $k$ which were used in the noise-free case. The results are shown in the bottom row of Figure 2. In contrast to the noise-free case all obtained parameter vectors are located close to the point $(3, 0)$. Hence a user might assume that the problem can be solved uniquely. Even more critical is the fact that the obtained parameter vectors differ significantly from $k^*$, namely the "true" solution. The reason for this is mainly due to the calculated spectra of the second component. The red spectra in the two right plots of Figure 2 are compared. In the band of possible spectra indicated in the top plot, the spectra shown in the bottom are more distant from the baseline in a larger channel range; especially between 0 and 25. Thus, negative entries caused by noise will result in smaller residuals for the spectra shown in the bottom when applying an MCR approach.

Hence, we introduce generalizations of the sets $\mathcal{K}$ and $\mathcal{K}^+$ in order to overcome these difficulties. Their computation enables the user to decide whether an ambiguity of the model parameters has to be expected or not.

**Definition 3.1** (Set of $D$-approximate parameters). *Let a matrix $D \in \mathbb{R}^{m \times n}$ with $\mathrm{rank}(D) \geq s \geq 1$ and a truncated singular value decomposition $D \approx D_s = U \Sigma V^T$ be given. The set of $D$-approximate parameters is defined by*

$$\mathcal{K}_\varepsilon := \left\{ k \in \mathbb{R}^q : \mathrm{rank}(T(k)) = s \quad \text{and} \quad \frac{\|C(k) - C^{\mathrm{ode}}(k)\|_F}{\|C(k)\|_F} \leq \varepsilon \right\}$$

*with a tolerance $\varepsilon \geq 0$ allowing deviations from a perfect kinetic fit.*

The necessary condition $\mathrm{rank}(T(k)) = s$ guarantees that $D \approx D_s = C(k)(S(k))^T$ holds. Furthermore, an intentionally small kinetic fit error implies that $C \approx C^{\mathrm{ode}}(k)$. In summary, this leads to $D \approx C^{\mathrm{ode}}(k)(S(k))^T$ and thus a similar expression as in (2).
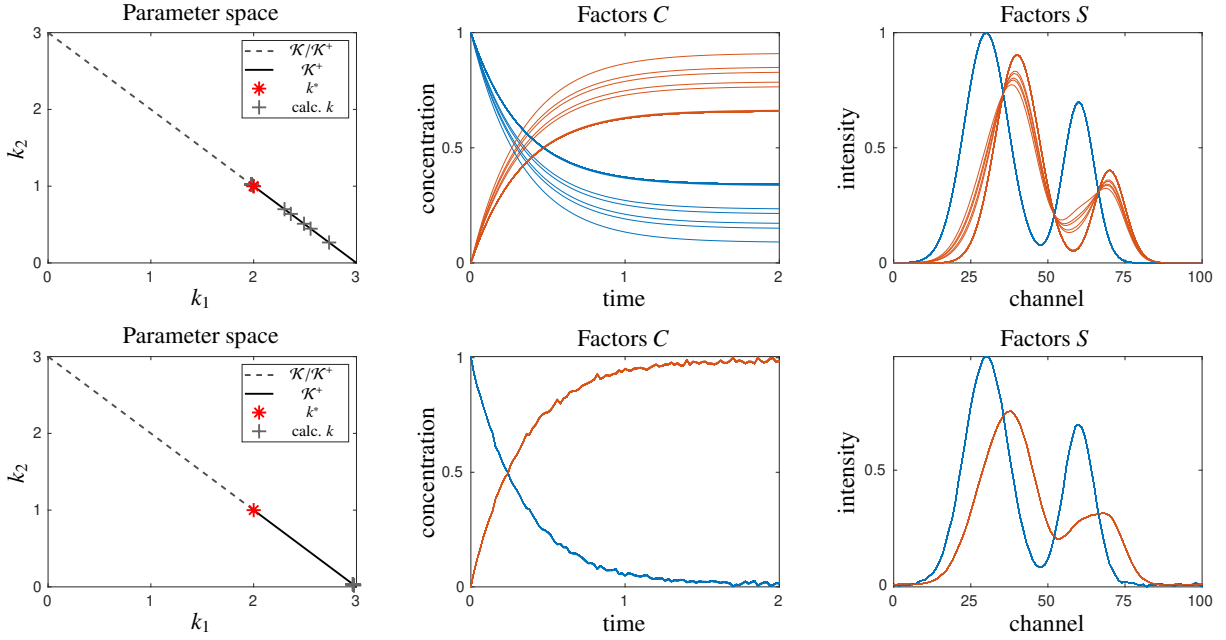
4

Figure 2: Results of the kinetic hard-modeling approach from [7] for the simulated data sets, shown in Figure 1; noise-free (top) and noisy case (bottom). Plots of the results in the parameter space are shown in the first column. The set $\mathcal{K}^+$ is plotted as solid line. It contains $k^*$ (red asterisk) that is used for data generation, as well as the optimal parameter vectors $k$ (grey crosses), that are obtained by kinetic hard-modeling for 30 different initializations of the underlying minimization. The same initializations have been used for the noise-free and noisy case. The factors $C$ and $S$ that correspond to the 30 optimal parameter vectors $k$ are shown in the second and third column.

**Definition 3.2** (Set of feasible *D*-approximate parameters). *On the assumptions of Definition 3.1 we define the set of feasible D-approximate parameters*

$$\mathcal{K}^+_{\varepsilon,\theta} := \left\{ k \in \mathcal{K}_\varepsilon : \frac{(S(k))_{i,j}}{\max_l(|(S(k))_{l,j}|)} \geq -\theta \text{ for all } i, j \right\}$$

*with a tolerance $\theta \geq 0$ that allows small negative matrix elements of $S$.*

The definition of $\mathcal{K}^+_{\varepsilon,\theta}$ only requires an approximate nonnegativity for the factor $S(k)$ (respectively $S$ in (1)) in terms of a relative lower bound for the matrix elements. To this end, each column of $S(k)$ is scaled to have the absolute maximum 1 and the scaled profiles are bounded from below by $-\theta$ for a small $\theta \geq 0$. Summarizing, the principle underlying the reduction of $\mathcal{K}_\varepsilon$ to $\mathcal{K}^+_{\varepsilon,\theta}$ is similar to the set reduction from $\mathcal{K}$ to $\mathcal{K}^+$ in the idealized case. The equalities $\mathcal{K} = \mathcal{K}_\varepsilon$ and $\mathcal{K}^+ = \mathcal{K}^+_{\varepsilon,\theta}$ hold for $\varepsilon = \theta = 0$ in the noise-free case. Again, $\mathcal{K}_\varepsilon$ is often used as an intermediate step in theoretical considerations and numerical computations.

Figure 3 shows an approximation of the set $\mathcal{K}^+_{0.01,0.01}$ in green for the model problem in combination with the previous results. The computational details for this approximation are explained in Section 4. The parameter vector $k^*$ as well as all calculated vectors $k$ are contained in $\mathcal{K}^+_{0.01,0.01}$. The set serves as a reliable indicator on a potential rate constant ambiguity.

In these computations it is important to work with meaningful values for the error tolerances $\varepsilon$ and $\theta$. The authors recommend to select them according to the following guidelines:

- The error tolerances $\varepsilon$ and $\theta$ should have the same magnitude as the signal-to-noise ratio in $D$. If major perturbations due to background corrections or other systematic errors are present, then the values can be set higher. We recommend values less than three times the signal-to-noise ratio of the data $D$. A good initial guess is the absolute value of the minimal entry of $D$.

- A more systematic approach for the determination of $\varepsilon$ is given below: Compute $\mathcal{K}$ as if there were no pertur-
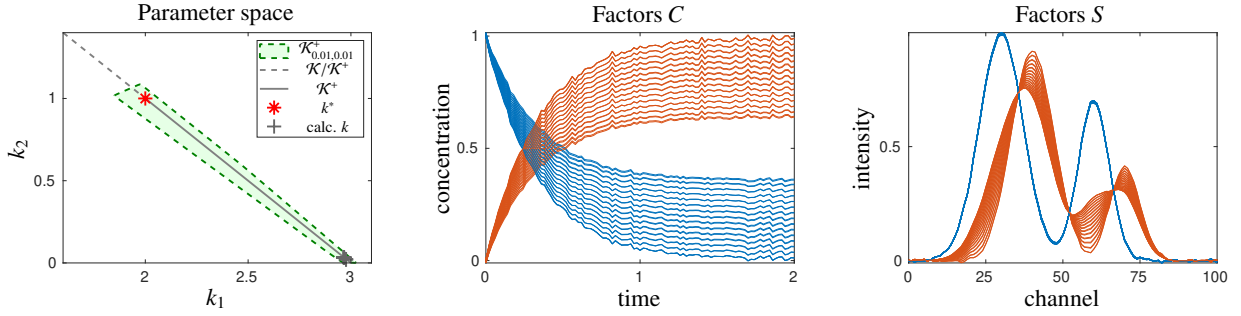
5

Figure 3: This figure shows the pendant of Figure 2 for the case of noisy data. The left subplot shows the set $\mathcal{K}^+_{0.01,0.01}$ (green). It contains $k^*$ (red asterisk) as well as the 30 calculated parameter vectors $k$ (gray crosses). The plotted concentration profiles (center) and spectra (right) correspond to equidistantly distributed parameter vectors on the boundary of $\mathcal{K}^+_{0.01,0.01}$.

bations in $D$ and then compute

$$\varepsilon = \max_{k \in \mathcal{K}} \frac{\|C(k) - C^{\text{ode}}(k)\|_F}{\|C(k)\|_F}.$$

- If $\mathcal{K}$ is not available, then any optimized parameter vector $k$ can be used to approximate the error tolerances:

$$\varepsilon = \alpha \frac{\|C(k^*) - C^{\text{ode}}(k^*)\|_F}{\|C(k^*)\|_F}, \quad \alpha \in [1, 1.2], \tag{6}$$

$$\theta = -\beta \min_{i,j} \frac{(S(k^*))_{i,j}}{\max_l(|(S(k^*))_{l,j}|)}, \quad \beta \in [1, 1.2]. \tag{7}$$

## 4. Cube enclosure algorithm

Next we present a numerical approximation method for the set of feasible $D$-approximate parameters $\mathcal{K}^+_{\varepsilon,\theta}$. This allows us to analyse the parameter ambiguities for experimental data sets. Our idea is to enclose the set by series of cubes in an iterative process. We call this approximation process the Cube Enclosure Algorithm.

The set $\mathcal{K}^+_{\varepsilon,\theta}$ can be represented in the form of a level set $\mathcal{N} := \{k \in \mathbb{R}^q : f(k) = 0\}$ with the function

$$f(k) = \|I - T(k)^+ T(k)\|^2_F + \max\left(\frac{\|C(k) - C^{\text{ode}}(k)\|^2_F}{\|C(k)\|^2_F} - \varepsilon, 0\right) + \sum_{i=1}^n \sum_{j=1}^s \min\left(\left(\frac{(S(k))_{i,j}}{\max_l(|(S(k))_{l,j}|)}\right) + \theta, 0\right)^2. \tag{8}$$

It holds that $f(k) = 0$ if and only if $k \in \mathcal{K}^+_{\varepsilon,\theta}$ for a given $\varepsilon$ and $\theta$. Functions $f$ for the representations of $\mathcal{K}, \mathcal{K}^+$ and $\mathcal{K}_\varepsilon$ can be derived by neglecting the last summand in (8) and/or setting $\varepsilon = \theta = 0$. All representations based on (8) involve the numerical solution of an initial value problem in order to compute $C^{\text{ode}}(k)$ which is numerically costly. A more efficient function for $\mathcal{K}$ in the context of first-order kinetic models is described in Appendix C.

In this section an algorithm is presented that generates an enclosing superset of $\mathcal{N}$ for a given $f(k)$ by a set of $q$-dimensional cubes. The idea is similar to the approximation of the boundary of subsets of the area of feasible solutions by the *triangle enclosure algorithm* for three-component systems as introduced by Golshan et. al. [18]. Here we use hypercubes instead of triangles, because of an easier algorithmic implementation for arbitrary dimensions $q \geq 2$. The only prerequisite for the execution of the algorithm is the knowledge of at least one initial $k' \in \mathcal{N}$. In the case of $\mathcal{N} = \mathcal{K}^+_{\varepsilon,\theta}$ such an element can easily be calculated with an MCR code that supports kinetic modeling.

Figure 4 illustrates the main steps of the cube enclosure algorithm:

1. An initial cube $W_0$ with an edge length $\omega$ is generated that contains the known element $k'$ in the center of $W_0$. The initial cube implicitly defines an underlying grid of cubes for the following steps. The set of enclosing cubes $\mathcal{W}$ (green) is initialized with $\{W_0\}$.

2. The set $C$ of all neighboring cubes (with at least a common edge) to the cubes in $\mathcal{W}$ (yellow) is determined. If one of these new cubes has already been tested concerning its feasibility, it can be ignored.
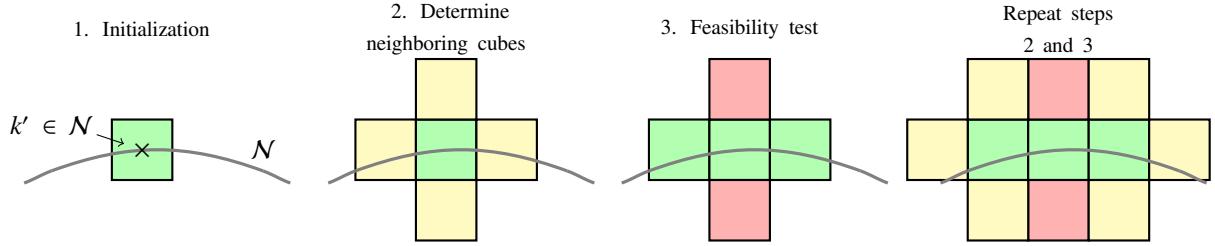
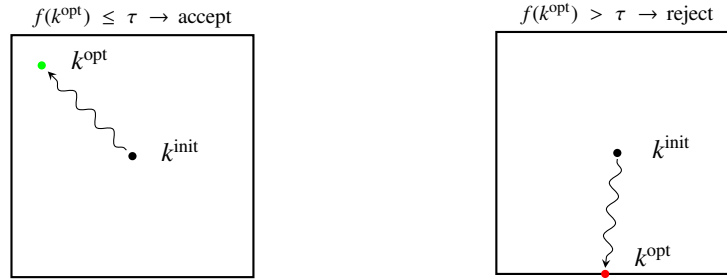Figure 4: Illustration of one iteration step of the cube enclosure algorithm.



Figure 5: 2D illustration of the two possible outcomes of the feasibility test of a cube.

3. Each cube in $C$ is tested for feasibility. This means that the function $f$ (which defines $\mathcal{N}$) is minimized within the boundaries of the currently tested cube. If the final function value $f(k^{\text{opt}})$ is close to zero, e.g., with $\tau \approx 10^{-7}$, then the cube is added to the current approximation of $\mathcal{N}$ (green). Otherwise it is rejected (red). This approximation process is illustrated in Figure 5. Here we use the *lsqnonlin* Matlab implementation of the *trust-region-reflective* algorithm [19] for the minimization. Typically, only few iterations are necessary. The obtained optimal parameter vector of the feasibility test for a cube $W$ is denoted by $k^{\text{opt}}(W)$.

4. Repeat steps 2 and 3. Cubes that have already been tested can be skipped. The algorithm stops if no further (not yet tested) neighboring cubes exist.

The optimized vectors $k^{\text{opt}}$ are known from the feasibility test for each cube whose union encloses $\mathcal{N}$. They are equally distributed with respect to the edge length $\omega$ and represent points of the set $\mathcal{N}$ with a high accuracy that can be controlled with the parameter $\tau$. The set of these vectors can be used as a representation of the set $\mathcal{N}$ and for additional post-processing steps, e.g. the application of further constraints. The detailed pseudocode of the cube enclosure algorithm is reproduced in Appendix B. The Matlab code can be found on the FACPACK homepage.

### 4.1. Refinement strategies

A given enclosure of $\mathcal{N}$ by equally-sized cubes can easily and adaptively be refined. For example, each cube can be subdivided. An obvious approach is to halve all edge lengths $\omega \to \omega/2$. Then a three-dimensional cube splits into eight sub-cubes. To obtain an improved approximation of $\mathcal{N}$ each sub-cube has to be tested for feasibility, see Step 3 and Figure 5.

This cube refinement procedure can be implemented in an adaptive way. We suggest to refine only those cubes which most likely increase the accuracy of the approximation of $\mathcal{N}$. To this end we use the following approach. Let $W$ be an approximating cube of $\mathcal{N}$. First, the algorithm has to decide whether $W$ should be subdivided or not. If $W$ is an inner cube, which means that all neighbors of $W$ are part of the approximation of $\mathcal{N}$, then $W$ remains undivided. Otherwise, the algorithm proceeds as follows: Let $W_1^N, \ldots, W_p^N$ be the neighboring cubes of $W$, that are also part of the approximation of $\mathcal{N}$. If

$$\max_{i=1,\ldots,p} f\left(0.5\left(k^{\text{opt}}(W) + k^{\text{opt}}(W_i^N)\right)\right) < \tau \tag{9}$$
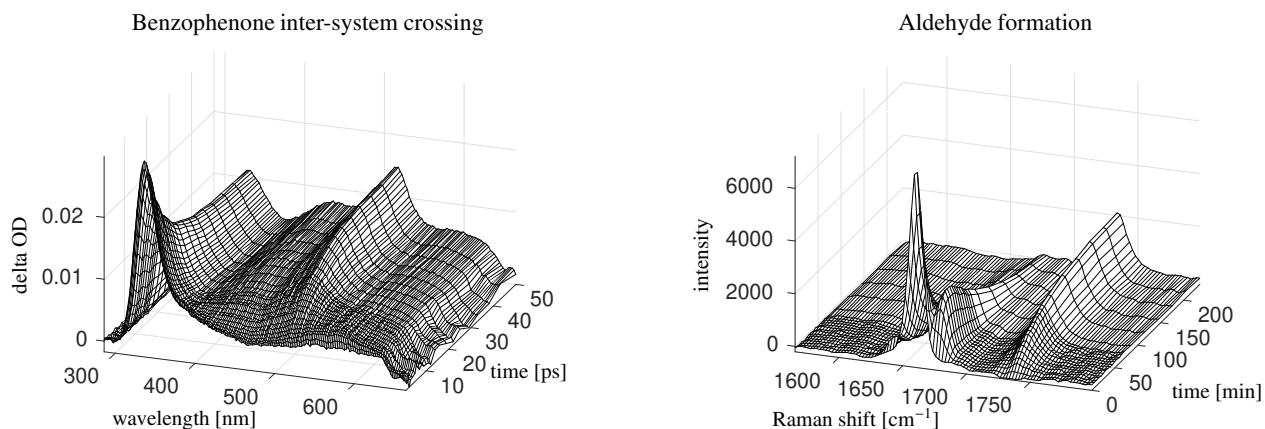
7

Figure 6: Representation of the data matrices $D$ of the data sets in Sections 5.1 and 5.2.
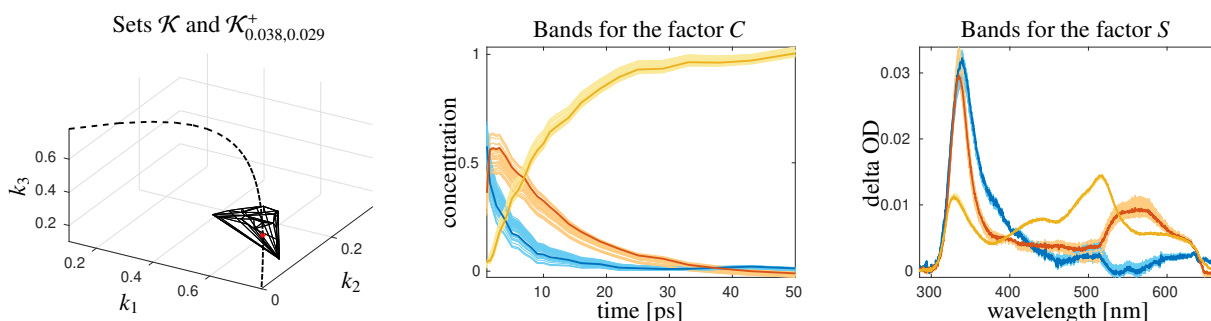


Figure 7: Results of the application of the cube enclosure algorithm for the data set from Section 5.1. The convex hull of the calculated elements of $\mathcal{K}_{0.038,0.029}^+$ by the cube enclosure algorithm is shown in the left plot in black. Additionally the parameter vector $k^*$ of the solution of the hard-modeling approach is plotted in red. In the center and right plot the bands of the factors $C = C(k)$ and $S = S(k)$ are shown for the determined elements $k \in \mathcal{K}_{0.038,0.029}^+$. The factors that correspond to $k^*$ are highlighted by bold colored lines.

is fulfilled for the error tolerance $\tau$ from step 3 of the cube enclosure algorithm, then the cube $W$ remains undivided. In other words, the mean values of the optimal parameter vector $k^{\mathrm{opt}}(W)$ and each $k^{\mathrm{opt}}(W_i^N)$ of the neighboring cubes are evaluated regarding their feasibility. If (9) holds, then it is assumed that a linear interpolation method is sufficient to determine vectors in a surrounding of $k^{\mathrm{opt}}(W)$ and a further subdivision of $W$ is not improving the local approximation quality.

### 4.2. Parallelization of the cube enclosure algorithm

Parallel computing by simultaneous execution of processes is a well-known and effective strategy in computer sciences in order to speed up computer programs. Typically, this requires that the computational problem can be divided into subproblems which can be solved independently. Here, the feasibility tests in each iteration can be parallelized since for every test the corresponding minimization is restricted to its associated cube. Also the refinement of an enclosing superset of $\mathcal{N}$ and the resulting feasibility tests can easily be parallelized.
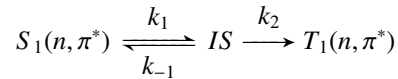
## 5. Numerical results

Next, two spectroscopic data sets are studied regarding the existence of parameter ambiguities.

### 5.1. Relaxation of the photoexcited $S_1(n, \pi^*)$ state of benzophenone (UV/Vis)

Benzophenone and benzophenone derivatives have been widely studied as they provide model compounds for the study of the photophysics of aromatic ketones. In the investigated picosecond time-range, the model generally admitted for benzophenone intersystem crossing would imply an intermediate state $IS$ in the $S_1$ to $T_1$ states relaxation

process. Characterizing these states can play a very important role in understanding the fundamental dynamics and the role they take in various applications.

Here, we analyze a data set that contains a series of transient absorption spectra of benzophenone. Transient absorption spectroscopy is a time-resolved optical spectroscopy technique by which a pump laser pulse triggers a photo-induced chemical process and the subsequent system evolution is monitored by a delayed probe pulse. The measured UV/Vis spectroscopic data set $D \in \mathbb{R}^{25 \times 1340}$ has $n = 1340$ data channels from 285 nm to 665 nm in each of the $m = 25$ measured spectra [20, 21]. The series covers a time span from 1 ps to 50 ps. The matrix $D$ is illustrated on the left in Figure 6. Relative concentration values are used in the analysis and thus the initial concentration values are set to $c_0 = (1, 0, 0)^T$. The assumed kinetic model is

$$S_1(n, \pi^*) \underset{k_{-1}}{\overset{k_1}{\rightleftharpoons}} IS \overset{k_2}{\longrightarrow} T_1(n, \pi^*)$$

with the excited $S_1(n, \pi^*)$ state, an intermediate step $IS$ and lowest excited triplet $T_1(n, \pi^*)$ state. The model is slightly more complex than the one proposed in [20, 21] and thus gives a more general approach.

First the hard-model approach [7] is used to determine optimized concentration and spectral factors $C^*$ and $S^*$ as well as the parameter vector $k^* = (0.6115 \text{ ps}^{-1}, 0.1421 \text{ ps}^{-1}, 0.1328 \text{ ps}^{-1})^T$. The factors are presented in the center and right plot of Figure 7 by bold colored lines. The obtained solution is optimal in a sense that the equally weighted sum of the following error indicators is minimal:

$$\frac{\|D - C^*(S^*)^T\|_F}{\|D\|_F} = 0.024 \quad \frac{\|C^* - C^{\text{ode}}(k^*)\|_F}{\|C^*\|_F} = 0.037, \quad \frac{\|\min(C^*, 0)\|_F}{\|C^*\|_F} = 0.003, \quad \frac{\|\min(S^*, 0)\|_F}{\|S^*\|_F} = 0.007.$$

Then the Cube Enclosure Algorithm is applied. The initial cube is generated with an edge length of $\omega = 0.07$ such that $k^*$ is located in its center. The tolerances were chosen according to the recommended settings at the end of Section 3. These values are

$$\varepsilon = 0.038 > 0.0371 = \frac{\|C^* - C^{\text{ode}}(k^*)\|_F}{\|C^*\|_F} \quad \text{and} \quad \theta = 0.029 > 0.0283 = -\min_{i,j} \frac{(S(k^*))_{i,j}}{\max_l(|(S(k^*))_{l,j}|)} . \tag{10}$$

One (nonadaptive) refinement step is applied. Figure 8 illustrates the evolution of the set of approximating cubes and the refinement step for the approximation of $\mathcal{K}^+_{\varepsilon,\theta}$.

The results are presented in Figure 7. We recommend to evaluate the feasible (optimal) parameter vectors instead of the cubes. For illustration purposes the left plot shows the convex hull of these parameter vectors by black solid lines. The band plots in the center and on the right are more descriptive. The factors $C(k)$ and $S(k)$ are plotted by pale colors for the aforementioned optimal parameter vectors $k$ that determine the approximation of $\mathcal{K}^+_{\varepsilon,\theta}$. Additionally the factors $C^* = C(k^*)$ and $S^* = S(k^*)$ are shown by bold colored lines. Especially the concentration factors show significant differences by just allowing slightly larger error tolerances than the ones obtained for $C^*$ and $S^*$, cf. (10). The indicated bands show all feasible concentration factors within the error tolerance $\varepsilon$. In contrast, the characteristic peaks of the spectral factors are quite similar. We conclude that a qualitative analysis, namely the assignment of the spectra to chemical species, can be done reliably. A quantitative analysis, however, will be affected by uncertainties. For a first order model, like the one in this example, such non-unique solutions have to be expected. Often a reduction of the solution set is only possible by adding further constraints for the factors. These solution sets usually do not occur in higher order models. Their occurrence can have different reasons, e.g. the chosen model or the time or spectral resolution.

## 5.2. Aldehyde formation from cis-2-butene (Raman)

Hydroformylation is one of the most important homogeneously catalyzed reactions in chemical industry. In this reaction, olefins are converted into aldehydes using synthesis gas, a mixture of carbon monoxide and hydrogen, and a homogeneous catalyst system. The olefins used in this process have different chain lengths. The raw materials used for the process include ethylene, propylene and a large number of other olefins with chain lengths well over ten carbon atoms. The reaction kinetic quantities here comprise several orders of magnitude depending on the olefin used. For this reason, precise knowledge of these quantities is essential for the design of the reactors and for the
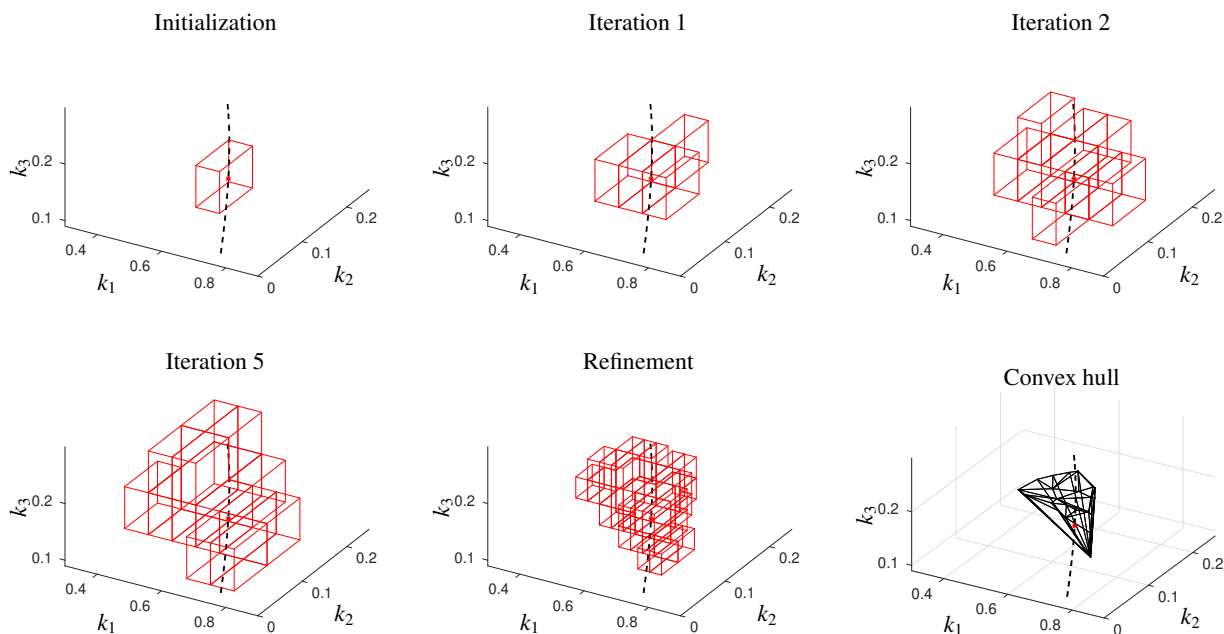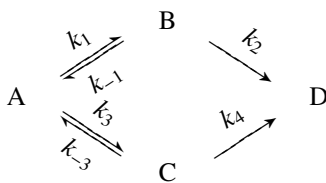
Figure 8: The steps of the cube enclosure algorithm to compute an approximation of $\mathcal{K}^+_{0.038,0.029}$ are illustrated for the data set from Section 5.1. This includes the initialization, multiple iteration steps, one refinement step and the convex hull of the calculated elements of $\mathcal{K}^+_{0.038,0.029}$.

planning of the technical process. In the case of liquid gases as feedstock (e.g. butenes), standardized analytical tools as gas chromatography often suffer from poor reproducibility or require complex sampling set-ups to gain accurate time-resolved data of the reaction progress. Another approach presented here is the usage of spectroscopic tools as infrared or Raman to detect the reaction progress online using an immersion probe inside the set-up without any offline sampling necessary.

A Raman spectroscopic data set $D \in \mathbb{R}^{202 \times 782}$ with $n = 782$ data channels from $1566\,\text{cm}^{-1}$ to $1800\,\text{cm}^{-1}$ in each of the $m = 202$ measured spectra is analyzed. The spectra cover a time span from 0 min to 238 min. The matrix $D$ is illustrated on the right in Figure 6. Relative concentration values are used in the analysis and thus the initial concentrations are set to $c_0 = (1, 0, 0, 0)^T$. We assume the kinetic model



with *cis*-2-butene (A), 1-butene (B), *trans*-2-butene (C) and a mixture of pentanal and 2-methyl-butanal (D). The model is simplified with regard to the two product species. Both, pentanal and 2-methyl-butanal, contribute in the spectral range from $1730\,\text{cm}^{-1}$ to $1750\,\text{cm}^{-1}$ and are modeled as one component. A separation of the two peaks was not possible.

Again the hard-model approach is used to determine optimized concentration and spectral factors $C^*$ and $S^*$ as well as the parameter vector

$$k^* = (k_1^*, k_{-1}^*, k_2^*, k_3^*, k_{-3}^*, k_4^*)^T = (0.033\,\text{s}^{-1}, 0.1604\,\text{s}^{-1}, 0.2008\,\text{s}^{-1}, 0.0362\,\text{s}^{-1}, 0.0309\,\text{s}^{-1}, 0.0110\,\text{s}^{-1})^T.$$

The factors are presented in Figure 9 by bold colored lines. The obtained solution is optimal in a sense that the equally
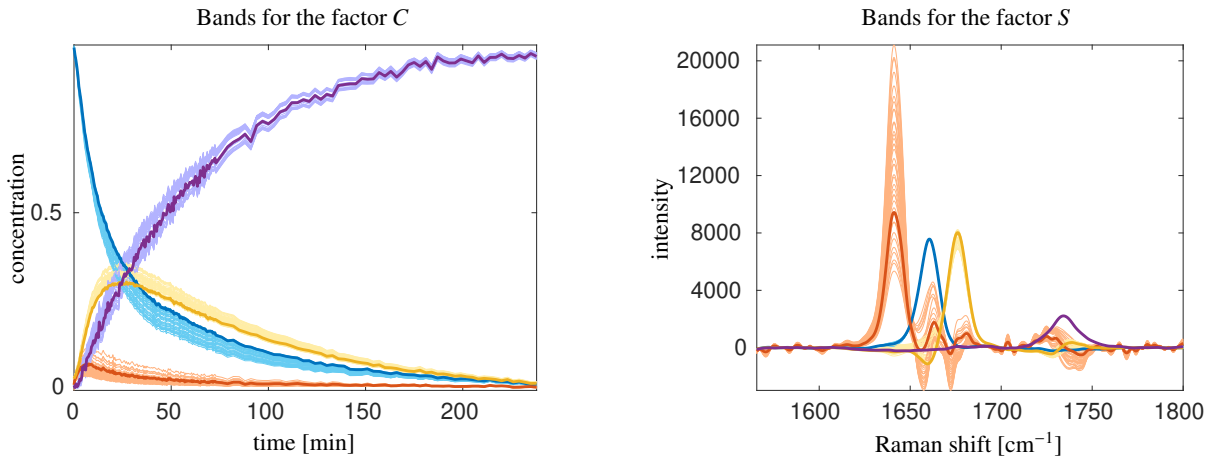
Figure 9: Results of the application of the cube enclosure algorithm for the data set from Section 5.2. The calculated elements of $\mathcal{K}^+_{0.018,0.141}$ are represented as band plots for $C$ (left) and $S$ (right).

weighted sum of the following error indicators is minimal:

$$\frac{\|D - C^*(S^*)^T\|_F}{\|D\|_F} = 0.023 \quad \frac{\|C^* - C^{\text{ode}}(k^*)\|_F}{\|C^*\|_F} = 0.0173, \quad \frac{\|\min(C^*,0)\|_F}{\|C^*\|_F} = 0.0002, \quad \frac{\|\min(S^*,0)\|_F}{\|S^*\|_F} = 0.097.$$

In particular, $S^*$ has relatively large negative matrix elements. This results from the spectrum of the species B, that has a rather negative contribution around 1660 cm$^{-1}$ in the highly overlapping range of the data set.

Starting from this solution the cube enclosure algorithm is applied. The initial cube is generated with a side length of $\omega = 0.06$ such that $k^*$ is located in its center. The tolerances were chosen according to the recommendation at the end of Section 3. The values are

$$\varepsilon = 0.018 > 0.017298 = \frac{\|C^* - C^{\text{ode}}(k^*)\|_F}{\|C^*\|_F} \quad \text{and} \quad \theta = 0.141 > 0.1405 = -\min_{i,j} \frac{(S(k^*))_{i,j}}{\max_l(|(S(k^*))_{l,j}|)}.$$

Only one (nonadaptive) refinement step is done.

The results are presented in Figure 9. Since the set $\mathcal{K}^+_{\varepsilon,\theta}$ has 6 dimensions, a graphical representation is not possbile. The factors $C(k)$ and $S(k)$ are plotted by pale colors for all elements of the set $\mathcal{K}^+_{\varepsilon,\theta}$ that were determined by the cube enclosure algorithm. Additionally the factors $C^* = C(k^*)$ and $S^* = S(k^*)$ are shown by bold colored lines. In contrast to the previous example the bands for the factor $C$ show only minor deviations. Also the spectra of the species A, C and D show only a minor ambiguity. The most significant differences can be found in the band for the factor $S$ of the species B (red). This is not surprising because of the low intensity of this peak in the data set, cf. Figure 6. Thus the species B is particularly affected by perturbations, for example baseline errors.

## 6. Summary and conclusion

MCR methods in combination with a kinetic modeling are sensitive with respect to various perturbations and/or noise in the spectral data. This might result in unreliable reaction rate constant and false conclusions regarding their uniqueness. In order to investigate the presence of solution ambiguities of kinetic model parameters, the set of feasible $D$-consistent parameters as well as a suitable approximation method were introduced. The set $\mathcal{K}^+_{\varepsilon,\theta}$ in the abstract space of reaction rate constants can easily be presented in terms of bands of feasible concentration profiles and bands of feasible spectra. The applicability is demonstrated for two experimental spectroscopic data sets.

The suggested techniques together with the provided Matlab code enable the readers to improve the reliability of their MCR results. The procedure is of a general nature because the computation and evaluation of $\mathcal{K}^+_{\varepsilon,\theta}$ can be applied as a post-processing step to any MCR method that supports kinetic modeling.

In perspective, the proposed algorithm can be easily adapted for the use of further constraints. For example this makes it possible to extract those model parameters that can be assigned to the smoothest or most unimodal concentration profiles or spectra.

## References

[1] R. Tauler, A. Smilde, and B. Kowalski. Selectivity, local rank, three-way data analysis and ambiguity in multivariate curve resolution. *J. Chemom.*, 9(1):31–58, 1995.

[2] M. Vosough, C. Mason, R. Tauler, M. Jalali-Heravi, and M. Maeder. On rotational ambiguity in model-free analyses of multivariate data. *J. Chemom.*, 20(6-7):302–310, 2006.

[3] H. Abdollahi and R. Tauler. Uniqueness and rotation ambiguities in Multivariate Curve Resolution methods. *Chemom. Intell. Lab. Syst.*, 108(2):100–111, 2011.

[4] H. Gampp, M. Maeder, C.J. Meyer, and A.D. Zuberbuehler. Quantification of a known component in an unknown mixture. *Anal. Chim. Acta*, 193:287–293, 1987.

[5] R. Tauler, A. Izquierdo-Ridorsa, and E. Casassas. Simultaneous analysis of several spectroscopic titrations with self-modelling curve resolution. *Chemom. Intell. Lab. Syst.*, 18(3):293–300, 1993.

[6] S. Vajda and H. Rabitz. Identifiability and distinguishability of first-order reaction systems. *J. Phys. Chem.*, 92(3):701–707, 1988.

[7] H. Schröder, M. Sawall, C. Kubis, D. Selent, D. Hess, R. Franke, A. Börner, and K. Neymeyr. On the ambiguity of the reaction rate constants in multivariate curve resolution for reversible first-order reaction systems. *Anal. Chim. Acta*, 927:21–34, 2016.

[8] H. Schröder, C. Ruckebusch, O. Devos, R. Métivier, M. Sawall, D. Meinhardt, and K. Neymeyr. Analysis of the ambiguity in the determination of quantum yields from spectral data on a photoinduced isomerization. *Chemom. Intell. Lab. Syst.*, 189:88–95, 2019.

[9] O. Devos, H. Schröder, M. Sliwa, J.P. Placial, K. Neymeyr, R. Métivier, and C. Ruckebusch. Photochemical multivariate curve resolution models for the investigation of photochromic systems under continuous irradiation. *Anal. Chim. Acta*, 1053:32–42, 2019.

[10] A. de Juan, M. Maeder, M. Martínez, and Tauler R. Combining hard- and soft-modelling to solve kinetic problems. *Chemom. Intell. Lab. Syst.*, 54:123–141, 2000.

[11] M. Maeder and Y.M. Neuhold. *Practical data analysis in chemistry*. Elsevier, Amsterdam, 2007.

[12] N.W. Alcock, D.J. Benton, and P. Moore. Kinetics of series first-order reactions. *Trans. Faraday Soc.*, 66:2210–2213, 1970.

[13] A.K. Smilde, H.C.J. Hoefsloot, H.A.L. Kiers, S. Bijlsma, and H.F.M. Boelens. Sufficient conditions for unique solutions within a certain class of curve resolution models. *J. Chemom.*, 15(4):405–411, 2001.

[14] S. Vajda and H. Rabitz. Identifiability and distinguishability of general reaction systems. *J. Phys. Chem.*, 98(20):5265–5271, 1994.

[15] J. Jaumot, A. de Juan, and R. Tauler. MCR-ALS GUI 2.0: New features and applications. *Chemom. Intell. Lab. Syst.*, 140(Supplement C):1 – 12, 2015.

[16] M Maeder and P King. Reactlab, 2009.

[17] J. Jaumot, P.J. Gemperline, and A. Stang. Non-negativity constraints for elimination of multiple solutions in fitting of multivariate kinetic models to spectroscopic data. *J. Chemom.*, 19(2):97–106, 2005.

[18] A. Golshan, H. Abdollahi, and M. Maeder. Resolution of rotational ambiguity for three-component systems. *Anal. Chem.*, 83(3):836–841, 2011.

[19] T.F. Coleman and Y. Li. An interior trust region approach for nonlinear minimization subject to bounds. *SIAM J. Optim.*, 6(2):418–445, 1996.

[20] S. Aloïse, C. Ruckebusch, L. Blanchet, J. Réhault, G. Buntinx, and J.-P. Huvenne. The benzophenone $S_1(n,\pi^*) \rightarrow T_1(n, \pi^*)$ states intersystem crossing reinvestigated by ultrafast absorption spectroscopy and multivariate curve resolution. *J. Phys. Chem. A*, 112(2):224–231, 2008.

[21] C. Ruckebusch, S. Aloise, L. Blanchet, J.-P. Huvenne, and G. Buntinx. Reliable multivariate curve resolution of femtosecond transient absorption spectra. *Chemom. Intell. Lab. Syst.*, 91(1):17–27, 2008.

## Appendix A. Mathematical analysis of the sets $\mathcal{K}$ and $\mathcal{K}^+$ as considered in Section 3

The formula for the determination of $\mathcal{K}$ and $\mathcal{K}^+$ in Section 3 are explained in this section based on [7]. The kinetic model in (4) is equivalent to the initial value problem

$$\dot{c}(t) = \underbrace{\begin{pmatrix} -k_1 & k_2 \\ k_1 & -k_2 \end{pmatrix}}_{M(k)} c(t), \quad c(0) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}. \tag{A.1}$$

For arbitrary $k$ the eigenvalues of $M(k)$ are 0 and $k_1 + k_2$. For a known $k^* \in \mathcal{K}$ with $\kappa = k_1^* + k_2^*$ the eigenvalues of $M(k^*)$ are denoted by $\lambda_1 = 0$ and $\lambda_2 = \kappa$. This results in an equivalent definition of $\mathcal{K}$ (see (2)):

$$
\begin{aligned}
\mathcal{K} &= \left\{ k \in \mathbb{R}^2 : M(k) \text{ has the eigenvalues } \lambda_1 \text{ and } \lambda_2 \right\} \\
&= \left\{ k \in \mathbb{R}^2 : \lambda_1 = 0 \text{ and } \lambda_2 = k_1 + k_2 \right\} \\
&= \left\{ k \in \mathbb{R}^2 : \lambda_2 = k_1 + k_2 \right\} \\
&= \left\{ \begin{pmatrix} k_1 \\ \kappa - k_1 \end{pmatrix} : 0 < k_1 \leq \kappa \right\}
\end{aligned}
\tag{A.2}
$$

For the kinetic model/initial value problem (A.1) the set $\mathcal{K}^+$ is a subset of $\mathcal{K}$. It is determined by the lower bound

$$
k_1^* \max_{i=1,\dots,n} \frac{S_{1i}^* - S_{2i}^*}{S_{1i}^*} = 1.9984
$$

for $k_1$ in (A.2). Here the reaction rate vector $k^*$ and the corresponding spectral factor $S^*$ are used, which are also employed in the data generation in Section 3. In fact every valid result of an MCR method can be taken.

## Appendix B. Pseudocode of the cube enclosure algorithm

**Input:** initial element $k^* \in \mathcal{N}$, edge length $\omega \in \mathbb{R}$, error tolerance $\tau \geq 0$
**Output:** set of enclosing cubes $\mathcal{W}$ of $\mathcal{N}$

Construct initial cube $W_0$ with center $k^*$ and edge length $\omega$
$\mathcal{W} = \{W_0\}$
$C = \text{neighbour}(\mathcal{W})$
$\mathcal{A} = \mathcal{W} \cup C$
**while** $C \neq \emptyset$ **do**
    $C^+ = \emptyset$
    **for all** $W \in C$ **do**
        **if** feasibility_test$(W, \tau)$ **then**
            $\mathcal{W} = \mathcal{W} \cup \{W\}$
            $C^+ = C^+ \cup \{W\}$
        **end if**
    **end for**
    $C = \text{neighbour}(C^+) \setminus \mathcal{A}$
    $\mathcal{A} = \mathcal{A} \cup C$
**end while**

## Appendix C. Level set representations of $\mathcal{K}$ and $\mathcal{K}^+$

For first-order kinetic models a function $f(k)$ can be used to define the level set $\mathcal{N}$ for $\mathcal{K}$ in Section 4 that results in a more efficient computation. Let

$$
\dot{c}(t) = M(k)c(t), \quad c(0) = c_0
$$

be the initial value problem that corresponds to a given kinetic model of first order with a coefficient matrix $M(k) \in \mathbb{R}^{s \times s}$ and initial concentrations $c_0$. It is known from [7] that $\mathcal{K}$ can be represented with the help of the eigenvalues $M(k^*)$ for a known $D$-consistent parameter vector $k^* \in \mathcal{K}$. Such a vector can be computed by the hard-model approach described in [7] or any other MCR methods that supports kinetic modeling (e.g. MCR-ALS or ReactLab). Now let

13

$\sigma(k)$ be a function that maps $k$ to a vector that contains the sorted (in ascending order) eigenvalues of $M(k)$. Then the alternative representation reads

$$\mathcal{N} = \mathcal{K} = \left\{ k \in \mathbb{R}^q : \underbrace{\frac{\|\sigma(k) - \sigma(k^*)\|_2}{\|\sigma(k^*)\|_2}}_{f(k)} = 0 \right\}$$

Based on the function $f(k)$ a test whether a $k$ belongs to the set $\mathcal{N}$ or not just involves solving an $(s \times s)$-eigenvalue problem instead of solving an initial value problem (for $C^{\text{ode}}(k)$ in (8)). This is significantly faster for all application cases with $m, n \gg s$.