

On Polynomial Gaussian Least-Squares Fits and Interpretation of the Resulting Fit-Parameters

Uwe Hohm

Institut für Physikalische und Theoretische Chemie der Technischen Universität Braunschweig, FRG

Z. Naturforsch. **48a**, 878–882 (1993); received May 8, 1993

Fitting of experimental data with a polynomial determined by Gaussian least-squares analysis is a common and well-established praxis. However, there seems to be a lack in consistency if the resulting fit parameters are to be interpreted as physically meaningful quantities. Therefore, a detailed investigation of the fitting procedure itself as well as of the measuring process, including possible improvements of the resulting fit-parameters are presented in this paper.

Key words: Least-squares fit, evaluation of measurements.

1. Introduction

To describe measured data-sets $\{x, f(x)\}$ a simple polynomial calculated from Gaussian least-squares analysis is often used. However, a problem which may arise is the interpretation of the resulting fit-parameters. If a polynomial of finite degree m is fitted to data which can be represented in terms of physically meaningful parameters f_k only by a power series $f(x)$ of infinite degree, erroneous and spurious interpretations of the resulting fit-parameters may be given on account of the applied Gaussian least-squares fit. This problem is known for quite a long time [1] and discussed in the framework of analyzing pVT -data of CO_2 , but obviously no further treatment was given up to now.

Let us first assume, that the “correct” and physically meaningful power-series is given by

$$f(x) = \sum_{k=0}^{\infty} f_k x^k. \quad (1)$$

Examples may be the thermodynamic equation of state ($x \equiv \rho = \text{molar density}$, $f(x) \equiv Z = \rho_0/\rho = (p/RT)/\rho$, $\rho_0 = \text{molar density of the perfect gas}$, $f_k \equiv B_k = \text{thermodynamic virial coefficient, describing intermolecular interactions between } k+1 \text{ particles}$), or the evaluation of the dipole polarizability in terms of Cauchy moments ($x \equiv \omega^2 = \text{squared frequency}$, $f(x) \equiv \alpha(\omega) = \text{linear dipole polarizability}$, and $f_k \equiv S(-2k-2) = \text{dipole oscillator sums, which are related to excita-}$

tions of all degrees of freedom), respectively. Usually, the N measured data points $\{x_i, f(x_i)\}$ are fitted by the following polynomial of finite degree $m \leq N-1$:

$$f(x) \approx \tilde{f}(x) = \sum_{k=0}^m a_k^{(m)} x^k, \quad (2)$$

where the $a_k^{(m)}$ now are the fit parameters of the applied Gaussian least-squares fit. It is important to note that in general $a_k^{(m)} \neq f_k$. But instead of this trivial statement, it is common praxis to interpret $a_k^{(m)}$ as the corresponding physically meaningful parameters f_k . Of course, in many cases this may be done without significant loss of accuracy. However, sometimes controverse and erroneous interpretations of the $a_k^{(m)}$ are given. Based on these statements it will be shown how the measurements themselves may be performed (free choice of x_i is provided) to get $a_k^{(m)} \approx f_k$. Moreover, if any relationship between the physically meaningful parameters f_k is known, an iterative improvement of the $a_k^{(m)}$ can be applied to come close to the desired f_k , as will be shown below.

2. Formulation of the Procedure

Assuming that the measured N data-points $\{x_i, f(x_i)\}$ are correctly described by (1), the problem is to minimize the squared deviation

$$Q = \sum_{j=1}^N [\tilde{f}(x_j) - f(x_j)]^2 \quad (3)$$

with respect to the “unknown” $a_k^{(m)}$ of (2). Well-known and straightforward analysis leads to the matrix equa-

Reprint requests to Dr. U. Hohm, Institut für Physikalische und Theoretische Chemie, Technische Universität Braunschweig, Hans-Sommer-Straße 10, 38106-Braunschweig, FRG.

0932-0784 / 93 / 0800-0878 \$ 01.30/0. – Please order a reprint rather than making your own copy.



Dieses Werk wurde im Jahr 2013 vom Verlag Zeitschrift für Naturforschung in Zusammenarbeit mit der Max-Planck-Gesellschaft zur Förderung der Wissenschaften e.V. digitalisiert und unter folgender Lizenz veröffentlicht: Creative Commons Namensnennung-Keine Bearbeitung 3.0 Deutschland Lizenz.

Zum 01.01.2015 ist eine Anpassung der Lizenzbedingungen (Entfall der Creative Commons Lizenzbedingung „Keine Bearbeitung“) beabsichtigt, um eine Nachnutzung auch im Rahmen zukünftiger wissenschaftlicher Nutzungsformen zu ermöglichen.

This work has been digitalized and published in 2013 by Verlag Zeitschrift für Naturforschung in cooperation with the Max Planck Society for the Advancement of Science under a Creative Commons Attribution-NoDerivs 3.0 Germany License.

On 01.01.2015 it is planned to change the License Conditions (the removal of the Creative Commons License condition “no derivative works”). This is to allow reuse in the area of future scientific usage.

tion (set of so-called "normal equations")

$$XA = S, \quad (4)$$

where the elements of the $(1+m) \times (1+m)$ matrix X are given by

$$x_{ik} = x_{ki} = \sum_{j=1}^N x_j^{i+k-2}, \quad i, k = 1, \dots, 1+m. \quad (5)$$

The elements of the $1 \times (1+m)$ column-vector S are given by

$$s_i = \sum_{j=1}^N x_j^{i-1} \sum_{k=0}^{\infty} f_k x_j^k = \sum_{j=1}^N \sum_{k=0}^{\infty} f_k x_j^{k+i-1}, \quad (6)$$

and A is the $1 \times (1+m)$ solution vector with elements $a_{l-1}^{(m)}$, $l = 1, \dots, 1+m$. The solution of this problem can be formulated as the matrix equation

$$A = X^{-1} S, \quad (7)$$

which gives a relation between the fit-parameters $a_{l-1}^{(m)}$ and the desired parameters f_k of the polynomial (1):

$$a_{l-1}^{(m)} = \sum_{k=0}^{\infty} f_k \sum_{i=1}^{1+m} \sum_{j=1}^N \tilde{x}_{li} x_j^{k+i-1} = \sum_{k=0}^{\infty} f_k \varphi_{lk}^{(m)}, \quad (8)$$

\tilde{x}_{li} being the elements of the inverse matrix X^{-1} . Equation (8) separates into

$$a_{l-1}^{(m)} = f_{l-1} + \sum_{k=l}^{\infty} f_k \varphi_{lk}^{(m)}, \quad (9)$$

where the $\varphi_{lk}^{(m)}$ can be identified to be weighting coefficients for the exact (and physically meaningful) parameters f_k in the case of fitting a polynomial of finite degree m to the N measured data-points. Therefore, it is easy to see that the $a_{l-1}^{(m)}$ are only estimates of the true parameters f_{l-1} . As can be seen from standard textbooks on numerical mathematics [2], the error in $a_{l-1}^{(m)}$ for a Gaussian least-squares fit is given by

$$\delta a_{l-1}^{(m)} = [\tilde{x}_{li} Q / (N - m - 1)]^{1/2}. \quad (10)$$

If now $\delta a_{l-1}^{(m)} < \sum_{k=l}^{\infty} f_k \varphi_{lk}^{(m)}$, it is advantageous to rewrite (9) as

$$f_{l-1} = a_{l-1}^{(m)} - \sum_{k=l}^M f_k \varphi_{lk}^{(m)}, \quad (11)$$

where $M < \infty$ has to be chosen for practical purposes. Remember that f_k is determined by the physics of the system only, and that $\varphi_{lk}^{(m)}$ depends only on the choice of x_k . Minimization of $|\sum f_k \varphi_{lk}^{(m)}|$ will lead to the desired result $a_{l-1}^{(m)} \approx f_{l-1}$. Of course, this seems to be a formidable task, but the problem may be treated as follows.

2.1. General Formalism

If nothing is known about the f_k , one way to minimize the sum term in (11) may be given by the approach

$$\sum_{k=l}^M |\varphi_{lk}^{(m)}| \rightarrow \text{Minimum, or equivalently}$$

$$P \equiv 1 / \sum_{k=l}^M |\varphi_{lk}^{(m)}| \rightarrow \text{Maximum}. \quad (12)$$

To deal with this problem in a systematic manner, the measuring interval $t_{\text{Min}} \leq t \leq t_{\text{Max}}$ is first transformed via

$$x_k \equiv (t_k - t_{\text{Min}}) / (t_{\text{Max}} - t_{\text{Min}}) \quad (13)$$

into the interval $0 \leq x \leq 1$. The N -abscissa values x_k are then calculated via

$$x_k = \left\{ 1 - \exp \left[- \left(\frac{k-1}{N-1} \right)^\lambda \right] \right\} / [1 - \exp(-1)], \quad (14)$$

$k = 1, \dots, N.$

As is shown in Fig.1, this one-parametric function allows the realization of the following distributions of the x_k :

$$\lambda = \begin{cases} < 1.3: & \text{accumulation at large } x, \\ \approx 1.3: & \text{nearly equal spaced } x, \\ > 1.3: & \text{accumulation at low } x. \end{cases}$$

According to these distributions, $P(N, \lambda)$ is calculated by (12) with $M = 10$, as a function of N and λ for $m = 3$.

The resulting surfaces $P(N, \lambda)$ are shown in Figs. 2a–d, concerning the fit-parameters $a_0^{(3)}$, $a_1^{(3)}$, $a_2^{(3)}$, and $a_3^{(3)}$, respectively. In Figs. 2b–d, an increase of P with increasing λ is shown. This indicates a diminution of

$\sum_{k=l}^M |\varphi_{lk}^{(m)}|$ in the case of accumulation of x_k in the vicinity of $x = 0$. This is also the case for $a_0^{(m)}$, although

this cannot be seen in Fig. 2a on account of the special scaling. But additionally, in this case an accumulation of x_k at $x = 1$ (low λ) causes an increase of P . In all cases, the slight dependence of P on N is probably due to the special distribution chosen for calculation of the x_k according to (14). Although not shown, the same statements can be given for lower and higher values of m . It is interesting to note that, for given N and λ , P increases with decreasing k of the fit-parameters $a_k^{(m)}$. However, this special feature has already been found empirically [1], since the advice has been given to use high polynomial degrees m in order to get reliable $a_k^{(m)}$.

Therefore, one should measure (if possible) as many data points as possible at low values of x in the mea-

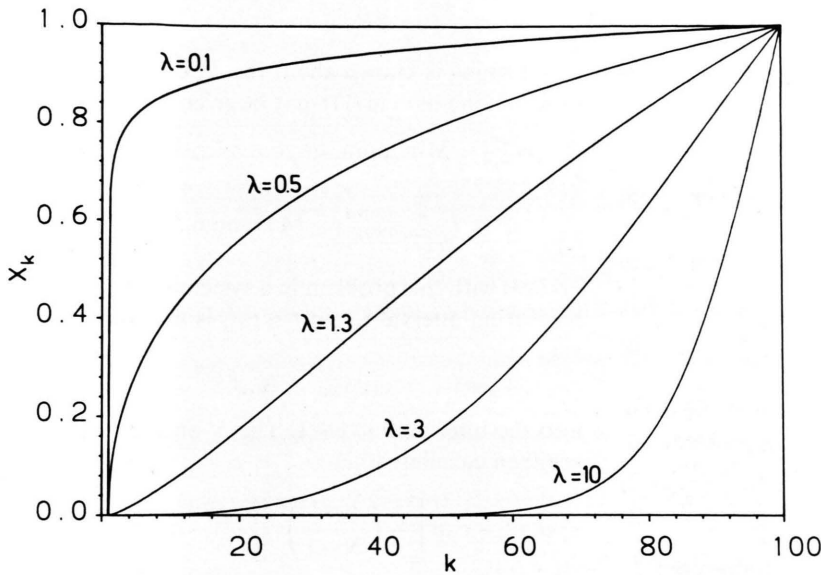


Fig. 1. Distribution of abscissa-values x_k according to (14). For the sake of clearness full curves are shown, although they are only defined for integer values of k .

Function	λ	$a_0^{(3)}$	$a_1^{(3)}$	$10a_2^{(3)}$	$10a_3^{(3)}$
exp(x)	0.1	1.0000000 (48)	1.0613 (10)	3.148 (22)	3.422 (11)
	0.3	0.999999 (90)	1.0382 (22)	3.653 (53)	3.147 (31)
	1.0	0.99967 (52)	1.0163 (46)	4.21 (11)	2.805 (72)
	3.0	0.99975 (27)	1.0133 (37)	4.289 (98)	2.762 (67)
	10.0	0.999922 (85)	1.00285 (30)	4.669 (10)	2.485 (74)
	20.0	0.999999950 (51)	1.0000791 (15)	4.91772 (12)	2.26430 (11)
		\downarrow $f_0 = 1.0$	\downarrow $f_1 = 1.0$	\downarrow $f_2 = 0.5$	\downarrow $f_3 = 0.1\bar{6}$
exp(-x)	0.1	1.0000000 (11)	-0.98012 (23)	4.3192 (48)	-0.8392 (26)
	0.3	1.000000 (24)	-0.98594 (59)	4.447 (14)	-0.9087 (83)
	1.0	0.99988 (10)	-0.9930 (17)	4.631 (26)	-1.022 (26)
	3.0	0.99990 (10)	-0.9942 (14)	4.660 (37)	-1.039 (25)
	10.0	0.9999959 (45)	-0.99846 (16)	4.8153 (54)	-1.1518 (39)
	20.0	0.999999968 (32)	-0.99995018 (96)	4.948002 (76)	-1.269706 (67)
		\downarrow $f_0 = 1.0$	\downarrow $f_1 = -1.0$	\downarrow $f_2 = 0.5$	\downarrow $f_3 = -0.1\bar{6}$

Table 1. Fit parameters $a_0^{(3)}$, $a_1^{(3)}$, $a_2^{(3)}$, $a_3^{(3)}$ of Gaussian least-squares fit of exp(x) and exp(-x) with $N=10$ data-points x_k calculated according to (14).

suming interval in order to get as close as possible to the desired approximation $a_k^{(m)} \approx f_k$. This may be demonstrated using the following "polynomials" of infinite degree with known coefficients f_k :

$$\exp(x) \equiv \lim_{M \rightarrow \infty} \sum_{k=0}^M f_k x^k, \quad f_k \equiv \frac{1}{k!}, \quad (15)$$

$$\exp(-x) \equiv \lim_{M \rightarrow \infty} \sum_{k=0}^M f_k x^k, \quad f_k \equiv \frac{(-1)^k}{k!}. \quad (16)$$

In Table 1 the resulting fit parameters $a_k^{(3)}$ are shown for $N=10$ data points, calculated via (14) in the λ -

range between 0.1 and 20. The general statements about the distribution of the x_k , visualized in Figs. 2a–d. are fully confirmed by these concrete examples. Of course, the distributions with $\lambda=0.1$ and $\lambda=20$ are pathological and can never be realized in true experiments, but they clearly confirm the trend that accumulation at low x in general improves the desired relationship $a_k^{(m)} \approx f_k$. This behaviour is reflected by the approach of the fit-parameters $a_k^{(m)}$ given in Table 1 towards the exact polynomial coefficients f_k . This behaviour seems to be very reasonable, since a Taylor-series expansion of $f(x)$ around $x=0$ also

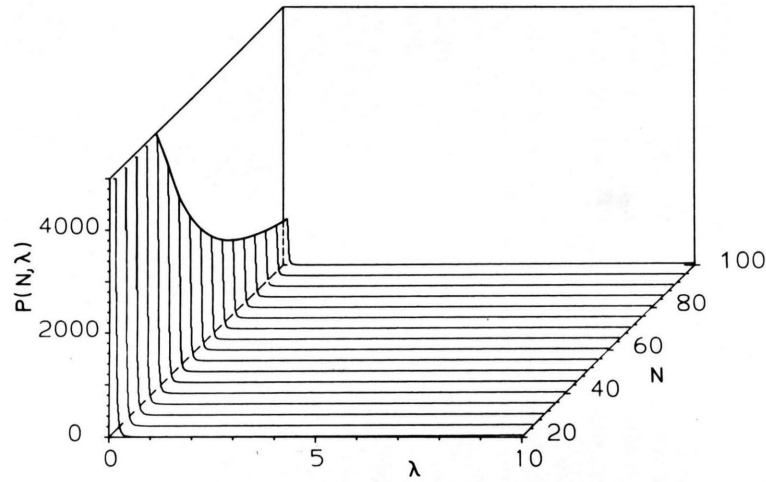


Fig. 2a. $P(N, \lambda)$ -surface of $a_0^{(3)}$, calculated according to (12).

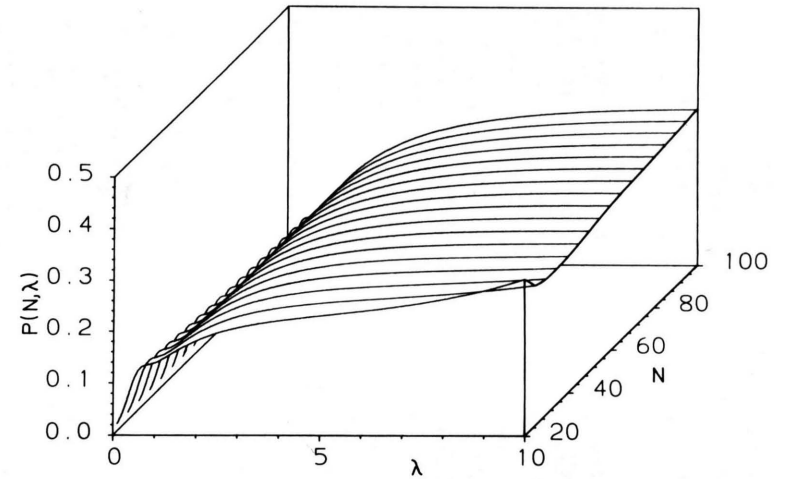


Fig. 2b. $P(N, \lambda)$ -surface of $a_1^{(3)}$, calculated according to (12).

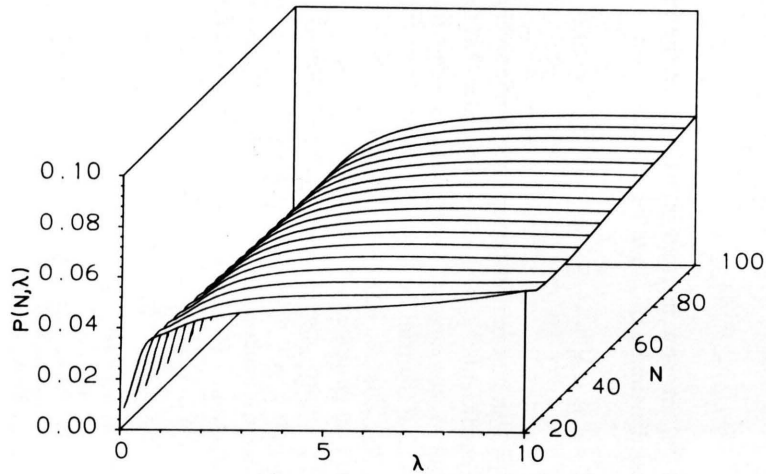


Fig. 2c. $P(N, \lambda)$ -surface of $a_2^{(3)}$, calculated according to (12).

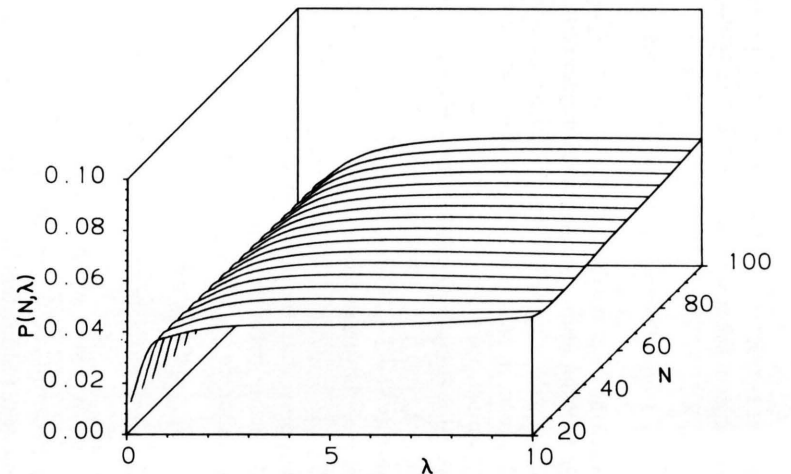


Fig. 2d. $P(N, \lambda)$ -surface of $a_3^{(3)}$, calculated according to (12).

gives the correct f_k . This approach can also be observed, e.g., in the case of $\cos(x)$ or $\sin(x)$, respectively.

2.2. The Case of Known Relationships between the f_k

After the basic mathematical formalism, which tells us something about the desirable distribution of the abscissa-values x_k , is known, a further application to real measurements is developed now. In some cases, the f_k must obey certain relationships on account of physical reasons. This can be, for instance, an inequality like $f_k < \omega f_{k+1}$, which is the case for the Cauchy-moments $S(-2-2k)$ describing the frequency dependence of the mean linear dipole polarizability $\alpha(\omega)$ [3, 4].

Suppose that the function F describes the exact relationship between f_k and f_{k+1} , and that F can be approximated by a trial-function \tilde{F}_k :

$$f_{k+1} = F(f_k) \approx \tilde{F}_k. \quad (17)$$

The zeroth estimate of f_k is given by

$$^{(0)}f_k = a_k^{(m)}, \quad k=0, \dots, m. \quad (18)$$

A superscript on f_k now always denotes an estimate of the true parameters f_k . Using the trial-function \tilde{F} , higher coefficients f_k can be approximated according to

$$1. \quad ^{(0)}f_{k+1} = \tilde{F}_k(^{(0)}f_k), \quad k=m, \dots, M-1. \quad (19)$$

All $^{(0)}f_k$ are now inserted into (11), and the first improved estimate of f_k can be calculated via

$$2. \quad ^{(1)}f_k = a_k^{(m)} - \sum_{j=k+1}^M \varphi_{k+1,j}^{(m)} ^{(0)}f_j, \quad k=0, \dots, m. \quad (20)$$

Steps 1 and 2 can now be repeated, until no change in the coefficients $^{(k)}f_k$ is observed, provided that the algorithm is convergent. Divergence is sometimes observed, if the trial function (17) yields the wrong sign of $^{(k)}f_k$. This is, however, the worst case where the iteration described above does not work.

The improvement of the $a_k^{(m)}$ by iteration with a trial function \tilde{F} is not an academic example, as has been shown with success in the case of the dispersion of the polarizability [4]. In this special case a suitable trial-

function has been

$$\tilde{F}_{l-1} = {}^{(k)}f_l = {}^{(k)}f_{m-1} \left(\frac{{}^{(k)}f_m}{{}^{(k)}f_{m-1}} \right)^{l-m+1} [1 + q(l-m)]^{c(l-m+1)}, \quad l > m-2, \quad (21)$$

which is exact for $l=m-1$ and $l=m$. In the case of increasing f_k , $c = +1$, otherwise $c = -1$ has to be chosen. With a suitable choice of q and M (e.g. $q=0.05$ and $M \geq 7$ in the case of experimentally determined Cauchy-moments), not only for the polarizability-data, but also for the exponential-functions (15) and (16) (in this case $q=0.02$ and $M=5$ has been chosen), often very reasonable final approximations $^{(2)}f_k \approx f_k$ can be obtained [4].

3. Conclusions

In this paper it has been shown, that an interpretation of the fit parameters of a Gaussian least-squares analysis is in general not a trivial task. Therefore, the following advices may be given:

a) Be careful with the interpretation of $a_k^{(m)}$, (2), if f_k , (1), has to be determined! Only if $a_k^{(m)}$ is independent of m , the right (physically meaningful) coefficients seem to be determined [1].

b) If it is possible to make a free choice of x_k in the measuring process, *provided the errors of x_k and $f(x_k)$ are unaffected by this choice* (see [4]), accumulation of x_k at $x \approx 0$ with one point at $x = 1$ (reduced quantities as in (14)), seems to be a suitable way to determine $a_k^{(m)} \approx f_k$ by Gaussian least-squares analysis.

c) If any relation is known between the parameters f_k and if the rms-errors of the applied Gaussian least-squares fit are very small (see (10)), iterative improvement of $a_k^{(m)}$ with respect to f_k is possible in certain cases [4]. But the iterative procedure has to be done very carefully, since misleading results may occur, which, however, can be detected by thorough inspection of the resulting quantities. But one should bear in mind that the resulting fit-parameters are highly correlated, as is always the case in Gaussian least-squares analysis.

[1] A. Michels, J. C. Abels, C. A. Ten Seldam, and W. De Graaff, *Physica* **26**, 381 (1960).

[2] W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. Vetterling, *Numerical Recipes in Pascal*, Cambridge University Press, Cambridge 1989.

[3] P. W. Langhoff, *J. Chem. Phys.* **57**, 2604 (1972).

[4] U. Hohm, *Mol. Phys.* **78**, 929 (1993).