# SIMULATION STUDIES OF SELF-ASSOCIATING SYSTEMS; DISCRIMINATION BETWEEN SPECIFIC AND ISODESMIC ASSOCIATIONS

Marc S. LEWIS and Gary D. KNOTT

*Laboratory of Vision Research, National Eye Institute,
and Laboratory of Statistical and Mathematical Methodology, Division of Computer Research and Technology,
National Institutes of Health, Department of Health, Education and Welfare, Bethesda, Maryland 20014, USA*

The possibilities of discriminating between definite and indefinite (isodesmic) modes of self-association are explored by fitting the different models to simulated data, using non-linear least-squares curve-fitting to determine the fitting parameters for real and impostor models. It was found that over an extensive range of values for the equilibrium constant of a non-ideal isodesmic generating model, only a non-ideal monomer–dimer–tetramer–octamer was a successful impostor model. Some criteria for rejecting inappropriate models are discussed.

## 1. Introduction

The meaningful analysis of a self-associating system basically depends upon two factors: the acquisition of precise data and the determination of the nature of the association and the numerical values of the parameters which describe it. The quality of the data will be dependent upon the skill of the investigator and the homogeneity and perfect reversibility of the associating system, and it will be ultimately limited by the available instrumentation. Since it is implicit that the determination of the type of the association and the association constants depends upon the quality of the data, computer simulation of data frees us from the inherent limitations of experimentally obtained data and allows us to more readily investigate the problems involved in establishing a unique description of an associating system.

Simulation of data has been a very commonly used technique in the study of associating systems. The literature relevant to this is so extensive that its citation would be more appropriate to a review article; hence the reader is referred to several recent monographs and reviews [1—4] for more complete discussions of the various methods which have been used both for simulation and for data analysis.

The choice of models for investigation was suggested by the past research experience of one of us. Adams and Lewis [5], using sedimentation equilibrium, des-

cribed the self-association of $\beta$-lactoglobulin A as an isodesmic association with a monomer molecular weight of 18000 daltons. An isodesmic association is a type of unlimited self-association where the addition of each successive monomer to the polymer involves an equal change in free energy [6]. Timasheff and his associates [7—9], using light scattering and sedimentation velocity, have described the association as a monomer–dimer–trimer–tetramer with a monomer molecular weight of 36000 daltons. Roark and Yphantis [10], using sedimentation equilibrium, found a similar model for this associating system. Van Holde, Rossetti, and Dyson [11], studying the self-association of cytidine found that this association could be equally well fit as a monomer–dimer–trimer, a non-ideal monomer–dimer–trimer–tetramer or as a non-ideal isodesmic association. These studies suggest that it may be very difficult to discriminate between these different modes of self-association, and this is what will be examined in the following sections of this paper.

## 2. Methods

Curve fitting is a useful analytical tool in many diverse disciplines, and in some form it is crucial to the processing and interpretation of ultracentrifuge data.

The basic notion is easily described. Given data, say various points in the plane $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$,

and a function $y = f(x)$ where $f$ involves some param-
eters, say $a$ and $b$, as in $f(x) = ax^b + 1$, we wish to cal-
culate values for the parameters $a$ and $b$ so that the
function $f$ well-predicts the observed data, that is, so
that $f(x_i) = y_i$ for $1 \leqslant i \leqslant n$. In this case, we say we
have fit the *model f* to the data. The end result is
merely the values obtained for the initially unknown
parameters.

The notion of well-prediction which is generally
used is that of minimizing the so-called sum of squares,
$S(a,b) = \Sigma_{1 \leqslant i \leqslant n} [f(x_i) - y_i]^2$ by appropriately choos-
ing values for $a$ and $b$. This is because least-squares
minimization is a well-studied method which allows
various theorems involving error and uniqueness to be
invoked. Analogous results do not exist for other meas-
ures of goodness of fit.

There are various algorithms which may be employed
to minimize a sum-of-squares value. Many of these are
reviewed by Magar [12]. One of the most robust is the
magnified diagonal form of the Marquardt—Levenberg
method [13—15] which has been used for the work
described in this paper.

When curve-fitting is considered, there are various
problems which must somehow be dealt with. First,
there must exist some way the calculations can be car-
ried out. This generally means that a digital computer
must be used, and that some programming is needed.
Curve-fitting is generally a complex iterative activity.
The modeler will wish to try various models with var-
ious initial guesses applied to various sets of data. Dif-
ferent weights may be tried, and data may need to be
scaled or otherwise transformed, as may the model.
Generally the most convenient way to view the results
is in the form of a graph of the best-fit model (the theo-
retical curve) and the plotted data points (the observed
curve). Often the computer programs employed are
neither sufficiently general nor easy to use, and the
curve-fitting process becomes a tedious exercise.

A general system called MLAB (for modeling labora-
tory) has been developed and is in use at the N.I.H.. It
is particularly suited for quick and convenient curve-
fitting [16,17], and is the program which has been em-
ployed in obtaining the results reported here. It runs
on a DECSYSTEM-10 time-sharing system and is avail-
able for public distribution. MLAB provides facilities
for data manipulation and graphics as well as curve-fit-
ting. The figures presented here were produced using
MLAB.

But, no matter how convenient the curve-fitting
process is, there remains the problem of judging and
validating the results. Is the model incorrect? Is it cor-
rect? If correct, are the parameter values obtained by
curve-fitting near the true values underlying the data?
What is the standard error associated with the computed
parameter values? These questions are often unanswer-
able, or at best answerable only by means of subjective
judgments. On the other hand, there are various facts
and heuristics which can sometimes be applied to aid
in the analysis of the results. Indeed this paper is con-
cerned with using existing heuristics and developing
new rudimentary "rules-of-thumb" which apply par-
ticularly to judging fits for self-association data ob-
tained from the ultracentrifuge.

In particular, the following "principles" are of inter-
est.

1. If the model is linear in the parameters, and the
error in the data is normally-distributed, and weights for
each data point equal to the reciprocal of the variance
for that point are used, then the fitted parameter values
are maximum-likelihood estimates and their standard
deviation estimates are correct. When any of the above
hypotheses fail, the standard deviation estimates are
merely suggestive.

2. The RMS (root-mean-square) error of a fit is a
dimensional measure of the goodness of fit. It is not
an absolute criterion, however, since equal RMS values
could be obtained for a fit that has randomly distribu-
ted deviations and a fit that shows systematic deviations.
The RMS values here take into account the number of
degrees of freedom dictated by the number of fitting
parameters used..

3. One of the best ways to judge a fit is by examin-
ing a plot of the theoretical minus observed deviations
vs. the independent variable. If these points show any
systematic deviation about zero, there may be reason
to be suspicious.

4. Some fits are very precise. Small variations in the
best-fit parameter values no longer allow the model to
be a good fit. Other fits may permit any of a wide range
of "answers". In this case, one or more parameters can
generally be changed to compensate for any small change
imposed on various other parameters. Dependency values
provide a measure of the preciseness of a fit. Parameters
with high dependency values form a set which can be
varied widely so as to maintain a good fit. One may ex-
plore the significance of dependency values by varying

the value of a parameter above and below the optimum value obtained when fitting, allowing the other parameters to vary to give the best fit under these circumstances, and then plotting the sum of squares of the deviations as a function of the value of the parameter being varied [18,19]. Since, except as noted, the dependency values obtained in this study were not large, this procedure was judged to be too costly in computer time for the value of the information it would yield.

Now we shall consider below the fitting of various models for self-associating systems to data generated by exact evaluation of the equations for an isodesmic system with a particular $k$ value. The object is to study to what extent we can reject models which in fact are incorrect in that they do not reflect physical reality. We also consider this question in the presence of error as might appear in real data.

The same equations were, of course, used for generating the data and for carrying out the analyses of the generated data. The isodesmic association is described by

$$M_w = M_1 (1 + kc_1)/(1 - kc_1),$$  (1)

and

$$c_1 = \underset{0<x<c}{\text{Root}} \{c - [x/(1 - kx)^2]\},$$  (2)

where $M_w$ is the weight-average molecular weight, $M_1$ is the molecular weight of the monomer, $c$ is the total concentration and $c_1$ is the concentration of monomer, both on a gram per liter scale, $k$ is the intrinsic association constant on that concentration scale, subject to the constraint that $kc_1 < 1$. The notation $\text{Root}_{a<x<b}(f(x))$ denotes a value, $x_0$, in the interval $a$ to $b$ which is a root of the expression $f(x)$, i.e., $x_0$ is a value such that $f(x_0) = 0$; MLAB supports this operator directly.

The monomer–dimer–hexamer (1–2–6) association is described by

$$M_w = M_1 (c_1 + 2k_{12}c_1^2 + 6k_{26}k_{12}^3 c_1^6)/c$$  (3)

and

$$c_1 = \underset{0<x<c}{\text{Root}} [c - (x + k_{12}x^2 + k_{26}k_{12}^3 x^6)],$$  (4)

where $k_{12} = c_2/c_1^2$ and $k_{26} = c_6/c_2^3$. In general, the gram per liter concentration will be denoted by $c_n$.

The monomer–dimer–trimer–hexamer (1–2–3–6)

association is described by

$$M_w = M_1 (c_1 + 2k_{12}c_1^2 + 3k_{23}k_{12}c_1^3 + 6k_{36}k_{23}^2 k_{12}^2 c_1^6)/c$$  (5)

$$c_1 = \underset{0<x<c}{\text{Root}} [c - (x + k_{12}x^2 + k_{23}k_{12}x^3$$
$$+ k_{36}k_{23}^2 k_{12}^2 x^6)],$$  (6)

where $k_{23} = c_3/c_1 c_2$ and $k_{36} = c_6/c_3^2$.

The monomer–dimer–octamer (1–2–8) association is described by

$$M_w = M_1 (c_1 + 2k_{12}c_1^2 + 8k_{28}k_{12}^4 c_1^8)/c$$  (7)

and

$$c_1 = \underset{0<x<c}{\text{Root}} [c - (x + k_{12}x^2 + k_{28}k_{12}^4 x^8)],$$  (8)

where $k_{28} = c_8/c_2^4$.

The monomer–dimer–tetramer–octamer (1–2–4–8) association is described by

$$M_w = M_1 (c_1 + 2k_{12}c_1^2 + 4k_{24}k_{12}^2 c_1^4$$
$$+ 8k_{48}k_{24}^2 k_{12}^4 c_1^8)/c$$  (9)

and

$$c_1 = \underset{0<x<c}{\text{Root}} [c - (x + k_{12}x^2 + k_{24}k_{12}^2 x^4$$
$$+ k_{48}k_{24}^2 k_{12}^4 x^8)],$$  (10)

where $k_{24} = c_4/c_2^2$ and $k_{48} = c_8/c_4^2$. All of the equations above may be readily derived from the definition of the weight-average molecular weight, $M_w = \Sigma c_i M_i/\Sigma c_i = \Sigma c_i M_i/c$, by substitution from the relationships $M_i = iM_1$ and those wherein the $c_i$'s may be defined in terms of the products of the appropriate $k$'s and powers of $c_1$.

Thermodynamic non-ideality is described in terms of a single virial coefficient $B$ [20], and the apparent weight-average molecular weight is given by

$$M_{w,app} = M_w/(1 + BM_w c).$$  (11)

Monomer–dimer, monomer–dimer–trimer, and monomer–dimer–tetramer associations may be considered to be subsets of the more extended associations

described above, and the appropriate equations obtained simply by setting the unneeded equilibrium constants equal to zero.

One may readily convert the equilibrium constants from a gram per liter scale to a mole per liter scale by recalling that the concentration on the molar scale, $C_i = c_i/iM_1$, and making appropriate substitutions. Thus, for the isodesmic association, the molar intrinsic equilibrium constant, $K = kM_1$. Similarly, $K_{12} = k_{12}M_1/2$ $K_{24} = k_{24}M_1$, $K_{48} = 2k_{48}M_1$, $K_{23} = 2k_{23}M_1/3$, $K_{36}$ $= 3k_{36}M_1/2$, $K_{26} = 4k_{26}M_1^2/3$, $K_{28} = 2k_{28}M_1^3$, and the molar virial coefficient, $B' = BM_1$.

By means of these equations, we may study to what extent one of these models may fit data generated by another.

## 3. Results

We will first examine the fitting of various imposter models to perfect isodesmic association data generated with a value of $M_1$ of 20000 daltons, a value of the intrinsic association constant, $k$, of 0.400 liters/gram and a value of $B$ of $1.000 \times 10^{-7}$ mole liters/gram$^2$. These latter values are nominally those found by Adams and Lewis [5] for an isodesmic association of $\beta$-lactoglobulin A. Fig. 1A demonstrates the fit of a non-ideal 1—2—8 model. While the quality of the fit is superior to that which Adams and Lewis obtained fitting this model to their data for $\beta$-lactoglobulin A, it can be clearly seen that this is not a satisfactory fit to these data. When the model is changed by the incorporation of a tetramer in the association reaction, a very good fit is obtained, as can be seen in fig. 1B. The obtained parameter values, their standard errors, and the root-mean-square (RMS) errors are given in table 1. It should be noted that, indicative of the quality of the two fits, the standard errors and also the RMS error are significantly larger for the 1—2—8 association than for the 1—2—4—8 association. Table 1 also shows the expected result that, when perfect 1—2—4—8 model data is generated using the parameters obtained for the best fit to the isodesmic model data, fitting the isodesmic model to the data gives values for the equilibrium constant and virial coefficient virtually identical to those used in the original data generating isodesmic model and moreover, the RMS error is essentially the same as that obtained in fitting the 1—2—4—8 model to isodesmic model data.

Since the values of the apparent weight-average molecular weight generated for the isodesmic model appear to be approaching a maximum below the molecular weight of a hexamer, this appeared to offer the possibility of being a reasonable terminal oligomer.
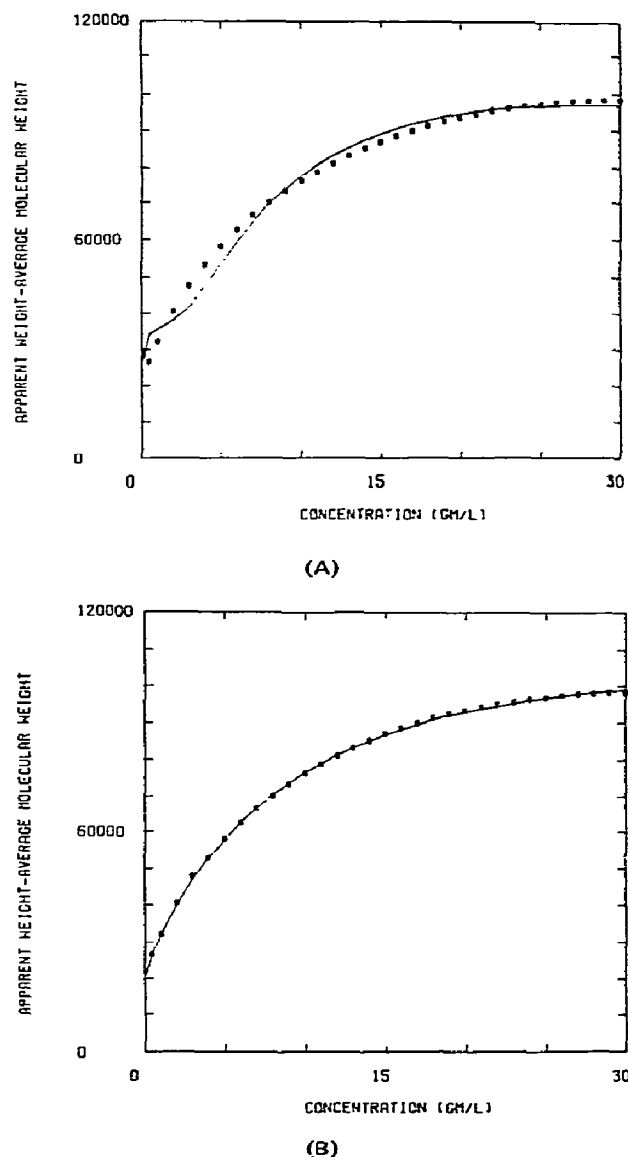


(A)



(B)

Fig. 1. Fits of models with the octamer as a terminal oligomer to data generated with an isodesmic model. (A) Non-ideal monomer—dimer—octamer. (B) Non-ideal monomer—dimer—tetramer—octamer. All parameter values are given in table 1.

Table 1
Fit of different models to isodesmic data ($k$ = 0.400, $B$ = 1.000 × $10^{-7}$)

| Model | Parameter values | RMS error |
|---|---|---|
| 1–2–8 (non-ideal) | $k_{12}$ = 18.612 ± 14.018<br>$k_{28}$ = (0.3088 ± 0.0656) × $10^{-2}$<br>$B$ = (0.6920 ± 0.0513) × $10^{-7}$ | 2850.5 |
| 1–2–4–8 (non-ideal) | $k_{12}$ = 1.2593 ± 0.0969<br>$k_{24}$ = 0.5385 ± 0.0460<br>$k_{48}$ = 0.1583 ± 0.0095<br>$B$ = (0.3047 ± 0.0131) × $10^{-7}$ | 347.5 |

Fit of isodesmic model to data generated with the 1–2–4–8 model parameters given above:

| | | |
|---|---|---|
| | $k$ = 0.3997 ± 0.0015<br>$B$ = (0.9986 ± 0.0079) × $10^{-7}$ | 335.0 |

Because the 1–2–4–8 model gave a very good fit, a 1–2–3–6 association seemed to be a logical model to examine. The results of fitting this model using eq. (5) and (6) are given in table 2. Initially, eq. (11) was used as well, but when a negative virial coefficient with a value about one-tenth that of the virial coefficient used for generating the data was obtained, it was decided to fit the data with an ideal model. The difference in the RMS error resulting from this was small. It should be noted that in working with simulated data generated with the equations used here, the only interpretation that can be given to a negative virial coefficient is that the fitting model does not go to a high enough order of association. In these circumstances, the virial coefficient is simply another fitting parameter and has no other physical significance.

The large values of the standard errors, the small value of $k_{23}$, and the large values of the parameter dependencies suggested, however, that the 1–2–3–6 association might not be an appropriate model. Since $c_6 = k_{36}k_{23}^2k_{12}^2c_1^6$, a very small value of $k_{23}$ forces $k_{36}$ to be very large if any significant amount of hexamer is required for the fit. This situation, where the dependency values are quite high, can be alleviated to some extent since the model equations may be written in an alternative form where the formation of each species of oligomer is written as occurring only from the monomer rather than involving any intermediate species. Thus,

$$M_w = M_1(c_1 + 2k_{12}c_1^2 + 3k_{13}c_1^3 + 6k_{16}c_1^6)/c \qquad (12)$$

and

$$c_1 = \underset{0<x<c}{\text{Root}} \; [c - (x + k_{12}x^2 + k_{13}x^3 + k_{16}x^6)]. \qquad (13)$$

It may be readily seen that $k_{13} = k_{12}k_{23}$ and that $k_{16} = k_{12}^2k_{23}^2k_{36}$. The results obtained by using these equations also appear in table 2. It seems obvious that since $k_{13}$ is virtually equal to zero and has a large standard error, the presence of trimer is quite unnecessary for obtaining the best fit, and fitting with a 1–2–6 model gives equally good results. When this is done, it can be seen in table 2 that the values of $k_{12}$ for both fits are virtually identical; $k_{26}$ has a predicted value of 0.1208 from the relationship $k_{16} = k_{12}^3k_{26}$, in reasonably good agreement with the fit value of 0.1116. The non-agreement of these values and the better value for the RMS error may be explained by the fact that the minimum of the sum of squares in parameter space is sufficiently shallow that the convergence factor used permitted this variation in parameter values and quality of fit, particularly since the sum of squares surfaces must be different for the different models used. The fits of all of the models described in table 2 to the isodesmic model data is shown in fig. 2; the slight differences between them cannot be discerned in a plot of this type, indeed, all of the theoretical curves appear as superimposed in fig. 2. While the quality of fit is better than for the 1–2–8 model, it is still obvious that these are not appropriate models for these data.

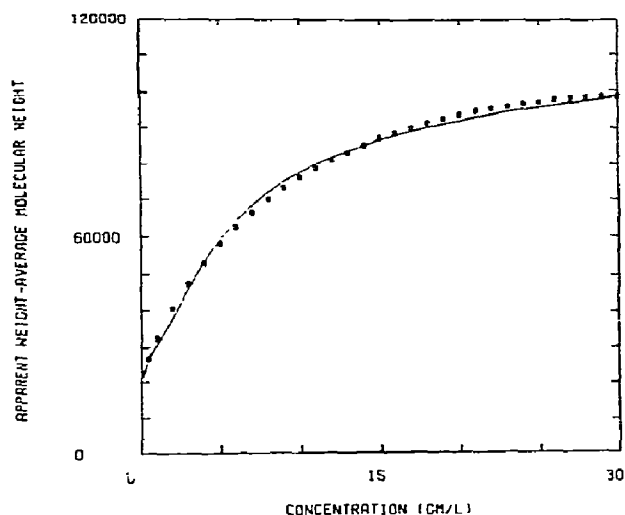We would next like to examine the effects of the addition of error to the isodesmic model data on the quality

Fig. 2. Fits of models with the hexamer as a terminal oligomer to data generated with an isodesmic model. Ideal monomer–dimer–hexamer and monomer–dimer–trimer–hexamer models are shown, but the difference between them cannot be discerned. All parameter values are given in table 2.

of fit obtainable with the different imposter models. The initial problem here is to determine the nature of the distribution and the magnitude that the error should have. We may reasonably expect that in the absence of systematic error, the solute concentration measured as a function of radial position in the ultracentrifuge cell should have normally distributed error. However, the concentration-dependent apparent weight-

Table 2
Fit of different models to isodesmic data ($k$ = 0.400, $B$ = 1.000 $\times 10^{-7}$)

| Model | Parameters | RMS error |
|---|---|---|
| 1–2–3–6 (ideal) | $k_{12}$ = 1.6568 ± 0.5952 $k_{23}$ = 0.0361 ± 0.0938 $k_{36}$ = 166.57 ± 814.00 | 1334.7 |
| 1–2–3–6 (ideal) | $k_{12}$ = 1.8691 ± 0.6453 $k_{13}$ = 0.1120 × $10^{-8}$ ± 0.1614 $k_{16}$ = 0.7889 ± 0.4766 | 1288.4 |
| 1–2–6 (ideal) | $k_{12}$ = 1.8695 ± 0.3005 $k_{26}$ = 0.1116 ± 0.0117 | 1266.0 |

average molecular weight calculated from this may be expected to have an error which must be, to some extent, dependent on the method used to calculate the molecular weight. One could not expect to obtain the same error distribution if the $M_{w,app}$'s were calculated from a sliding fit of ln $c$ vs. $r^2$ data as compared to those obtained by calculations using numerical differentiation to obtain $dc/dr$ as a function of $r$ and combining this with $c$ vs. $r$ data. Since these are but two of several methods for obtaining $M_{w,app}$ as a function of concentration, we have chosen the simple expedient of assuming that the error in $M_{w,app}$ is normally distributed and that the magnitude of the standard error is a constant fraction of the monomer molecular weight. This latter assumption appears reasonably justified by the published data of Adams and Lewis [5]. Based on this, we have chosen a standard error of 1000, 5% of the monomer molecular weight of 20000, as a minimum error. The error was generated using computer-generated random numbers provided by MLAB having a normal distribution with a mean of zero and a standard error of one. These were multiplied by 1000 and added to the isodesmic model data to obtain simulated noisy data. This was done five times independently, so that five new data sets were created. These five data sets could be grouped together for analysis or treated singly as data obtained from individual experiments. Another five data sets were generated with a standard error of 2000 and were used for individual and grouped analyses.

Table 3 gives the best-fit parameter values for the different models which were fit to the combined sets of isodesmic data with normal error having a standard error of 1000; the quality of these various fits may be seen in figs. 3A–3D. It should be noted that the values of the parameters are rather close to those obtained from fitting the perfect data although the RMS error values are larger, reflecting the greater dispersion of the data. The addition of error makes it more difficult to judge the quality of fit from the graphs. For this reason, we have also presented the fits in the form of difference graphs, wherein the differences between the data points and the fitting line are plotted as points about $\Delta M_{w,app} = 0$. Figs. 4A–4D show this for the various fits shown in figs. 3A–3D. These quite clearly show the random distribution of data around the fitting line for the isodesmic and 1–2–4–8 models, while the non-random distributions of the data around the fitting lines for the ideal and the non-ideal 1–2–6 models are equally ob-

Table 3
Fit of different models to combined isodesmic data with added
error: $k = 0.400$, $B = 1.000 \times 10^{-7}$, SE $= \pm 1000$

| Model | Parameters | | RMS error |
|---|---|---|---|
| Isodesmic | $k$ = | $0.3985 \pm 0.0214$ | 1057.7 |
| (non-ideal) | $B$ = | $0.9805 \pm 0.0111) \times 10^{-7}$ | |
| 1–2–6 | $k_{12}$ = | $1.7408 \pm 0.1703$ | 1714.6 |
| (ideal | $k_{26}$ = | $0.1184 \pm 0.0076$ | |
| 1–2–6 | $k_{12}$ = | $2.5897 \pm 0.2456$ | 1406.2 |
| (non-ideal) | $k_{26}$ = | $0.0802 \pm 0.0052$ | |
| | $B$ = | $(-0.1171 \pm 0.0125) \times 10^{-7}$ | |
| 1–2–4–8 | $k_{12}$ = | $1.2501 \pm 0.1368$ | 1107.6 |
| (non-ideal) | $k_{24}$ = | $0.5362 \pm 0.0654$ | |
| | $k_{48}$ = | $0.1582 \pm 0.0135$ | |
| | = | $(0.2883 \pm 0.0186) \times 10^{-7}$ | |

vious. Such systematic patterns in the differences may indicate non-independent error or the use of an incorrect model; thus they are valuable graphs to examine in trying to distinguish between various models. Here we may clearly reject the ideal and non-ideal 1–2–6 models.

It is now apparent that the principal problem that we face is that of attempting to discriminate between an isodesmic model and a 1–2–4–8 model of self-association when fitting data generated for the isodesmic model, and that neither the fitting statistics or the various graphical presentations has enabled us to do this. We have examined the fitting of each of the five individual data sets and compared the values of the fitting parameters obtained with those obtained from the combined data set for the isodesmic model and that of the 1–2–4–8 model, all with a standard error of 2000. For data which was generated with the isodesmic model, the values of the parameters obtained when fitting individual isodesmic models are in quite good agreement with the values obtained when fitting the combined data, while if the 1–2–4–8 model is used for fitting, the values of the parameters for the individual data sets vary quite widely from those obtained when fitting the combined data. Unfortunately, this cannot be used as a means of discriminating between the two models, because when the 1–2–4–8 model is used for generating the data, exactly the same results are obtained, and this it would appear that the

lesser dependancy of the isodesmic model fitting parameter values on the effects of error in individual data sets is simply due to the lesser number of parameters being fitted and to the nature of the model.

Thus far, only the 1–2–4–8 model has proved to be a successful imposter model for fitting data generated using an isodesmic model of self-association. What we would now like to examine is how uniquely defined in terms of parameter values the 1–2–4–8 model must be so that it can not be readily differentiated from the isodesmic model. Our initial approach to this has been to generate data using different parameter values for the 1–2–4–8 model and fit with an isodesmic model, using difference plots as a means of determining the quality of the fits.

In figs. 5A, B, and C we demonstrate the effects of varying the values of $k_{12}$, $k_{24}$, and $k_{48}$ respectively of the 1–2–4–8 generating model upon the quality of fit obtained for the isodesmic fitting model. The initial values of the parameters are those which gave the best fit of the 1–2–4–8 model to the isodesmic data. The individual parameters were then doubled or halved while holding the other parameter values constant. $k_{12}$ appears to have the greatest effect on quality of fit of the isodesmic model, as either doubling or halving its value results in difference plots where the lack of fit would be very clearly discernable even in the presence of added error. While $k_{24}$ appears to have the least effect, doubling or halving its values also clearly produces difference plots which would clearly demonstrate a lack of fit. $k_{48}$ is intermediate between $k_{12}$ and $k_{24}$ in terms of its effects on the difference plots. It is of interest to note that the locations, the magnitudes, and the directions of the deviations from fitting are very clear functions of which parameter is involved.

Since doubling or halving each parameter while leaving the others unperturbed gave results where the lack of fit was clearly discernable, we next explore the effects of doubling or halving $k_{12}$, then seeking new values of $k_{24}$, $k_{48}$, and $B$ that give the best fit for the original 1–2–4–8 model, and then fitting the isodesmic model. The results of this are shown in table 4. It is apparent that doubling the value of $k_{12}$ causes the value of $k_{24}$ to be halved and the values of $k_{48}$ and $B$ are increased, but to a lesser extent, with the effect of the increase in the value of $B$ compensating for the effect of the increase in the value of $k_{48}$. Halving the value of $k_{12}$ produces exactly the reverse of this.
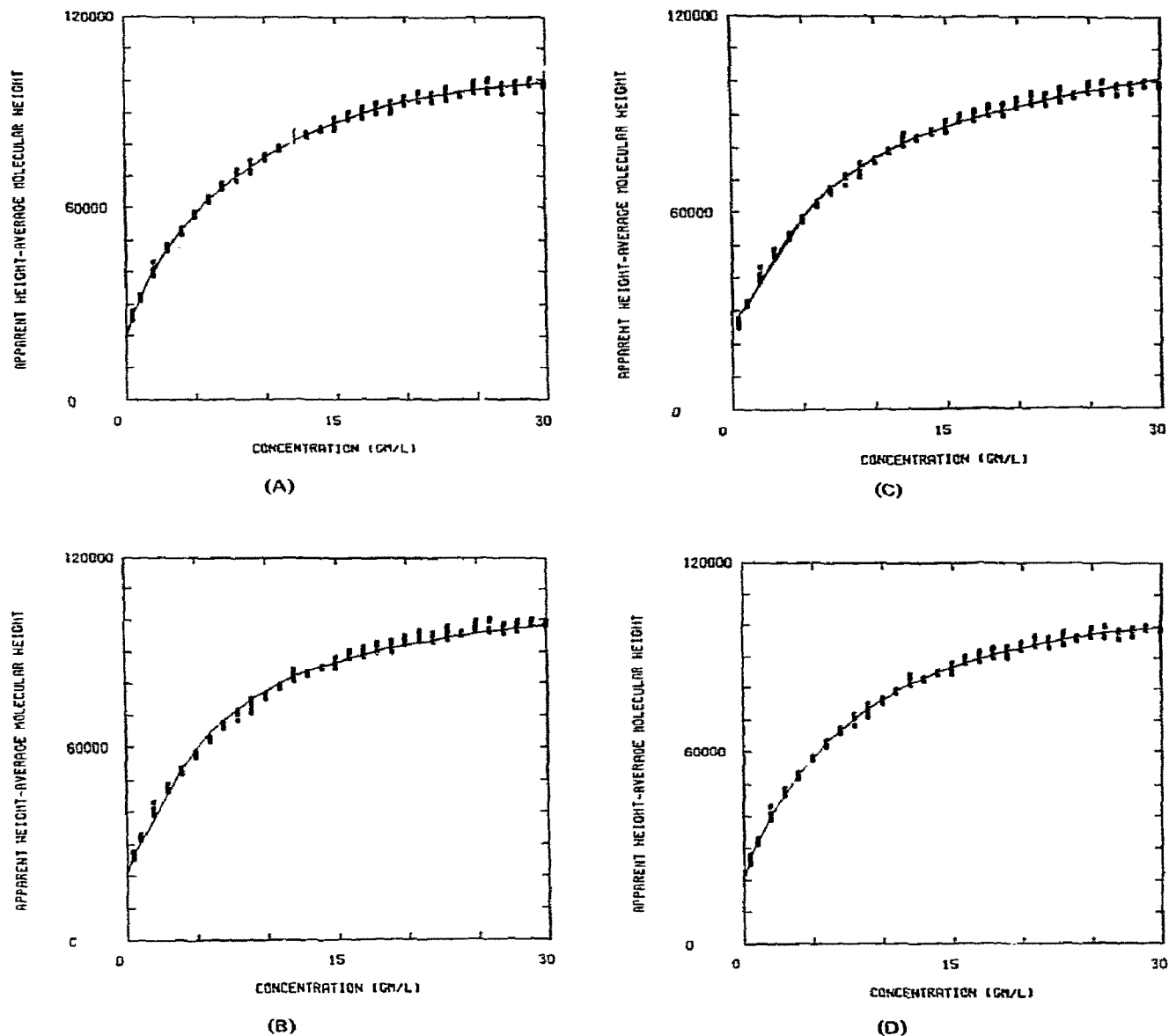
(A)



(C)



(B)



(D)

Fig. 3. Fits of different models to five combined sets of data generated with an isodesmic model and having added error. (A) Non-ideal isodesmic model. (B) Ideal monomer—dimer—hexamer model. (C) Non-ideal monomer—dimer—hexamer model. (D) Non-ideal monomer—dimer—tetramer—octamer model. All parameter values are given in table 3.

In both cases the RMS error is small enough that these would be called good fits by that criterion. The effects on the values of $k$ and $B$ for the isodesmic model is insignificant, and while the RMS errors are larger, the in-

crease is not significant. Thus, it is particularly interesting to observe in fig. 6 that while the deviations from fitting are too small to be significant in the presence of added error over most of the concentration range,
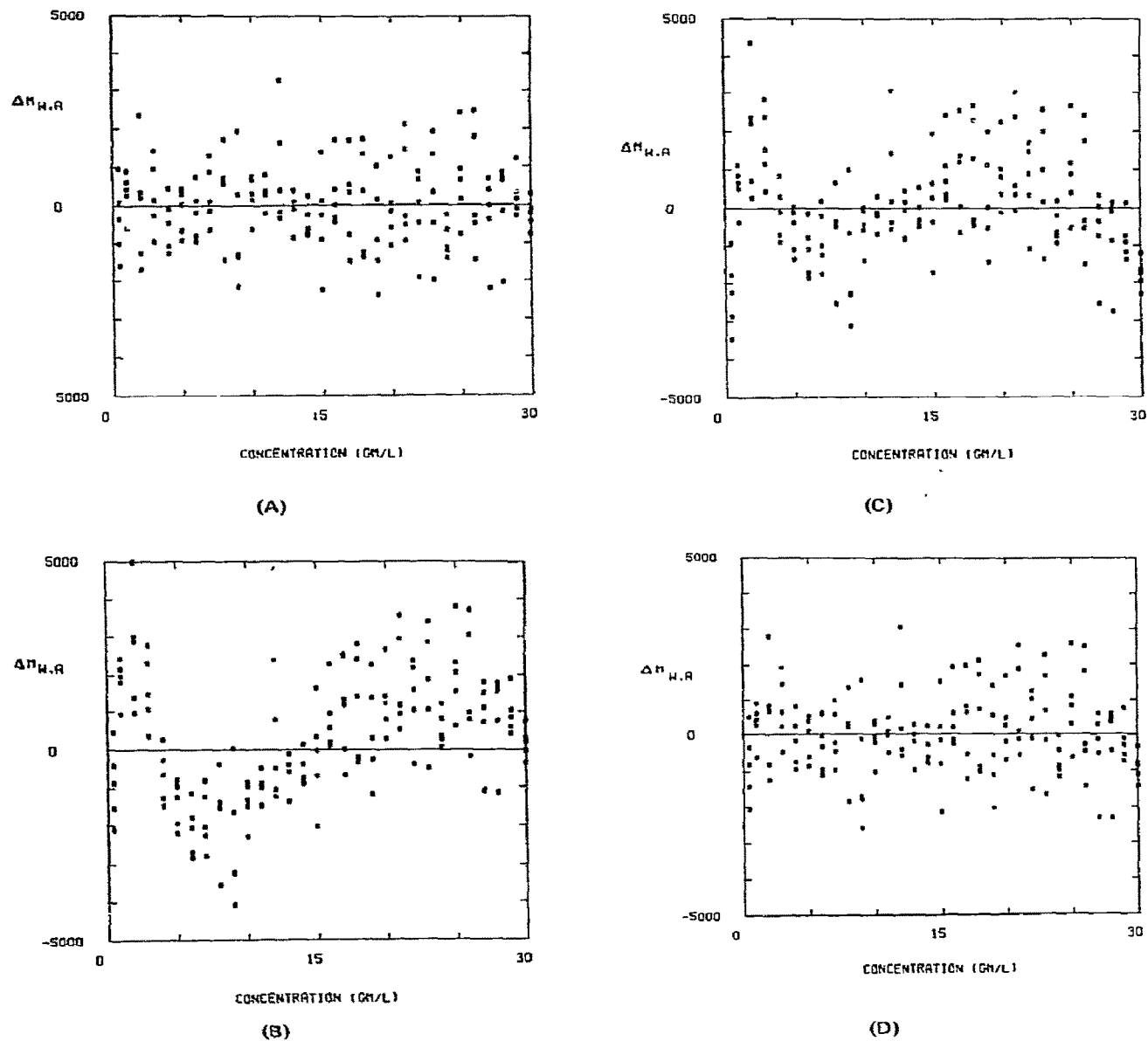
Fig. 4. Plots of the differences between the data and the fitting curves for the different models shown in fig. 3. (A) Non-ideal isodesmic model. (B) Ideal monomer—dimer—hexamer model. (C) Non-ideal monomer—dimer—hexamer model. (D) Non-ideal monomer—dimer—tetramer—octamer model. All parameter values are given in table 3.

there is a marked increase in the magnitude of the deviations near zero concentration, and these are of such magnitude as to have a reasonable probability of detection. From th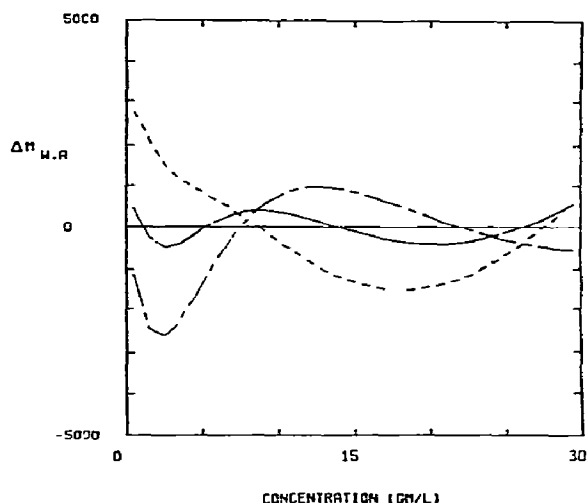is result we may surmise that there may be a fairly specific relationship between the value of the equilibrium constant of an isodesmic model and the values of the equilibrium constants of an equivalent 1—2—4—8 model.

In order to explore this possibility, we have generated data with the isodesmic model for a range of values of $k$ with a constant value for $B$, and have then fit this with the 1—2—4—8 model. The results of this are shown in table 5. Examination of the RMS error values show that the quality of fit is quite good up through a value
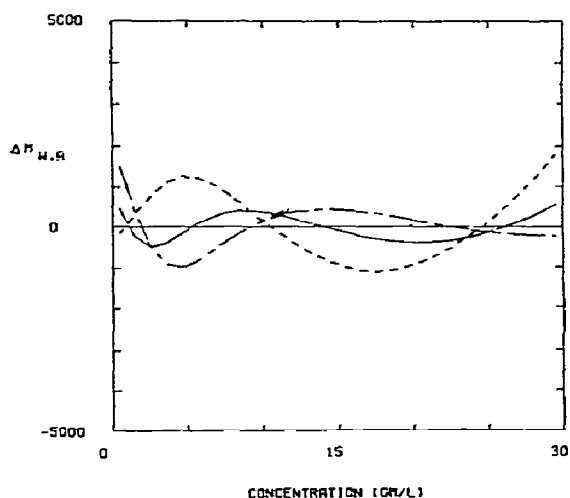
Table 4
Effect of variation of equilibrium constants of 1—2—4—8 model on 1—2—4—8 and isodesmic model fitting parameters
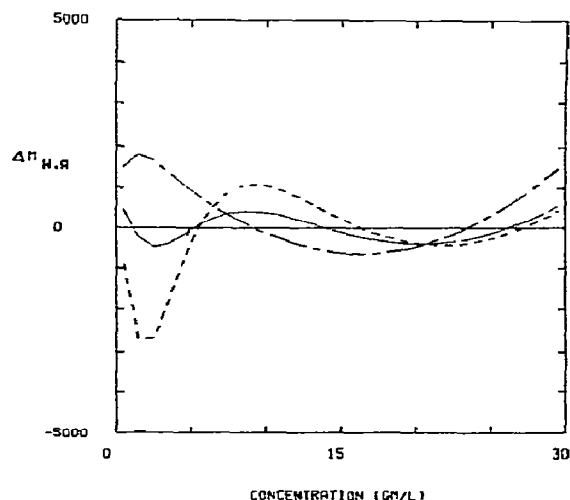
| | | | |
|---|---|---|---|
| $k_{12}$ | 1.2593 | 2.5186 | 0.6296 |
| $k_{24}$ | 0.5385 | 0.2834 | 1.1762 |
| $k_{48}$ | 0.1583 | 0.2381 | 0.1090 |
| $B \times 10^7$ | 0.3048 | 0.3545 | 0.2510 |
| RMS | – | 531.6 | 501.3 |
| $k$ | 0.4000 | 0.4011 | 0.3991 |
| $B \times 10^7$ | 1.0015 | 1.0069 | 0.9977 |
| RMS | 325.7 | 611.9 | 599.5 |

of 0.400 for $k$ for the isodesmic model and that it deteriorates very rapidly at higher values. Examination of table 5 shows that the region where the RMS error values are low is also a region where the ratio of $k_{12}/k_{24}$ is small and is increasing slowly and linearly from 1.35 to 2.00, and shows that the ratio of $k_{12}/k$ is also linear and increasing slowly from 2.5 to 3.0. At higher values of $k$ both of these ratios start increasing rapidly, the value of $k_{24}$ decreases, and the RMS error values are such that there is little question concerning our ability to discriminate between the two models.

At the lowest values of $k$ the small values of $k_{48}$ and



(A)



(B)



(C)

Fig. 5. Plots showing the effects of specific association constants on the differences between the data generated for a non-ideal monomer—dimer—tetramer—octamer model ar d the fit with a non-ideal isodesmic model. The initial values of the generating parameters are those given in table 1 and are indicated by the solid line in each figure. The short dashed line shows the effect of doubling a particular equilibrium constant and the alternate long and short dashed line shows the effect of halving that equilibrium constant. (A) Effect of variations of $k_{12}$. (B) Effect of variations of $k_{24}$. (C) Effect of variations of $k_{48}$.

Table 5
Effect of variation of $k$ ($B$ = 1.000 × $10^{-7}$) on fitting parameters of 1–2–4–8 model

| $k$ | 0.0500 | 0.1000 | 0.2000 | 0.3000 | 0.4000 | 0.6000 | 0.8000 |
|---|---|---|---|---|---|---|---|
| $k_{12}$ | 0.1273 | 0.2544 | 0.5048 | 0.8033 | 1.1707 | 2.3074 | 5.1325 |
| $k_{24}$ | 0.0944 | 0.1914 | 0.3828 | 0.5176 | 0.5866 | 0.5228 | 0.2816 |
| $x_{48}$ | 0.0099 | 0.0243 | 0.0520 | 0.0913 | 0.1488 | 0.3932 | 1.4132 |
| $B \times 10^7$ | 0.7709 | 0.7238 | 0.5243 | 0.3880 | 0.2884 | 0.1558 | 0.0745 |
| RMS | 68.7 | 70.5 | 56.7 | 155.9 | 356.4 | 921.4 | 1569.0 |

Table 6
Effect of variation of $k$ ($B$ = 1.000 × $10^{-7}$) on fitting parameters of 1–2–3 and 1–2–4 models

| $k$ | 0.0500 | 0.0500 | 0.1000 | 0.2000 |
|---|---|---|---|---|
| $k_{12}$ | 0.0655 | 0.1272 | 0.2317 | 0.2087 |
| $k_{23}$ | 0.1702 | — | — | — |
| $k_{24}$ | — | 0.0835 | 0.2026 | 1.4164 |
| $B \times 10^7$ | − 0.8137 | 0.2111 | − 0.3569 | − 0.3381 |
| RMS | 171.9 | 75.1 | 190.8 | 1060.8 |

the large values of $B$ suggest the possibility that a monomer—dimer—trimer model or a monomer—dimer—tetramer model might fit well enough to be indistinguishable from the generating isodesmic model. The results of investigating these possibilities are given in table 6. It is quite apparent from the large RMS error and the negative value of $B$ that a 1–2–4 model is a poor fit for $k$ = 0.2000; when $k$ = 0.1000 only the negative value of $B$ indicates that the model may be unsuitable. When $k$ = 0.0500 it is quite obvious that the 1–2–4 model now gives an excellent fit, but the negative value of $B$ again indicates that the 1–2–3 model does not have a high enough order of association.

## 4. Discussion

In the preceding section it was shown that, over about a ten-fold range of the value of the equilibrium constant for the generating non-ideal isodesmic model, the non-ideal 1–2–4—8 model was the only imposter model which gave a satisfactory fit to the generated data. The quality of the fit as indicated by the value of the RMS error and with even greater sensitivity by the difference plots, was shown to be a suitable criterion for the rejection of other possible models. It must be re-
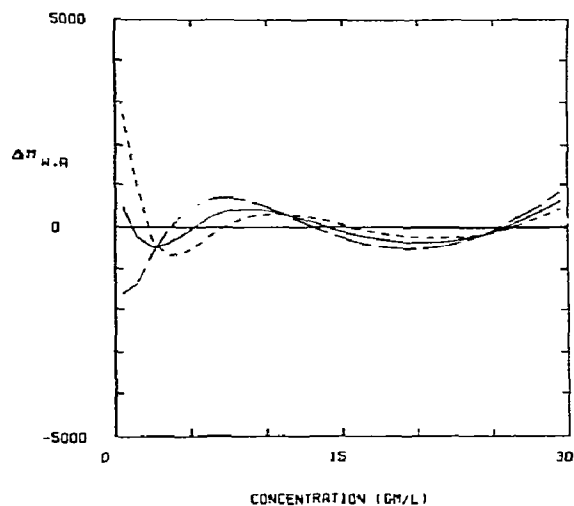


Fig. 6. Difference plots showing the effect of the adjustment of $k_{24}$, $k_{48}$, and $B$ compensating alteration of $k_{12}$ on the fit of a non-ideal isodesmic model to non-ideal monomer—dimer—tetramer—octamer data. The solid line is the fit obtained for the parameters given in table 1; the short dashed line is the fit obtained when $k_{12}$ is doubled; the long dashed line is for the fit obtained when $k_{12}$ is halved.

membered that the successful impostor model was rather rigidly constrained in terms of the relationships of the values of its equilibrium constants to the value of the equilibrium constant of the generating isodesmic model. This would tend to imply that it is more probable that an isodesmic association can be fit with an impostor specific association model than that a specific association can be fit with an impostor isodesmic association model.

In dealing with real systems, we must consider if the criteria which were used to reject unsuitable models are still appropriate as well as what additional criteria are

required. It is obvious that the criteria regarding quality of fit are as appropriate here as for model systems, but in addition, the values obtained for the fitting parameters must have some significance in physical reality. Specifically, the values of the virial coefficients and the values of the changes of free energy, enthalpy, and entropy for the associations, as calculated from the equilibrium constants, must be appropriate. Previously we have rejected negative values for the virial coefficients on the basis that these were simply fitting parameters like the association constants and that they indicated that a higher order of association than was described by the model was required. This may not be an entirely valid assumption in a real situation. The use of a single term to describe the effects of the multi-body interactions which are involved in thermodynamic non-ideality is quite simplistic, but it has generally proved to be reasonably adequate for the practical analysis of associating systems [1,4]. Positive virial coefficients are generally attributed to excluded volume and repulsive charge effects while negative virial coefficients are attributed to attraction between molecules [21]. It might be possible to discriminate between a higher order of association and weak intermolecular attractive forces as the cause of a negative virial coefficient by analyzing the system as if the higher order association occurred, determining the equilibrium constant, and then using the magnitude of the calculated change in free energy as a criterion for choosing between these two possibilities.

The question which we must now attempt to answer is what we can do when we are confronted with the problem of distinguishing between a specific association and an isodesmic association which appear to be equally good fits for a given set of data. There are several possible answers, none of which is necessarily sufficient in itself, but taken together might suggest that one model could be preferred over the other. One of these is to obtain apparent weight-average molecular weight versus concentration data at a variety of temperatures and attempt fitting both models to all of this data. If one model can be fit and the other cannot, or if one can be fit significantly better than the other, this could be valid reason for preferring that model. Assuming that both models can be fit equally well, the physical significance of the values of the changes of the free energy, the enthalpy and the entropy for the different associations might be of assistance in decid-

ing between models. It is particularly important that the apparent weight-average molecular weight versus concentration data be obtained to as high a concentration as is practical since the high concentration data are the most significant in determining the extent of the association and the value of the virial coefficient.

Another approach that might be of assistance in determining the most appropriate model is the study of the association by means of sedimentation velocity. Holloway and Cox [22] and Cox [23,24] have demonstrated by computer simulation that an isodesmic association gives a boundary that differs significantly in shape from that obtained for definite associations because of the presence of significant concentrations of larger oligomers than would be observed in the model involving a definite association. Thus, the boundaries observed in sedimentation velocity studies could be compared with those which can be predicted by simulation methods for the different models, and a choice might be made between them. Gel filtration studies can be treated in a similar manner and should be equally useful [25,26]. Additionally, the optimum model should be consistent with what is known regarding the structure and function of the associating molecule.

It must be kept in mind that it is quite realistic to expect that it may not be possible to obtain an unequivocal answer and that one will have to be satisfied with describing two or more models of association which are equally valid by present criteria. Improvements in the precision of analytical techniques and a better understanding of the mechanisms involved in the process of self-association are essential but not necessarily sufficient for any resolution of this problem.

References

[1] H. Fujita, Foundations of ultracentrifugal analysis (Wiley, New York, 1975).
[2] J.R. Cann, Interacting macromolecules (Academic Press, New York, 1970).
[3] D.C. Teller, in: Methods in enzymology, Vol. 27, eds. C.H.W. Hirs and S.N. Timasheff (Academic Press, New York, 1973) p. 346.
[4] E.T. Adams, Jr., W.E. Ferguson, P.J. Wan, J.L. Sarquis and B.M. Escott, Separation Science 10 (1975) 175.
[5] E.T. Adams, Jr. and M.S. Lewis, Biochemistry 7 (1968) 1044.
[6] K.E. Van Holde and G.P. Rossetti, Biochemistry 6 (1967) 2189.

[7] R. Townend and S.N. Timasheff, J. Am. Chem. Soc. 82 (1960) 3168.

[8] S.N. Timasheff and R. Townend, J. Am. Chem. Soc. 83 (1961) 464.

[9] T.T. Herskovits, R. Townend and S.N. Timasheff, J. Am. Chem. Soc. 86 (1964) 4445.

[10] D.E. Roark and D.A. Yphantis, Ann. N.Y. Acad. Sci. 164 (1969) 245.

[11] K.E. Van Holde, G.P. Rossetti and R.D. Dyson, Ann. N. Y. Acad. Sci. 164 (1969) 279.

[12] M.E. Magar, Data analysis in biochemistry and biophysics (Academic Press, New York, 1972).

[13] R.I. Shrager, J. Assn. Comp. Mach. 17 (1970) 446.

[14] D.W. Marquardt, Soc. Indust. Appl. Math. J. App. Math. 11 (1963) 431.

[15] K. Levenberg, Quart. App. Math. 2 (1955) 164.

[16] G.D. Knott and D.K. Reece, in: Proceedings of the International Conference ONLINE '72, Vol. 1, p. 497.

[17] G.D. Knott and R.I. Shrager, SIGGRAPH Notices 6 (1972) 138.

[18] L. Endrenyi and F.H.F. Kwong, Acta Biol. Med. Germ. 31 (1973) 495.

[19] K.C. Ingham, H.A. Saroff and H. Edelhoch, Biochemistry 14 (1975) 4745.

[20] E.T. Adams, Jr., Biochemistry 4 (1965) 1646.

[21] C. Tanford, Physical chemistry of macromolecules (Wiley, New York, 1961).

[22] R.R. Holloway and D.J. Cox, Arch. Biochem. Biophys. 160 (1974) 595.

[23] D.J. Cox, Arch. Biochem. Biophys. 142 (1971) 514.

[24] D.J. Cox, Arch. Biochem. Biophys. 146 (1971) 181.

[25] J.K. Zimmerman, D.J. Cox and G.K. Ackers, J. Biol. Chem. 246 (1971) 4242.

[26] G.K. Ackers, in: The proteins, 3rd Ed., Vol. 1, eds. H. Neurath and R.L. Hill (Academic Press, New York, 1975) p. 1.