

Statistical analysis of data pertaining to complex state systems by stepwise regression with reformulated parameters; Application to spectroscopically monitored hemoglobin oxygen binding data¹

William G. Gutheil *

Department of Biochemistry, Meharry Medical College, Nashville, TN 37208, USA

Received 14 July 1997; revised 12 September 1997; accepted 12 September 1997

Abstract

A method is described for the statistical analysis of data pertaining to complex state systems, based on the concept of reformulating the parameters describing the system as a hierarchy of interactions, and this method demonstrated on the analysis of spectroscopically monitored hemoglobin oxygen binding data [K. Imai, *Biophys. Chem.* 37 (1990) 197–210]. The concept of reformulation was first extended to state parameters other than ΔG° s, such as the extinction coefficients (ϵ s) associated with different ligation states during hemoglobin oxygen binding. The reformulated parameters are incrementally allowed to vary in the data fitting procedure, and the statistical significance of the added parameters tested by F and Kolmogorov–Smirnov tests. The result of this method is the minimal set of statistically significant parameters required to describe the data. The hierarchical nature of reformulated parameters allows the physical significance of the subset of statistically significant parameters to be discussed even when all reformulated terms may not be statistically significant. Applying this method to hemoglobin oxygen binding data with the reformulated Adair model demonstrated that at least two, and at most three, of the four reformulated Adair constants are statistically significant. A reformulated square model was found to give a statistically indistinguishable fit from the Adair model, with the statistically significant thermodynamic terms essentially those proposed by Linus Pauling in 1935. A change in $\Delta\epsilon$ with subsequent oxygen binding events was found to be significant in both models. These results are consistent with a model for hemoglobin oxygen binding where a subunit changes its conformation upon oxygen binding, and affects the conformation of adjacent subunits. © 1998 Elsevier Science B.V.

Keywords: Hemoglobin; Thermodynamics; Statistical analysis

* Corresponding author. Tel.: +1-615-327-6749; fax: +1-615-327-6442.

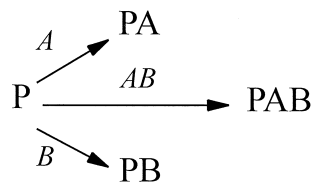
¹ This work was supported by Grant 5D34 MB00003 from the Human Resources and Services Administration. This manuscript is dedicated to the memory of Linus Pauling (1901–1994), and to his contributions towards our understanding of macromolecular structure and function.

1. Introduction

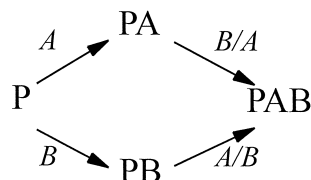
How a mathematical problem is formulated often determines the success or failure of finding a solution. This also applies to the statistical analysis of data where how the physical model of the system in question is formulated can influence both the success of the statistical analysis as well as the utility of the derived parameters for subsequent interpretation. Proteins which can bind multiple ligands with thermodynamic interactions between separate ligand binding events represent examples of complex systems with several alternatives for how the systems are formulated. We describe here a general method for the statistical analysis of data pertaining to complex state systems based upon the method of formulating such systems as a hierarchy of interactions (reformulation),² and demonstrate this method on the statistical analysis of hemoglobin oxygen binding data.

A simple hypothetical example of a protein-ligand binding system, a protein which can bind two different ligands on two separate sites, is shown in Scheme 1. Three ways of formulating the ΔG° s describing this system are given. Panel A shows this system formulated in terms of ΔG° s of assembly of the complexes PA, PB, and PAB from the free components of P, A, and B, where A represents the ΔG° of assembly of PA from P and A, B represents the ΔG° of assembly of PB from P and B, and AB represents the ΔG° of assembly of PAB from P, A, and B. Panel B shows this system formulated in terms of stepwise constants, using the concept of conditional ΔG° s [1] where A and B are as defined above and A/B represents the binding of A with B already bound and B/A represents the binding of B with A already bound. Panel C shows this system formulated using the concept of interaction free energy, where a and b are identical to A and B described above, and ab is the ΔG° of interaction and represents the mutual effect (interaction) of A binding on

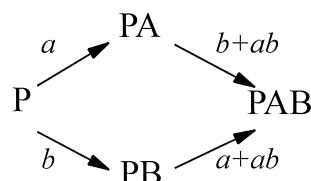
A. ΔG° s of assemble for PAB system from free P, A and B.



B. Stepwise ΔG° s for PAB system.



C. Reformulated ΔG° s for PAB system.



Scheme 1.

B binding and vice versa. Panel A requires three terms to describe the assembly of the three complexes from the free components. Panel B is described by four terms. However, the ΔG° of assembly of a complex is independent of path (it is a function of state) and this requires that $A + B/A = B + A/B$, i.e., these four terms overdetermine this system. Three of these parameters can be considered independent, with the fourth defined by the relationship just given. Panel C is also described by three terms.

For the statistical analysis of data pertaining to the system shown in Scheme 1, the model could in principle be described using any one of the three formulations shown, and the question is—is one of these formulations preferable to the others? The formulation shown in Panel A is the most direct way of describing this system, since each complex is associated with a ΔG° of assembly, and this requires three parameters. The formulation shown in Panel B is described by four parameters, of which three are independent. The formulation shown in Panel C also

² This method was originally given the name of the unique and independent parameters (UIP) formulation. Because the full scope of this method, and its potential applications, are as yet unknown such a name may be premature. Therefore, this method will simply be referred to here as reformulation.

requires three parameters to describe the full *potential* complexity of this system. However, if data for this system demonstrates no discernable effect of A binding on B binding, and vice versa, then fixing the value of ab to zero will not result in a significantly worse fit to the data. Data pertaining to this system could in this case be described by the two parameters a and b , and the constraint $ab = 0$. If however, the same data is fit with the three parameters a , b , and ab as adjustable parameters, or with the model formulated in Panels A and B, the fit parameters will have larger standard errors and lower statistical significance than from the fit with only two parameters. Because the data is adequately described by the two parameters a and b , no further information is available in the data to define the added third parameter ab . In order to demonstrate the statistical significance of $ab \neq 0$, the fit to the data would be performed with ab constrained to zero, and then with the value of ab as an adjustable parameter. Appropriate statistical tests can then determine the significance of $ab \neq 0$. In general, it is desired to find a minimal set of parameters which can adequately fit the data in question—additional parameters have little or no statistical significance. An effective way of performing this task for systems described by state parameters is presented in the following analysis.

The system shown in Scheme 1 is based upon the well known concept of pairwise interactions between two ligands, a concept which can be traced back at least to the classic analysis of hemoglobin oxygen binding data presented by Pauling [2], and this concept has been developed in detail by Weber [1,3]. Although the concept of pairwise interactions is very useful for treating systems with only two interacting ligands, it cannot be used to treat the full potential complexity of systems with more than two interacting ligands. A general method for formulating complex thermodynamic models has been developed which takes a system of N complexes, described by $N \Delta G^\circ$ s for assembly of each complex from its free components, and redefines it in terms of a new set of $N \Delta G^\circ$ s which reflect the complete hierarchy of interactions possible between individual ligand binding events [4,5]. This method for reformulating² complex thermodynamic models logically extends the concept of interaction energy from the well

known simple case of two interacting ligands as defined above to the general case of N interacting ligands as described previously for systems not showing monomer–multimer equilibria [4] and for systems showing monomer–multimer equilibria [5]. Reinhart [6] had previously introduced the concept of higher than second order interaction terms for systems with more than two interacting ligands, although the mathematical definition of higher order terms proposed by Reinhart is fundamentally different than that used here. Reformulated parameters are useful for both the theoretical analysis of complex physical systems as previously demonstrated by the analysis of cooperativity in a dimeric protein [7], as well as for the statistical analysis of data pertaining to such systems as demonstrated in the following analysis of hemoglobin oxygen binding data.

The oxygen binding properties of hemoglobin serve as the classic example of allosteric regulation of an important biological phenomenon [8,9]. Experimental methods for studying hemoglobin oxygen binding can be divided into two classes; those in which the α and β subunits of the $\alpha_2\beta_2$ hemoglobin tetramer are not distinguishable, and those in which the α and β subunits are distinguishable. Spectroscopically (UV–Vis) monitored hemoglobin oxygen binding data does not allow α and β subunits to be distinguished and is the basis of the analysis presented below. Experiments using NMR monitored titrations of hemoglobin (reviewed in Ref. [10]) and experiments using ligation inert hemoglobin variants (reviewed in Ref. [11]) are capable of distinguishing between oxygen binding to α and β subunits. Although it is expected that the method of reformulation will prove useful in these cases as well, this is beyond the scope of the present analysis. In spite of the vast amount of research devoted to the elucidation of the thermodynamic properties of hemoglobin oxygen binding using spectroscopically monitored titrations, the statistical and physical significance of the fit Adair binding constants still generates considerable debate (Ref. [11–17] and references therein). The physical significance of experimentally determined parameters cannot reasonably be discussed until the statistical significance of the parameters has been established. Comparison of alternative physical models also requires the application of statistical tests to establish the statistical confidence in one

model over alternative models. As a demonstration of the application of reformulated physical models to the analysis of data pertaining to a complex physical system, it is used here to analyze the hemoglobin oxygen binding data of Imai [16].

2. Theory

The method described here is related to statistical methods where the interaction between variates is analyzed [18]. In standard two-way ANOVA pairwise interaction terms between variates (categories) are determined and their statistical significance is tested. The total variation of the measurements in two-way ANOVA is described by the expression

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \epsilon_{ijk} \quad (1)$$

where $\alpha\beta_{ij}$ represents interaction between two variates, and is analogous to the term ab in Scheme 1C describing the pairwise interaction between ligands. More complex systems with three or more factors are described by expression with higher order terms, and data pertaining to such systems can be analyzed by converting the data into categorical frequencies followed by fitting with log-linear models (Ref. [18]; BMDP program 4F, Ref. [19]). The log-linear equation for a three way analysis is

$$\ln F_{ijk} = \theta + \lambda_{Ai} + \lambda_{Bj} + \lambda_{Ck} + \lambda_{ABij} + \lambda_{ACik} \\ + \lambda_{BCjk} + \lambda_{ABCijk} \quad (2)$$

where A, B, and C refer to different treatments, i , j , and k refer to subgroups (levels) of treatments, and $\ln F_{ijk}$ is the natural log of the frequency of a given cell of the three subgroups. The terms containing an A, B, or C are first order effects, terms containing AB, AC, or BC are second order effects and the term containing ABC is a third order effect. These second and third order interaction terms are analogous to the second and third order interaction terms for the PABC protein-ligand binding system described previously (Figure 5, Ref. [4]). The analysis of data with log-linear models involves the sequential addition or elimination of second and higher order interaction terms in the fitting process followed by statistical testing to determine the significance of the added or removed parameter.

Reformulated parameters are mathematically analogous to the statistical concept of hierarchical interactions, except that they represent the hierarchy of physical interactions which might occur in a system. In the present case, the models describing the system are nonlinear, and the principal advance in methodology is the use of reformulated parameters as the regressors in the nonlinear fitting of data. This method is based upon the property that higher order reformulated parameters reflect increasingly complex interactions between events and may or may not be physically and/or statistically significant. For example, the term ab in Scheme 1 will be statistically indistinguishable from 0 in two cases; (1) if $ab = 0$, and (2) if the data provides little or no information concerning the value of ab . In both cases, fitting the data with the constraint $ab = 0$ would provide estimates for a and b . Attempting to include ab in the fit would in the first case provide an estimate for ab close to 0. In the second case, there would be no information concerning ab and its value could not accurately be determined. The statistical analysis of data pertaining to this system would involve fitting the data with a and b as adjustable parameters, first with ab fixed to its null value³ (0, no interaction) and then with ab as an adjustable parameter. A statistical test can then be used to determine the significance of $ab \neq 0$. (It is tempting to ask how we can prove that $ab = 0$, however, the difference of ab from 0 will generally only be one of degree.) The statistical analysis of data using reformulated models

³ Reformulated parameters represent the hierarchy of interactions which can occur between events in a system. The terms primary and first order are used interchangeably to denote those parameters associated with a protein-ligand binding event in the absence of effects from other events. The term second order refers to parameters associated with pairwise interactions between events. The term third order refers to parameters describing the mutual interaction between three events, etc. The term zero order refers to state parameters other than ΔG° 's, such as extinction coefficients, where parameters are required to describe properties of free components. Base terms refer to zero and first order terms. The term null value is used to denote the value for a parameter corresponding to no interaction between events and as such is a concept which can only be applied to second and higher order parameters. The null value is 0 for a ΔG° of interaction but would be 1 for the equivalent equilibrium constant interaction term such as the term α used by Pauling [2].

begins by setting all the reformulated parameters except the base terms to their null values, and fitting the data with only the set of base terms allowed to vary. The fit is then repeated with individual second order terms also allowed to vary in the fitting procedure. This process is repeated, adding one parameter at a time, for all of the parameters up to highest order.

The standard method for determining the statistical significance of the sequentially added terms during stepwise regression with linear models is the *F*-test, which is based upon the assumption of normally distributed errors. When the data are well fit by the fitting function, the residuals (errors) are generally expected to be normally distributed, however, when the data are poorly fit then the residuals will probably not be normally distributed. A non-parametric (normal distribution independent) test of differences in distribution between two sets of values is the Kolmogorov–Smirnov (K–S) test [18,20]. Both the *F* and K–S methods are implemented and compared below.

In the analyses presented below, all possible combinations of hierarchically added parameters were tested. The term hierarchical is used here in the same sense as this term is used in connection in stepwise log–linear regression [18,19], where all lower order terms modulated by the higher order term must be significant before the higher order term can be included in the analysis. The methodology adopted here is analogous to that for log–linear regression and the reader may wish to consult Sokal and Rohlf [18] and Dixon [19] for additional details. The reader may also find the following analysis more tractable if reference is made to Figure 5 from Gutheil and McKenna [4], using *abc* as representative of a high order term and *ab*, *ac* and *bc* as representative of lower order terms. When terms above second order are found to be statistically significant it is necessary to consider the possibility that a high order term (above second order) may be statistically different than its null value while a related lower order term (one modulated by the high order term) may be statistically indistinguishable from its null value. This is a consequence of the fact that the lower order term, although real in the sense that such an interaction can occur in the system, may have a value fortuitously close to its null value. To detect this

possibility, it might be required to consider all possible assortments of parameters in the fitting procedure, rather than hierarchically as described above, a generally more arduous process. Two possibilities must be considered: (1) the sequential addition of the high order parameter does not have a statistically significant effect on the fit, and (2) the sequential addition of the high order parameter does have a statistically significant effect on the fit. In both cases, restriction of the lower order term to its null value with the high order term allowed to vary cannot possibly significantly improve the fit over that obtained with both parameters allowed to vary. In the second case, if the submodel with the lower order term allowed to vary and the high order term constrained to its null value has the same statistical significance as the submodel with the high order term allowed to vary and the lower order term constrained to its null value, is there any reason to choose one of these submodels over the other? Consideration of the physical significance of the reformulated parameters provides a resolution to this dilemma. Higher order terms represent the modulation of lower order terms. For a high order interaction to exist physically requires all the lower order interactions which it modulates to exist physically. For example, for a second order (pairwise) interaction to exist requires both first order interactions to exist (i.e., both ligands must be able to bind to the protein before interaction can occur between binding events). Similarly for a third order interaction to exist requires all second order interactions which it modulates to exist. In the first case cited above, it would be inconsistent with the physical properties of the reformulated parameters not to conclude that the submodel with the lower order interaction was the most appropriate. In the second case, it can simply be concluded that the term for the third order interaction is statistically significant. A subsequent fit with the second order term constrained to its null value can be performed to determine if the second order term is not statistically different than its null value.

Based upon these arguments the logical course for the application of reformulated models to the statistical analysis of data is to sequentially add parameters in a hierarchical fashion. If third and higher order terms are statistically significant as determined by *F* and/or K–S tests then lower order terms can be

Table 1

Reformulated Adair model for hemoglobin oxygen binding for both ΔG° s and $\Delta \epsilon$ s

	T ^a	→	TO ^b	→	TO ₂	→	TO ₃	→	TO ₄
Stepwise reformulated ΔG° s		$o - RT \ln 4$		$o + oo - RT \ln 3/2$		$o + 2oo + ooo - RT \ln 2/3$		$o + 3oo + 3ooo + oooo - RT \ln 1/4$	
Reformulated ΔG° s of assembly ^c	0		$o - RT \ln 4$		$2o + oo - RT \ln 6$		$3o + 3oo + ooo - RT \ln 4$		$4o + 6oo + 4ooo + ooooo$
Stepwise reformulated $\Delta \epsilon$ s		$\Delta \epsilon_o$		$\Delta \epsilon_o + \Delta \epsilon_{oo}$		$\Delta \epsilon_o + 2\Delta \epsilon_{oo} + \Delta \epsilon_{ooo}$		$\Delta \epsilon_o + 3\Delta \epsilon_{oo} + 3\Delta \epsilon_{ooo} + \Delta \epsilon_{oooo}$	
Reformulated total ϵ of species	ϵ_t		$\epsilon_t + \Delta \epsilon_o$		$\epsilon_t + 2\Delta \epsilon_o + \Delta \epsilon_{oo}$		$\epsilon_t + 3\Delta \epsilon_o + 3\Delta \epsilon_{oo} + \Delta \epsilon_{ooo}$		$\epsilon_t + 4\Delta \epsilon_o + 6\Delta \epsilon_{oo} + 4\Delta \epsilon_{ooo} + \Delta \epsilon_{oooo}$

^aT is used to denote tetrameric hemoglobin.^bO is used to denote an oxygen molecule, O₂.^cFrom free components of T and O.

constrained to determine if they are indistinguishable from their null values if this appears warranted or useful.

3. Materials and methods

3.1. Hemoglobin oxygen binding data

The hemoglobin oxygen binding data set BISTRIS of Imai was used for the analysis presented here [16]. This data set was obtained by monitoring the absorbance of 1.2 mM adult human hemoglobin A as a function of oxygen concentration at 620 nm, 25°C, and pH 7.4 in 0.1 M Bis–Tris. The relatively high hemoglobin concentration of 1.2 mM was used to minimize the effects of dimer–tetramer equilibria on the observed oxygenation curves. Analogous results with somewhat different parameter values were obtained with the PHOS data set of Imai [16], which was obtained using phosphate as a buffer in place of Bis–Tris. The results from the analysis of the PHOS data set are not reported here. The pO_2 values of Imai, given in mm Hg (Torr), were converted to units of atm using the relation (1 atm = 760 mm Hg) for the analysis presented below.

3.2. Reformulation of the Adair model

The Adair model [21] is the basic model for most analyses of hemoglobin oxygen binding data. This model makes no distinction between the α and β subunits or of the geometric arrangement of subunits. The reformulated form of the Adair model has been derived previously [4] and is given in Table 1. Note that the reformulated Adair model has a single reformulated ΔG° (oo) describing the interaction between oxygen binding to different subunits, and is therefore equivalent to a tetrahedral model of subunit interactions.

3.3. Reformulation of extinction coefficients for the Adair model

The concept of reformulation can be extended to any state (path independent) parameter by logical analysis using previously described principles [4]. The Adair model for oxygen binding to hemoglobin is formally described by the four Adair binding

constants. In practice, hemoglobin oxygen binding is monitored by some physical method, such as the absorbance change associated with hemoglobin oxygen binding in the present case. To date, all analyses of spectroscopically monitored hemoglobin oxygen binding data by fitting with the Adair model have been based upon the implicit assumption that each bound oxygen produces the same spectroscopic effect. However, several studies suggest that this is in fact not the case [22–24]. In order to address the possibility that in the sequential binding of oxygen to hemoglobin each subsequent oxygen binding event may not produce the same spectroscopic change requires reformulation of the extinction coefficients associated with oxygen binding, and inclusion of these reformulated parameters in the stepwise regression procedure. Unlike reformulated ΔG° s where multiple ways for forming a given complex require statistical correction factors, once a complex is formed there are not multiple ways it can exhibit a spectroscopic change and therefore there is no need for statistical correction factors for reformulated ϵ s. The observed absorbance of free hemoglobin must also be described by a parameter. Free tetrameric hemoglobin is described by its extinction coefficient ϵ_t , a zero order term³. The spectroscopic change associated with the first oxygen binding is assigned a value of $\Delta\epsilon_o$, a first order term. The second oxygen binding is described by the same spectroscopic change, $\Delta\epsilon_o$, but possibly including some perturbation described by the term $\Delta\epsilon_{oo}$, a second order spectroscopic interaction term. The third oxygen binding is described by the base term $\Delta\epsilon_o$, plus two second order interactions with already bound oxygens, $2 \times \Delta\epsilon_{oo}$, plus the third order interaction term $\Delta\epsilon_{ooo}$. The spectroscopic effect of the fourth oxygen binding is described by the base term $\Delta\epsilon_o$, plus three second order interactions with already bound oxygens, $3 \times \Delta\epsilon_{oo}$, plus three third order interactions with pairs of already bound oxygens, $3 \times \Delta\epsilon_{ooo}$, plus the fourth order interaction term $\Delta\epsilon_{oooo}$. The complete reformulated Adair model for both ΔG° s and $\Delta\epsilon$ s is given in Table 1. The pattern of terms describing the reformulated $\Delta\epsilon$ s is the same as that for the reformulated ΔG° s except for the lack of statistical correction factors and the requirement of a zero order term for the absorbance of free hemoglobin.

3.4. Reformulated square model for hemoglobin oxygen binding data

The reformulated version of the Adair model implies equivalence of interaction between all of the subunits. It is of interest to test the ability of a slightly more complicated model to fit the data. The known structure of hemoglobin [9] allows for a number of possible subunit–subunit interactions, and the reformulated version of this system with a distinction between α and β subunits (the microscopic model) has been presented previously (Fig. 7, Ref. [4]). Elimination of the distinction between α and β subunits in the microscopic model results in the reformulated square model for hemoglobin oxygen binding, with a distinction between adjacent ($\alpha\beta$ and $\beta\alpha$) and diagonal ($\alpha\alpha$ and $\beta\beta$) interactions, and the reformulated square model is shown in Table 2. The reformulated ΔG° s are o —the intrinsic affinity of a subunit for an oxygen molecule, oo —the effect of an oxygen bound to a subunit on the oxygen affinity of the directly adjacent (not diagonally adjacent) subunits, oxo —the effect of an oxygen bound to a subunit on the oxygen affinity of the diagonally adjacent subunit, and ooo and $oooo$ —the

corresponding third and fourth order interaction terms. A square model for hemoglobin oxygen binding with pairwise interactions between adjacent subunits was originally proposed by Pauling which required only two thermodynamic terms to account for the observed oxygen binding properties of hemoglobin [2]; K' —the intrinsic affinity of a binding site for an oxygen molecule, and α —the effect of an oxygen bound to a site on the affinity of adjacent sites. The terms o and oo in the reformulated square model are the ΔG° equivalents to the equilibrium constant terms K' and α used by Pauling to describe his square model. The spectroscopic constants associated with the square model are also reformulated and given in Table 2.

3.5. Data analysis

Data were analyzed by fitting with the appropriate mathematical models by derivative free nonlinear regression using the IBM PC based version of the BMDP program AR [19]. No weighting of data points was used in this analysis [25]. The BMDP file used for fitting data with the reformulated Adair model is given in Table 3.

Table 2

Reformulated square model for hemoglobin oxygen binding for both ΔG° s and $\Delta \epsilon$ s^a

<p>The diagram shows the sequential binding of oxygen to hemoglobin. It starts with state T (empty square). Binding to one site leads to TO (one 'O' in a square). From TO, two paths exist: one to TOXO (one 'O' in a square, one 'O' in a separate square) with term $o+oxo-RT\ln 1/2$, and another to TOO (two 'O's in a square) with term $o+oo$. From TOXO, binding to the second site leads to TOOO (two 'O's in a square, one 'O' in a separate square) with term $o+2oo+ooo-RT\ln 2$. From TOO, binding to the second site leads to TOOO with term $o+oo+oxo+ooo$. Finally, binding to the third site leads to TOOOO (three 'O's in a square, one 'O' in a separate square) with term $o+2oo+oxo+3ooo+oooo-RT\ln 1/4$.</p>		
Species	Reformulated ΔG° s of assembly	Reformulated total ϵ s of species
T	0	ϵ_t
TO	$o-RT\ln 4$	$\epsilon_t + \Delta\epsilon_o$
TOO	$2o + oo-RT\ln 4$	$\epsilon_t + 2\Delta\epsilon_o + \epsilon_{\Delta_{oo}}$
TOXO	$2o + oxo-RT\ln 2$	$\epsilon_t + 2\Delta\epsilon_o + \Delta\epsilon_{oxo}$
TOOO	$3o + 2oo + oxo + ooo-RT\ln 4$	$\epsilon_t + 3\Delta\epsilon_o + 2\Delta\epsilon_{oo} + \Delta\epsilon_{ooo}$
TOOOO	$4o + 4oo + 2oxo + 4ooo + oooo$	$\epsilon_t + 4\Delta\epsilon_o + 4\Delta\epsilon_{oo} + 2\Delta\epsilon_{oxo} + 4\Delta\epsilon_{ooo} + \Delta\epsilon_{oooo}$

^aSee footnotes to Table 1.

Table 3

BMDP file used to fit Imai's hemoglobin oxygen binding data with the reformulated Adair model (O2E3 submodel)

```

/PROBLEM TITLE IS 'FIT WITH REFORMULATED ADAIR HB OXYGEN BINDING MODEL'.
/INPUT
  VAR = 2.
  FORMAT IS FREE.
  FILE = 'HB.DAT'.
/VAR
  NAMES ARE PO, ABS.
/REGRESS
  DEPENDENT = ABS.
  PARAM = 9.
  ITER = 200.
  HALF = 8.
/PARAMETER
  NAMES =          O,          OO,          OOO,          OOOO,          E,          EO,          EOO,          EOOO,          EOOOO.
  INIT =          10583.9,      -4374.0,          0,          0,          619.3,      -170.3,          0,          0.0,          0.0.
  FIXED =          0,          0,          OOO,          OOOO,          0,          0,          0,          EOOO,          EOOOO.
  MAX =          2E4,          1E4,          4E4,          5E4,          1E4,          1E4,          1E4,          1E4,          1E4.
  MIN =          -1E4,          -1E4,          -1E4,          -5E5,          -1E4,          -1E4,          -1E4,          -1E4,          -1E4.
/FUNCTION
  T = 298.15.
  R = 8.314.
  RT = R*T.
  KTO = EXP(-O/RT)*4.
  KTOO = EXP(-(2*O + OO/RT)*6.
  KTOOO = EXP(-(3*O + 3*OO + OOO)/RT)*4.
  KTOOOO = EXP(-(4*O + 6*OO + 4*OOO + OOOO)/RT).
  ETO = E + EO.
  ETOO = E + 2*EO + EOO.
  ETOOO = E + 3*EO + 3*EOO + EOOO.
  ETOOOO = E + 4*EO + 6*EOO + 4*EOOO + EOOOO.
  TT = 1.2E - 3.
  TF = TT/(1 + KTO*PO + KTOO*PO**2 + KTOOO*PO**3 + KTOOOO*PO**4).
  TO = TF*KTO*PO.
  TOO = TF*KTOO*PO**2.
  TOOO = TF*KTOOO*PO**3.
  TOOOO = TF*KTOOOO*PO**4.
  F = E*TF + TO*ETO + TOO*ETOO + TOOO*ETOOO + TOOOO*ETOOOO.
/END

```

3.6. Analysis procedure for hemoglobin oxygen binding data

The analysis begins by fitting hemoglobin oxygen binding data with the minimal submodel of the Adair model described by the base terms o , ϵ_1 , and $\Delta\epsilon_o$. This submodel is designated as O1E2 where the one in O1 denotes the number of reformulated ΔG° terms and the two in E2 denotes the number of reformulated ϵ terms allowed to vary in the fit. All parameters not allowed to vary in the fitting process are fixed at their null values, which are all zero in

the present case. The fit is then repeated with the value of $\Delta\epsilon_{oo}$ also allowed to vary (submodel O1E3). This process is continued by the sequential addition of reformulated ΔG° s and $\Delta\epsilon$ s. Each fit provides a value for the residual sum of squares (RSS) and estimates and standard errors for the fit parameters. The results of this procedure are summarized in Table 4 for the reformulated Adair model and Table 5 for the reformulated square model. At the bottom of Table 4 is shown one submodel designated O3-1E2, where o and ooo are adjustable parameters and oo is fixed at its null value. Note that for the square

Table 4

Fit of hemoglobin oxygen binding data from Imai [16] with an Adair model with both ΔG° s and ϵ s reformulated^a

Submodel	RSS	AC	P (AC < 2)	ΔG° (kJ mol ⁻¹) (S.E.)				ϵ (M ⁻¹ cm ⁻¹) (S.E.)				
				o	oo	ooo	$oooo$	ϵ_t	$\Delta \epsilon_o$	$\Delta \epsilon_{oo}$	$\Delta \epsilon_{ooo}$	$\Delta \epsilon_{oooo}$
O1E2	1.72e-1	0.04	< 0.01	-13.0(0.4)				610(30)	-168(6)			
O1E3	8.14e-2	0.06	< 0.01	-9.5(0.3)				580(10)	-430(20)	210(20)		
O1E4	1.95e-2	0.10	< 0.01	-12.7(0.2)				380(30)	340(50)	-780(50)	740(40)	
O1E5	9.45e-3	0.14	< 0.01	-14.5(0.2)				800(100)	-900(200)	2400(500)	-5400(800)	1.0e4(1e3)
O2E2	4.74e-4	0.16	< 0.01	-5.87(0.08)	-4.36(0.06)			485.7(0.8)	-121.9(0.3)			
O2E3	2.29e-5	-1.47	0.03	-5.85(0.02)	-4.26(0.01)			491.3(0.3)	-211(2)	59(2)		
O2E4	2.25e-5	1.51	0.03	-5.78(0.07)	-4.30(0.05)			491.6(0.4)	-220(10)	80(20)	-20(20)	
O2E5	1.59e-5	2.06	0.59	-5.85(0.06)	-4.22(0.04)			489.8(0.5)	-150(20)	-350(90)	1100(200)	-2200(500)
O3E2	4.95e-5	0.79	< 0.01	-8.3(0.1)	1.1(0.3)	-5.9(0.4)		494.5(0.6)	-123.2(0.1)			
O3E3	2.26e-5	1.48	0.03	-5.7(0.2)	-4.6(0.4)	0.4(0.5)		491.0(0.5)	-214(6)	61(4)		
O3E4	1.53e-5	2.13	0.69	-3.4(0.4)	-8.4(0.7)	4.0(0.7)		490.0(0.4)	-450(60)	400(90)	-250(60)	
O3E5	—			—	—	—		—	—	—	—	—
O4E2	4.69e-5	0.81	< 0.01	-7.9(0.2)	-1.3(0.5)	31.5(0.6)	-137 (na)	493(1)	-122.7(0.2)			
O4E3	—			—	—	—	—	—	—	—	—	—
O4E4	—			—	—	—	—	—	—	—	—	—
O4E5	—			—	—	—	—	—	—	—	—	—
O3-1E2	6.17e-5	0.66	< 0.01	-7.95(0.03)		-4.63(0.03)		492.8(0.3)	-122.9(0.1)			

^aFor a given submodel adjustable parameters have either a value and standard error (S.E.) or a '—' in the appropriate column, whereas spaces are left for parameters fixed to their null values. Cases where the fitting procedure was unable to provide an S.E. are denoted by a 'na' (not available) and cases where the fitting procedure either failed completely or gave physically unrealistic parameter values are denoted by a '—'.

Table 5

Fit of hemoglobin oxygen binding data from Imai [16] with a square model with both ΔG° 's and ϵ 's reformulated^a

Submodel	RSS	AC	P (AC < 2)	ΔG° (kJ mol ⁻¹) (S.E.)					ϵ (M ⁻¹ cm ⁻¹) (S.E.)					
				<i>o</i>	<i>oo</i>	<i>oxo</i>	<i>ooo</i>	<i>oooo</i>	ϵ_t	$\Delta \epsilon_o$	$\Delta \epsilon_{oo}$	$\Delta \epsilon_{oxo}$	$\Delta \epsilon_{ooo}$	$\Delta \epsilon_{oooo}$
O1E2	1.72e-1	0.04	< 0.01	-13.0(0.4)					610(30)	-168(6)				
O1E3.1	9.71e-2	0.06	< 0.01	-16.6(0.5)					60(200)	500(200)	-500(100)			
O1E3.2	9.71e-2	0.06	< 0.01	-16.6(0.5)					60(200)	500(200)		-1e4(2e2)		
O1E4	—			—					—	—	—	—		
O1E5	—			—					—	—	—	—	—	
O1E6	—			—					—	—	—	—	—	—
O2.1E2	4.60e-4	0.16	< 0.01	-5.4(0.1)	-7.0(0.1)				485.2(0.8)	-121.7(0.3)				
O2.1E3.1	2.51e-5	1.35	0.01	-5.39(0.02)	-6.82(0.02)				490.6(0.3)	-220(3)	98(3)			
O2.1E3.2	2.17e-5	1.54	0.04	-5.38(0.02)	-6.84(0.02)				490.2(0.2)	-202(2)		160(5)		
O2.1E4	1.72e-5	1.91	0.37	-5.37(0.02)	-6.80(0.02)				489.1(0.4)	-150(20)	-300(80)	700(100)		
O2.1E5	1.67e-5	1.98	0.47	-5.30(0.07)	-6.88(0.07)				489.5(0.5)	-170(20)	-250(90)	600(100)	-20(20)	
O2.1E6	1.66e-5	1.98	0.47	-5.30(0.07)	-6.88(0.07)				489.5(0.5)	-170(20)	-230(90)	-260(na)	3e2(1e2)	-7e2(3e2)
O2.2E2	—			—		—			—	—				
O2.2E3.1	—			—		—			—	—	—			
O2.2E3.2	—			—		—			—	—		—		
O2.2E4	—			—		—			—	—	—	—	—	
O2.2E5	—			—		—			—	—	—	—	—	
O2.1E6	—			—		—			—	—	—	—	—	—
O3E2	—			—	—	—			—	—				
O3E3.1	—			—	—	—			—	—	—			
O3E3.2	2.14e-5	1.56	0.05	-5.6(0.3)	-6(1)	-1(2)			490.6(0.7)	-204(4)		160(7)		
O3E4	—			—	—	—			—	—	—	—		
O3E5	—			—	—	—			—	—	—	—		
O3E6	—			—	—	—			—	—	—	—	—	—
O4E2	—			—	—	—			—	—	—	—		
O4E3.1	—			—	—	—	—		—	—	—	—		
O4E3.2	2.14e-5	1.56	0.05	-5.5(0.4)	-6(3)	-2(9)0.2(2)			490(1)	-210(30)		170(60)		
O4E4	—			—	—	—	—		—	—	—	—		
O4E5	—			—	—	—	—		—	—	—	—	—	
O4E6	—			—	—	—	—		—	—	—	—	—	—
O5E2	—			—	—	—	—	—	—	—	—	—		
O5E3.1	—			—	—	—	—	—	—	—	—	—		
O5E3.2	—			—	—	—	—	—	—	—	—	—		
O5E4	—			—	—	—	—	—	—	—	—	—	—	
O5E5	—			—	—	—	—	—	—	—	—	—	—	
O5E6	—			—	—	—	—	—	—	—	—	—	—	—

^aSee footnote to Table 4.

model there are two possible second order interaction terms, oo and oxo , and two possible second order spectroscopic interaction terms, $\Delta\epsilon_{oo}$ and $\Delta\epsilon_{oxo}$, either or both of which in either or both cases could be included in the fit. These possibilities are distinguished by using the designation O2.1 to denote the submodel where o and oo are the adjustable parameters, O2.2 to denote the submodel where o and oxo are adjustable parameters, and O3 to denote the submodel where o , oo , and oxo are adjustable parameters. Similarly, E3.1 denotes the submodel where ϵ_t , $\Delta\epsilon_o$, and $\Delta\epsilon_{oo}$ are adjustable parameters, E3.2 the submodel where ϵ_t , $\Delta\epsilon_o$, and $\Delta\epsilon_{oxo}$ are adjustable parameters, and E4 the submodel where ϵ_t , $\Delta\epsilon_o$, $\Delta\epsilon_{oo}$, and $\Delta\epsilon_{oxo}$ are adjustable parameters.

3.7. Autocorrelation (AC) analysis of residuals from fits

The RSS value obtained from a fit is one measure of goodness of fit, with a smaller RSS corresponding to a better fit. It is possible that the model which

gives the lowest RSS may in fact represent a poor fit to the data. For data which are poorly fit by a model, the residuals (observed/predicted values) from the fit will generally show a systematic sequential deviation from zero, whereas if the data are well fit by the model the residuals should show a random deviation from zero. Examination of the residuals for such systematic deviations is a common subjective test for goodness of fit. A quantitative method for evaluating the sequential nonrandomness of the residuals about zero can be performed by calculating the autocorrelation of the residuals and testing the resulting value for significance as described in standard texts [18]. The autocorrelation is calculated using the formula

$$AC = \Sigma(r_{i+1} - r_i)^2 / \Sigma(r_i - \bar{r})^2 \quad (3)$$

where r_i are the residuals, and \bar{r} is the average of the residuals which should be very close to zero unless the fitting procedure has failed. For random sequential variation, a value of 2 is expected for the autocorrelation. Values less than 2 signify that sequential residuals are more likely to be of similar

Table 6

Summary of F and K–S tests of significance of incrementally added parameters for the reformulated Adair model

A. F – test P values based on RSS values ^a							
	o		$+oo$		$+ooo$		$+oooo$
$\epsilon_t + \Delta\epsilon_o$	1.72e – 01 4.2e – 03	5.0e – 52	4.74e – 04 7.4e – 21	1.4e – 13	4.95e – 05 3.5e – 03	4.3e – 01	4.69e – 05
$+ \Delta\epsilon_{oo}$	8.14e – 02 6.1e – 07	9.4e – 76	2.29e – 05 4.8e – 01	4.8e – 01	2.26e – 05 8.9e – 02		–
$+ \Delta\epsilon_{ooo}$	1.95e – 02 6.2e – 03	2.9e – 59	2.25e – 05 1.2e – 01	9.2e – 02	1.53e – 05		–
$+ \Delta\epsilon_{oooo}$	9.45e – 03	7.4e – 54	1.59e – 05		–		–
B. K–S P values based on residuals ^b							
$\epsilon_t + \Delta\epsilon_o$	O1E2 0.39	< 0.0001	O2E2 0.005	0.010	O3E2 0.17	0.49	O4E2
$+ \Delta\epsilon_{oo}$	O1E3 0.028	< 0.0001	O2E3 0.5	0.5	O3E3 0.08		–
$+ \Delta\epsilon_{ooo}$	O1E4 0.23	< 0.0001	O2E4 0.17	0.17	O3E4		–
$+ \Delta\epsilon_{oooo}$	O1E5	< 0.0001	O3E5		–		–

^aThe roman font numbers in this section are the RSS values from the Adair submodels from Table 4 laid out as a function of the incrementally added reformulated ΔG° s across the top and of the incrementally added reformulated $\Delta\epsilon$ s down the side. The italicized numbers are F -test probability (P) values that the RSS for the submodel with an added parameter is less than that without by chance (one-sided test).

^bThe larger case entries in this section are the Adair submodels from Table 4 laid out as function of the incrementally added parameters. The smaller case numbers are (one-sided) P values that the frequency distribution of the residuals for the submodel with an added parameter is different than that without by chance using the K–S test.

sign whereas values greater than 2 signify that sequential residuals are more likely to be of opposite sign. For the application described here, a poor fit is expected to give a value for $AC < 2$ whereas a good fit is expected to give a value for $AC \cong 2$. Following the recommended procedure for sample sizes > 25 , the probability (P) that the determined value for AC is less than 2 by chance (one-sided test) was determined by the t -test. These calculations were performed in a Lotus 123 spreadsheet, and these results are summarized in Table 4 for the reformulated Adair model and Table 5 for the reformulated square model.

3.8. Statistical F -tests

The standard procedure to determine the statistical significance of an added parameter in a regression analysis is to use an F -test [18,20,26]. The probability (P) that the variance (RSS) with the added parameter is less than that without the added parameter (one-sided test) is due to chance is calculated by the F -test using the incomplete beta function [20]. These calculations were performed in a Lotus 123 spreadsheet with the following function definition to compare the fit from model 1 with RSS1 and $DF1$ with the fit from model 2 with RSS2 and $DF2$

Table 7

Summary of F and K–S tests of significance of incrementally added parameters for the reformulated Adair model^a

A. F -test P values based on RSS values.

	o		$+oo$	OR + oxo		$+oo \& oxo$		$+ooo$	$+oooo$
$\epsilon_t + \Delta \epsilon_o$	1.72e-1 2.1e-2	1.2e-52	4.60e- 1.1e-19	-		-		-	-
$+ \Delta \epsilon_{oo}$	9.71e-2 2.1e-2 ^b	1.4e-74	2.50e-5 4.3e-20	-	-	-	-	-	-
OR + $\Delta \epsilon_{oxo}$	9.71e-2	2.9e-78	2.17e-5 1.0e-1 ^d 2.1e-1 ^e	-	4.8e-1 ^f	2.14e-5	0.5	2.14e-5	-
$+ \Delta \epsilon_{oo} \& \Delta \epsilon_{oxo}$	-		1.72e-5 4.6e-1	-		-		-	-
$+ \Delta \epsilon_{ooo}$	-		1.67e-5 0.5	-		-		-	-
$+ \Delta \epsilon_{oooo}$	-		1.67e-5	-		-		-	-

B. K-S P values based on residuals.

	o		$+oo$	OR + oxo		$+oo \& oxo$		$+ooo$	$+oooo$
$\epsilon_t + \Delta \epsilon_o$	O1E2 0.39	< 0.0001	O2.1E2 0.0013	-		-		-	-
$+ \Delta \epsilon_{oo}$	O1E3.1 0.39 ^b	< 0.0001	O2.1E3.1 0.0028 ^c	-	-	-	-	-	-
OR + $\Delta \epsilon_{oxo}$	O1E3.2	< 0.0001	O2.1E3.2 0.08 ^d 0.23 ^e	-	0.5 ^f	O3E3.2	0.5	O4E3.2	-
$+ \Delta \epsilon_{oo} \& \Delta \epsilon_{oxo}$	-		O2.1E4 0.45	-		-		-	-
$+ \Delta \epsilon_{ooo}$	-		O2.1E5 0.5	-		-		-	-
$+ \Delta \epsilon_{oooo}$	-		O2.1E6	-		-		-	-

^aThis Table is analogous to Table 6 using the results summarized in Table 5. The layout of this Table is complicated by the fact that for both the reformulated ΔG° s and ϵ s there are two second order terms, oo & oxo and $\Delta \epsilon_{oo}$ & $\Delta \epsilon_{oxo}$, respectively, either or both of which in either or both cases could be incrementally included in the fitting procedure. All these possibilities are included in this Table, with the footnoted values further defined between the respective RSS values (top) or submodels (bottom) as follows: b; 1.72e-1 to 9.71e-2 (O1E2 to O1E3.2), c; 4.60e-4 to 2.17e-5 (O2.1E2 to O2.1E3.2), d; 2.50e-5 to 1.72e-5 (O2.1E3.1 to O2.1E4), e; 2.17e-5 to 1.72e-5 (O2.1E3.2 to O2.1E4), f; 2.17e-5 to 2.14e-5 (O2.1E3.2 to O3E3.2). All other P values in this Table apply to the immediately adjacent entries.

(DF = degrees of freedom = no. of data points – no. of fit parameters).

$$P = @BETAI(DF2/2, DF1/2, \\ DF2 * RSS2 / (DF2 * RSS2 + DF1 * RSS1)) \quad (4)$$

Using this method, the statistical significance of sequentially added parameters was evaluated and these results are summarized in Table 6A for the reformulated Adair model and Table 7A for the reformulated square model.

3.9. Kolmogorov–Smirnov (K – S) test

The F -test is dependent on the assumption that the errors (residuals) are normally distributed. This assumption is likely to be true when the experimental data are well fit by the model, and less likely to be true when the data are not well fit by the model. The nonparametric (i.e., normal distribution independent) equivalent to the F -test is the K – S test which determines the statistical significance of the difference between two distributions [18,20]. The vectors of residuals from two fits were compared using the K – S feature of BMDP routine 3S. This test is a two-sided test, and the P values obtained from this routine were divided by two to yield the corresponding one-sided P values for direct comparison with the one-sided F -test P values. The results for sequentially added parameters for the reformulated Adair and square models are summarized in Table 6B and Table 7B, respectively.

4. Results

4.1. Statistical analysis of hemoglobin oxygen binding data using the reformulated Adair model

The results of fitting hemoglobin oxygen binding data with the reformulated Adair model (Tables 1 and 3) are summarized in Table 4. Table 6A summarizes the results of statistical F -tests on each added parameters statistical significance and Table 6B the results of the K – S tests. Examination of Table 6 reveals that the F and K – S tests demonstrate similar patterns of significance for the sequentially added

parameters, although the F -tests provide much higher estimates for the significance of the added parameters for poorly fitting models than the K – S test, for example when considering O1E2 to O2E2 and O2E2 to O3E2 submodels. When comparing the effect of sequentially added parameters for well fit models both methods provide similar estimates for the significance of the added parameters, for example, O2E3 to O2E4, etc., as expected if the errors from the well fit models are (nearly) normally distributed. Although the K – S test is therefore the more rigorous method in general, it is also more difficult to perform than the F -test and the F -test may be used for applications of this type if undue emphasis is not placed on the high levels of significance resulting from performing this test on RSS values from non-normally distributed residuals.

The interpretation of the results from fitting this data with the reformulated Adair model begins in the upper left hand corner of Table 6A,B with the O1E2 submodel. The most significant improvement in fit upon addition of a parameter (oo or the $\Delta \epsilon_{oo}$) occurs upon addition of the oo term to arrive at the O2E2 submodel. From the O2E2 submodel, the most significant improvement in fit occurs by addition of the $\Delta \epsilon_{oo}$ term to give the O2E3 submodel, although the addition of the ooo term to give the O3E2 submodel is also statistically significant. From the O2E3 submodel, no more statistical significance can be extracted at the $P < 0.05$ level (K – S test) by subsequent parameters. The fit with the O3E2 submodel gives a positive value of 1.1 kJ mol^{-1} for the oo term, suggesting weak positive cooperativity between binding to adjacent subunits. To test the significance of the difference of oo from zero in this submodel an additional fit was performed using the O3E2 submodel with the value for oo fixed to zero, and the results are given in Table 6 (O3-1E2 submodel). This fit with the O3-1E2 submodel ($RSS = 6.173e - 5$) is not statistically worse than the fit with the O3E2 submodel ($RSS = 4.95e - 5$) ($P = 0.22$ (F -test), 0.45 (K – S test)), however, the O3-1E2 submodel is marginally statistically worse than the O2E3 submodel ($RSS = 2.29e - 5$) ($P = 0.001$ (F -test), 0.11 (K – S test)). We are therefore left with two submodels, with O2E3 better than O3E2, but not at a high level of statistical significance ($P = 0.003$ (F -test), 0.17 (K – S test)). The addition of parame-

ters beyond the minimal set necessary to explain the data results in dilution of the statistical significance of the fit parameters, and this is evident in the observed S.E.s for the fit parameters (Table 4) which begin increasing with additional parameters beyond the O2E3 submodel. The AC value for the residuals from the O2E3 submodel of 1.47 ($P = 0.03$) demonstrates that this submodel accounts for nearly all of the significant variation in the data. Together these results demonstrate that the O2E3 submodel is the best statistically for the description of the hemoglobin oxygen binding data of Imai using a reformulated Adair model, although it is not possible using this data to strictly rule out the O3E2 submodel based on the statistical analysis presented here.

4.2. Statistical analysis of hemoglobin oxygen binding data using other models for hemoglobin oxygen binding: the reformulated square model

When considering alternative models to the reformulated Adair model which might account for the spectroscopically observed data analyzed here, we must remain cognizant of the fact that although hemoglobin is composed of α and β subunits, the data analyzed here provides no information on binding to these different subunits. Several alternatives exist to the Adair model for hemoglobin oxygen binding within this constraint and given the known symmetry of tetrameric hemoglobin (reviewed in Ref. [9]). One is that tetrameric hemoglobin acts as a dimer of cooperative dimers. This possibility has previously been considered and ruled out on the basis of its failure to account for the known value of the Hill constant [2] and a reformulated version of this model has also been tested and found inadequate. The reformulated Adair model presumes equal interactions between all four subunits, which is in essence a tetrahedral model. In order to abandon this model for a more complex model, we must demonstrate at some level of statistical confidence that the four subunits do not interact equally by the analysis of appropriate data, or we must use other evidence to infer some other pattern of interaction between subunits. The next level of complexity for a model of hemoglobin oxygen binding is the reformulated square model.

The analysis procedure was repeated using a reformulated square model (Table 2). The results of this analysis are summarized in Table 5 and Table 7. From the upper left of Table 7 (O1E2 submodel), the most significant improvement in fit occurs by the addition of the oo term to arrive at the O2.1E2 submodel. From this submodel, it is possible either to add the $\Delta\epsilon_{oo}$ term to arrive at the O2.1E3.1 submodel, or the $\Delta\epsilon_{oxo}$ term to arrive at the O2.1E3.2 submodel, with the O2.1E3.2 submodel giving only a slightly better fit than the O2.1E3.1 submodel (RSSs of $2.17e-5$ vs. $2.50e-5$). Additional parameters do not yield any significant improvement in fit. From a physical point of view, it is expected that thermodynamic and spectroscopic interactions between adjacent sites would share a common structural mechanism, and for this reason the O2.1E3.1 submodel is preferred in the following discussion. The thermodynamic parameters are essentially identical between these two submodels (Table 5).

4.3. Comparison of fits obtained with the Adair and square models using the F -test

The fits obtained with the Adair and square models can be compared using a statistical F -test. Such a comparison can also be made with the K–S test but the F -test should be adequate in this instance since we are comparing two models which fit the data well. The Adair model is found to be not statistically significantly better than the square model ($P = 0.40$) using the RSS value from fit with the O2.1E3.1 square submodel. Therefore, the fits with reformulated Adair and square models are statistically indistinguishable.

5. Discussion

The analysis of data pertaining to complex physical systems is an essential part of the scientific study of such systems. A systems of N complexes (states) formulated using ΔG° s of assembly, or stepwise ΔG° s, requires $N \Delta G^\circ$ s. No reduction in the number of terms required to describe such a system is possible without eliminating some of the complexes from the model. Systems described by reformulated ΔG° s also require $N \Delta G^\circ$ s to describe the full potential

complexity of such a system [4], however, because higher order reformulated terms represent increasingly complex interactions between individual events, it is less and less likely that successively higher order terms will be physically and/or statistically significant, and models described by reformulated ΔG° s can often be described with some subset of these parameters. The utility of this approach for the theoretical analysis of cooperativity in a hypothetical dimeric protein, where higher order terms were eliminated based upon physical arguments, has been presented previously [7]. In the present analysis, this concept is demonstrated for the statistical analysis of data, where the statistical significance of successively higher order terms is determined by stepwise regression with reformulated parameters. This method is related to other statistical methods such as ANOVA and stepwise log–linear regression where interactions between variates are analyzed and their significance assessed [18], except that in the present case the data is continuous, the models are nonlinear, and the regressors are reformulated state parameters. Reformulation was originally developed for thermodynamic models described by ΔG° s, however, this concept can be applied to any state parameter as demonstrated here for extinction coefficients, and this method can therefore in principle be applied to any system which can be cast in terms of state parameters.

Hemoglobin oxygen binding represents the classic example of allosteric interaction in a multi-subunit protein, and has been subjected to a number of detailed analyses based upon the fitting of spectroscopically monitored hemoglobin oxygen binding data with the four Adair constants. Fitting this data with the unreformulated Adair model is equivalent to fitting the data with the O4E2 reformulated Adair submodel, with four thermodynamic terms allowed to vary in the fitting procedure. The analysis presented here demonstrates that at least two, and at most three, of the four reformulated Adair ΔG° s are statistically significant (Tables 4 and 6). Therefore, the debate concerning the statistical and physical significance of the four Adair constants (Refs. [14–17] and references therein) is a debate about parameters with reduced statistical significance. Although a number of studies have indicated that subsequent oxygen binding events may produce a different spec-

troscopic change than the initial event [22–24], such a possibility has never been included together with binding constants for the analysis of spectroscopically monitored hemoglobin oxygen binding data. In order to address this possibility, the concept of reformulation was extended to extinction coefficients (ϵ s), and these parameters included in the analysis. In preceding from the statistically significant O2E2 submodel (Table 6) addition of the $\Delta\epsilon_{oo}$ term (O2E3 submodel) was found to be the most statistically significant added parameter, however, the addition of the *ooo* term (O3E2 submodel) was also statistically significant, but to a lesser extent. A comparison of the O2E3 vs. the O3E2 submodels indicates that the O3E2 submodel cannot be strictly rejected on this basis, but the O2E3 submodel is preferred on both physical and statistical grounds.

The known structure for hemoglobin (reviewed in Ref. [9]) is consistent with a square model for hemoglobin if no distinction between oxygen binding to the α and β subunits can be made, or of interactions at the $\alpha\beta$ and $\beta\alpha$ interfaces between subunits. Stepwise regression with reformulated parameters was performed with the reformulated square model and the results are summarized in Table 5 and Table 7. The O2.1E3.2 submodel was best statistically, although the O2.1E3.1 submodel was not statistically significantly worse. It is reasonable to expect that spectroscopic interactions between subunits would parallel thermodynamic interactions, and the O2.1E3.1 model is preferred on this basis. The fit with a reformulated square model did not give a statistically better fit than the reformulated Adair model, and it cannot be concluded from these results alone that the reformulated square model is better than the reformulated Adair model, although a square model is a better representation of the known structure of hemoglobin than a tetrahedral model. It is of historical interest that Pauling proposed over 60 years ago that pairwise interactions between adjacent subunits in a square model could account for the oxygen binding properties of hemoglobin [2] and this is confirmed in the present analysis. Pauling obtained a value for his association equilibrium constant of interaction of $\alpha = 12$, which corresponds to the value $K_{oo} = \text{EXP}(oo/RT)$ (Table 5, submodel O2.1E3.1 or O2.1E3.2) = 15.7 reported here. Although the data used in the two studies was obtained under different

experimental conditions, the apparent interaction constants given above are clearly of the same magnitude. It would be of considerable interest to determine if a simple interaction term such as the β term proposed by Pauling could also account for the pH dependence of hemoglobin oxygen binding.

That first and second order thermodynamic interaction terms are statistically significant, and third and fourth order terms are not, must be interpreted carefully. The present analysis does not, nor in general cannot, prove that third and fourth order interaction terms do not exist. A Monod-Wyman-Changeux type of model proposing a sudden conformational change after two or three ligands have bound [27], which would be expected to demonstrate significant third and fourth order thermodynamic interaction terms given the quaternary nature of this mechanism, cannot therefore be ruled out based upon these observations. Although the possibility of quaternary interactions as the basis for hemoglobin cooperativity still exists, at some level interaction between two adjacent subunits must occur and we are left with the conclusion that second order interactions between adjacent subunits are both *necessary* and *sufficient* to account for spectroscopically monitored hemoglobin oxygen binding data. The finding of both second order thermodynamic and spectroscopic interactions between adjacent subunits presents a mechanistic model for hemoglobin oxygen binding where oxygen binding to a subunit causes a conformational change in that subunit and induces a conformational change in adjacent subunits. Such a model is most consistent with a sequential Koshland-Nemethy-Filmer type of model [28] with the modification that each subunit can exist in more than two conformational states; the fully deoxygenated conformational state, the fully oxygenated conformational state, and intermediate states depending upon the oxygenation status of the adjacent subunits, a possibility suggested by Haber and Koshland [29]. NMR studies on partially ligated hemoglobin provide information concerning the distribution of ligands and the conformational states of the subunits (reviewed in Ref. [10]) which also demonstrate such none two state behavior of subunits [30]. A theoretical analysis of cooperativity in a dimeric protein using reformulated thermodynamic models has also demonstrated that the parameter

space associated with a sequential type of model is the most consistent with the observation of both positive and negative cooperativity [7]. The statistical analysis presented here demonstrates that a model with simple conformational interactions between adjacent subunits is consistent with spectroscopically monitored hemoglobin oxygen binding data, and provides further evidence of none two state behavior during hemoglobin oxygen binding.

Acknowledgements

I would like to thank Dr. Kiyohiro Imai for publishing the raw data upon which the present analysis is based, and Dr. Michael Otto for his comments and suggestions concerning this manuscript.

References

- [1] G. Weber, *Biochemistry* 11 (1972) 864–878.
- [2] L. Pauling, *Proc. Natl. Acad. Sci. U.S.A.* 21 (1935) 186–191.
- [3] G. Weber, *Adv. Protein Chem.* 29 (1975) 1–83.
- [4] W.G. Gutheil, C.E. McKenna, *Biophys. Chem.* 45 (1992) 171–179.
- [5] W.G. Gutheil, *Biophys. Chem.* 52 (1994) 83–95.
- [6] G.D. Reinhart, *Arch. Biochem. Biophys.* 224 (1983) 389–401.
- [7] W.G. Gutheil, *Biophys. Chem.* 45 (1992) 181–191.
- [8] S.J. Edelstein, *Ann. Rev. Biochem.* 44 (1975) 209–232.
- [9] M.F. Perutz, *Q. Rev. Biophys.* 22 (1989) 139–236.
- [10] C. Ho, *Adv. Protein Chem.* 43 (1992) 153–312.
- [11] G.K. Ackers, M.L. Doyle, D. Myers, M.A. Daugherty, *Science* 255 (1992) 54–63.
- [12] G. Weber, *Proc. Natl. Acad. Sci. U.S.A.* 81 (1984) 7098–7102.
- [13] S.J. Edelstein, J. Edsall, *Proc. Natl. Acad. Sci. U.S.A.* 83 (1986) 3796–3800.
- [14] S.J. Gill, E. Di Cera, M.L. Doyle, G.A. Bishop, C.H. Robert, *Biochemistry* 26 (1987) 3995–4002.
- [15] E. Di Cera, C.H. Robert, S.J. Gill, *Biochemistry* 26 (1987) 4003–4008.
- [16] K. Imai, *Biophys. Chem.* 37 (1990) 197–210.
- [17] D. Myers, K. Imai, T. Yonetani, *Biophys. Chem.* 37 (1990) 323–340.
- [18] R.R. Sokal, F.J. Rohlf, *Biometry*, Freeman, New York, 1995.
- [19] W.J. Dixon, *BMDP Statistical Software Manual*, University of California Press, Berkeley, 1992.
- [20] W.H. Press, B.P. Flannery, S.A. Teukolsky, W.T. Vetterling, *Numerical Recipes*, Cambridge Univ. Press, New York, NY, 1986, pp. 155–190.

- [21] G.S. Adair, A.V. Bock, H. Field, *J. Biol. Chem.* 63 (1925) 529–545.
- [22] J. Rifkind, R. Lumry, *Fed. Proc., Fed. Am. Soc. Exp. Biol.* 26 (1967) 2325, Abstr.
- [23] M.L. Doyle, E. Di Cera, S.J. Gill, *Biochemistry* 27 (1988) 820–824.
- [24] T.M. Larsen, T.C. Mueser, L.J. Parkhurst, *Anal. Biochem.* 197 (1991) 231–246.
- [25] M.L. Doyle, G.K. Ackers, *Biophys. Chem.* 42 (1992) 271–281.
- [26] B. Mannervik, *Contemporary Enzyme Kinetics and Mechanism*, in: D.L. Purich (Ed.), Academic Press, New York, NY, 1983, pp. 75–95.
- [27] J. Monod, J. Wyman, J.P. Changeux, *J. Mol. Biol.* 12 (1965) 88–118.
- [28] D.E. Koshland Jr, G. Nemethy, D. Filmer, *Biochemistry* 5 (1966) 365–385.
- [29] J.E. Haber, D.E. Koshland, *Proc. Natl. Acad. Sci. U.S.A.* 58 (1967) 2087–2093.
- [30] G. Viggiano, C. Ho, *Proc. Natl. Acad. U.S.A.* 76 (1979) 3673–3677.