



## Similarity transference of molecular parameters. I. The atomic polar tensors of cyanoacetylene

B. B. Neto, M. N. Ramos, and R. E. Bruns

Citation: *The Journal of Chemical Physics* **85**, 4515 (1986); doi: 10.1063/1.451772

View online: <http://dx.doi.org/10.1063/1.451772>

View Table of Contents: <http://scitation.aip.org/content/aip/journal/jcp/85/8?ver=pdfcov>

Published by the [AIP Publishing](#)

---

### Articles you may be interested in

[Rotational spectrum of cyanoacetylene solvated with helium atoms](#)

*J. Chem. Phys.* **125**, 144310 (2006); 10.1063/1.2357604

[Similarity transference of molecular parameters. II. The bond distances, force constants and polar tensors of HC<sub>3</sub>N and HC<sub>5</sub>N](#)

*J. Chem. Phys.* **90**, 6933 (1989); 10.1063/1.456268

[Spintensor interaction in polarization transfer reactions](#)

*AIP Conf. Proc.* **69**, 662 (1981); 10.1063/1.32624

[Integrated infrared intensities and atomic polar tensors in fluoroform](#)

*J. Chem. Phys.* **73**, 5591 (1980); 10.1063/1.440079

[Infrared intensities: Polar tensors and charge flux parameters](#)

*J. Chem. Phys.* **69**, 4403 (1978); 10.1063/1.436430

---



# Similarity transference of molecular parameters. I. The atomic polar tensors of cyanoacetylene

B. B. Neto<sup>a)</sup> and M. N. Ramos

*Departamento de Química Fundamental, Universidade Federal de Pernambuco, 50000 Recife, PE, Brasil*

R. E. Bruns

*Instituto de Química, Universidade Estadual de Campinas, 13100 Campinas, SP, Brasil*

(Received 2 June 1986; accepted 11 July 1986)

A similarity transference procedure for the calculation of molecular parameters is proposed. The theoretical relationships between the direct transference method normally used and the proposed procedure are discussed. The importance of adequately defining the similarity models used in the transference calculations is emphasized. As an example, similarity models are constructed using the experimental and the STO-3G and 4-31G molecular orbital values of the atomic polar tensors of the HCN, C<sub>2</sub>H<sub>2</sub>, CH<sub>3</sub>CN, C<sub>4</sub>H<sub>2</sub>, CH<sub>3</sub>CCH, and C<sub>2</sub>N<sub>2</sub> molecules. Partial least squares calculations based on these similarity models and using the STO-3G and 4-31G values of the atomic polar tensors of HC<sub>3</sub>N result in estimates of the experimental tensors of cyanoacetylene which have about one-half or less the root mean square error of the molecular orbital values. Also the partial least squares detection of potentially unreliable estimates of polar tensor elements is illustrated.

## INTRODUCTION

The transference of localized parameters is a common technique used by chemists in attempting to estimate molecular properties when direct experimental observation is difficult. Several examples easily come to mind: isolated atom polarizabilities are often transferred to the molecular environment to calculate molecular polarizabilities<sup>1</sup>; bond energy tables, with average or representative values obtained from a series of molecules for which sufficient thermodynamic data are available, are used to calculate internal energy and enthalpy values<sup>2</sup>; standard bond distances and angles are transferred to predict geometries of molecules for which structural data are incomplete<sup>2</sup>; force constants of molecules with completely analyzed spectra are frequently transferred to molecules for which experimental observations are difficult or band assignments are uncertain.<sup>3</sup> In all such cases it is assumed that the parameters to be transferred are constant and insensitive to the differences between the electronic environments of the original molecules and those of the molecules to which the parameters are transferred. Eventual discrepancies are then explained on the basis of the chemist's knowledge of the differences in the electronic structures of the molecules involved.

For example, the limitations of direct transference procedures are striking for molecules exhibiting delocalization effects, such as conjugated carbon-carbon double bond systems. This is to be expected, since electronic and thermodynamic parameters evaluated using isolated double bonds seldom agree with the values obtained for the double bonds of conjugated systems, and the extent of the difference is not easy to predict.

In this work we apply a multivariate approach to the general transference of molecular parameters in which, in-

stead of assuming a parameter to have the same value for a group of similar molecules, we try to allow for its variation in a natural way, building a model of the similarity holding between the members of the group. The method we employ is based on the well-established principal component analysis,<sup>4</sup> where direct transference corresponds to a zero order principal component model, and variations in parameters for different molecular environments are described by increasing the number of principal components included in the model.

The similarity models used in the calculations are actually determined using the partial least squares (PLS) method introduced by Wold,<sup>5</sup> which is closely related to principal component regression. Linear models are defined which describe variations in atomic polar tensor parameters with differing electronic environments and also provide theoretical estimates of polar tensor elements for molecules where experimental intensity values are lacking.

As an example of the proposed procedure the atomic polar tensors of cyanoacetylene are calculated using experimental and quantum chemical polar tensors of several molecules containing carbon-carbon and carbon-nitrogen triple bonds. These values are then compared with those calculated from the experimental gas phase vibrational intensities and normal coordinates of this molecule.<sup>6</sup>

## SIMILARITY TRANSFERENCE OF POLAR TENSORS

The  $M$  elements of an atomic polar tensor can be regarded as the coordinates of a point in  $M$ -dimensional space. Depending on the situation  $M$  can go up to nine (all elements different and not null), but usually it is a much smaller value. To illustrate, seven diagonal ( $M = 3$ ) polar tensors are plotted in Fig. 1(a).

The points representing similar polar tensors naturally form a cluster in  $M$ -dimensional space. If the similarity is large, that is, if the points are clustered tightly enough, then

<sup>a)</sup> Present address: Instituto de Química, Unicamp, 13100 Campinas, SP, Brasil.

the whole set of polar tensors can be adequately represented by a linear *similarity model*,<sup>7</sup> much in the same way as the truncation of higher order terms in a Taylor series is not prejudicial to dipole moment evaluations for small displacements from the equilibrium geometry. In general, similarity transference of molecular parameters only assumes that the reference parameters are similar to those expected for the molecule whose properties are to be evaluated.

In the simplest of such models each polar tensor element is replaced by the average over the entire set of similar elements. This corresponds to a SIMCA (similarity modeling by class analogy)<sup>8</sup> model with zero principal components:

$$p_{k,\sigma\nu} = \bar{p}_{\sigma\nu} + e_{k,\sigma\nu}, \quad \sigma, \nu = x, y, \text{ or } z, \quad (1)$$

where  $p_{k,\sigma\nu} = \partial p_{\sigma} / \partial \nu$  for the  $k$ th polar tensor in the set, and  $\bar{p}_{\sigma\nu}$  is the average of all such elements:

$$\bar{p}_{\sigma\nu} = \frac{1}{n} \sum_k \frac{\partial p_{\sigma}}{\partial \nu}, \quad (2)$$

$n$  being the total number of similar polar tensors.

On using  $\bar{p}_{\sigma\nu}$  to represent  $p_{k,\sigma\nu}$  the residuals  $e_{k,\sigma\nu}$  are left, which contain both experimental and modeling errors, and are a measure of how well the model fits the experimental data. This zero order model corresponds to a single point in  $M$ -dimensional space, and the residuals can be used to define a confidence region around it, in which one expects to observe the points for other polar tensors which belong to the same similarity model.<sup>8</sup> In  $M$  dimensions this would be a hypersphere [a sphere in Fig. 1(a)].

Direct transference, which corresponds to our zero order model, thus assumes that the polar tensor for the atom being studied is equivalent to an average of polar tensors for atoms having similar electronic environments. Usually, due to scarcity of data, an atomic polar tensor from a selected reference molecule is chosen for transference, instead of an average. For example, the fluorine polar tensor of methyl fluoride has often been chosen to predict intensity values of other fluorine-containing molecules, such as SF<sub>6</sub>, UF<sub>6</sub>, and UF<sub>5</sub>.<sup>9</sup>

The zero order similarity model allows no parameter variation from one molecule to another. Of course it can be improved by the inclusion of principal components. With one component the SIMCA model becomes

$$p_{k,\sigma\nu} = \bar{p}_{\sigma\nu} + t_k b_{\sigma\nu} + e_{k,\sigma\nu}, \quad (3)$$

$$k = 1, 2, \dots, n, \quad \sigma, \nu = x, y, \text{ or } z.$$

Here the  $\sigma\nu$  element for the  $k$ th atom is represented by the average  $\bar{p}_{\sigma\nu}$  plus a product term,  $t_k b_{\sigma\nu}$ , calculated so as to minimize the new residuals  $e_{k,\sigma\nu}$ . In principal component-factor analysis terminology the  $b_{\sigma\nu}$  are the loadings of the different variables, the  $p_{\sigma\nu}$ , and determine the orientation of the principal component (PC) axis with respect to the initial coordinate system. Of all possible axes, the PC axis is the one that contains the maximum value of variance and so best represents the data, according to the least-squares criterion. The  $t_k$ , called scores, are the projections of the  $k$  polar tensor points on the PC axis. The more similar the actual atomic polar tensors, the more similar their scores. Based on the residuals one can draw a (hyper) cylinder as a region of confidence around the PC axis [Fig. 1(b)].

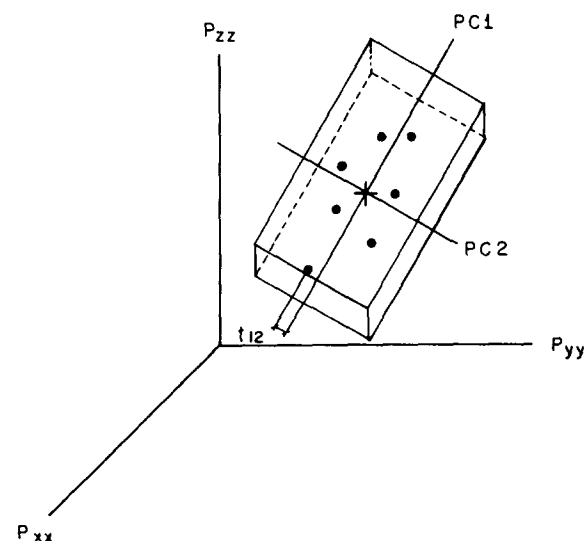
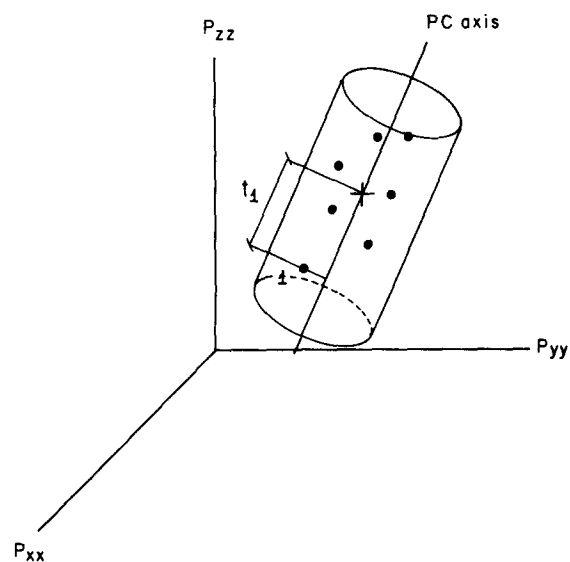
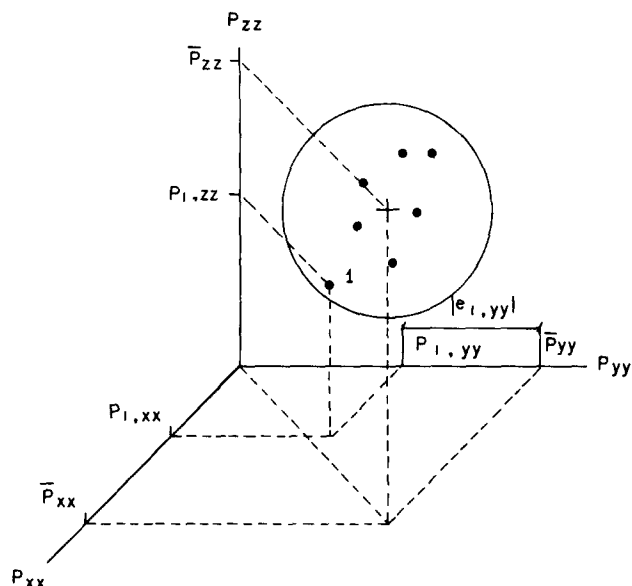


FIG. 1. Principal component models in three dimensions and respective confidence regions. (a) Zero components—a sphere; (b) one component—a cylinder; (c) two components—a box.

If a one component model is not sufficient to give an accurate representation of the data points a two-component model [Fig. 1(c)], in which each point is represented by two scores, can be used:

$$p_{k,\sigma\nu} = \bar{p}_{\sigma\nu} + t_{k1}b_{1,\sigma\nu} + t_{k2}b_{2,\sigma\nu} + e_{k,\sigma\nu},$$

$$k = 1, 2, \dots, n, \quad \sigma, \nu = x, y, \text{ or } z. \quad (4)$$

The second PC axis, whose orientation is given by  $b_{2,\sigma\nu}$ , contains the maximum amount of remaining variance, considering all possible axes orthogonal to the first one. The score of the  $k$ th polar tensor on this axis is given by  $t_{k2}$ , and the confidence region in this case is a (hyper)box [see Fig. 1(c)].

Usually the residuals  $e_{k,\sigma\nu}$  decrease as more components are added to the model, and additional PC's can be introduced until the data are represented adequately. The more complicated the data structure, the larger the number of PC's required to reduce the residuals to acceptable levels.

For the general case where the expansion is carried through to the  $A$ th principal component the equations become

$$p_{k,\sigma\nu} = \bar{p}_{\sigma\nu} + \sum_{a=1}^A t_{ka}b_{a,\sigma\nu} + e_{k,\sigma\nu},$$

$$k = 1, 2, \dots, n, \quad \sigma, \nu = x, y, \text{ or } z. \quad (5)$$

The scores  $t_{ka}$  of each set of polar tensors are taken as measures of the degree of similarity between the members of the set. The number of components kept in the model  $A$  can be determined statistically by means of the cross-validation method.<sup>10</sup> Good fitting models are obtained when the residuals contain little modeling error, i.e., when they reflect essentially experimental error.

The first (and major) step in PC analysis consists of deciding which atomic polar tensors are similar enough to be included in one model. Chemical arguments can be used, e.g., hydrogen polar tensors in one category or model and carbon tensors in another. More objectively, perhaps, one can resort to pattern recognition methods, such as the nearest neighbor rule (KNN),<sup>11</sup> hierarchical cluster techniques,<sup>12</sup> or the SIMCA method itself. The SIMCA method is especially suitable, since it provides a PC formalism for the different levels of approximation involved in the transference procedure. Confidence regions, given by hyperspheres, hypercylinders, or hyperboxes, can be constructed around the PC axes, and an atomic polar tensor is included in one model if its representative point falls within the respective confidence region.

Once similarity models are defined for the various types of atomic polar tensors, they can be used in polar tensor calculations for molecules not included in the model-building phase. PC regression<sup>4</sup> or the related PLS modeling technique<sup>5</sup> can then be used to make predictions from the similarity models.

If theoretical methods are reasonably successful in reflecting the similarities in electronic structures found by experimental observation, as is more or less the case with *ab initio* quantum chemical methods, they can be included in the model calculation. One then arrives at a model that connects experimental values with theoretical results, allowing

one to use the latter to make predictions about the former.

Two PC models are then constructed: one for the experimental polar tensors, given by Eq. (5), and the other for the corresponding theoretical values, represented by the analogous equation

$$\pi_{k,\sigma\nu} = \bar{\pi}_{\sigma\nu} + \sum_{a=1}^A \theta_{ka}\beta_{a,\sigma\nu} + e_{k,\sigma\nu},$$

$$k = 1, 2, \dots, n, \quad \sigma, \nu = x, y, z. \quad (6)$$

Here the Greek symbols correspond exactly to their Latin counterparts in Eq. (5), the only difference being that they are theoretical, rather than experimental, in origin. Both models—experimental and theoretical—are represented by dotted lines in Fig. 2, for a case in which one component suffices for an adequate description of the data.

The experimental and theoretical scores,  $t_{ka}$  and  $\theta_{ka}$ , respectively, can be related by the linear regressions

$$t_{ka} = \rho_a \theta_{ka} + \Delta_k, \quad a = 1, 2, \dots, A, \quad (7)$$

where  $\rho_a$  is the regression coefficient and  $\Delta_k$  measures the deviation of the  $k$ th polar tensor from the regression line. If the theoretical scores mirror reasonably well the similarities in the experimental scores, good regression lines will be obtained in Eq. (7).

This procedure, PC regression, ignores possible correlations between the experimental and theoretical blocks of values, which could improve the regressions in Eqs. (7). This is taken into account by the PLS technique, in which one block is weighted by the other and the regression is performed over latent variables, rather than over principal components. This amounts to relaxing the maximum variance criterion within one single block in order to improve the regression relating one block to the other. The latent variable axes (solid lines in Fig. 2) are tilted with respect to the principal component axes, but lead to better interblock relations.

To predict polar tensor values for atoms absent from the model but nevertheless similar to the ones that were actually employed, one starts by doing quantum chemical calculations to determine the theoretical polar tensor, whose elements are substituted into Eqs. (6) to yield the theoretical scores,  $\theta_{ka}$ . Experimental scores,  $t_{ka}$ , are then obtained from Eqs. (7), and in turn used in Eqs. (5) to predict "experimental" polar tensor elements. Of course, prior to modeling and prediction care must be taken to perform the rotations required to make all polar tensors consistent with one another. The whole procedure is depicted in Fig. 2 by the arrows connecting the three plots.

Misuse of the similarity modeling technique can also be detected with the PLS formalism. If theoretical values are not accurately described by a theoretical similarity model, then inaccurate predictions of experimental results may follow. This can be quantified by means of a statistical  $F$  test, as described below.<sup>13</sup>

A PLS prediction for the  $n'$ th polar tensor (not included in the model) will yield from Eq. (6) the residuals  $\epsilon_{n',\sigma\nu}$  which define for this tensor a residual variance with respect to the model:

$$S_{n'}^2 = \sum_{\sigma,\nu} \frac{\epsilon_{n',\sigma\nu}^2}{M - A}. \quad (8)$$

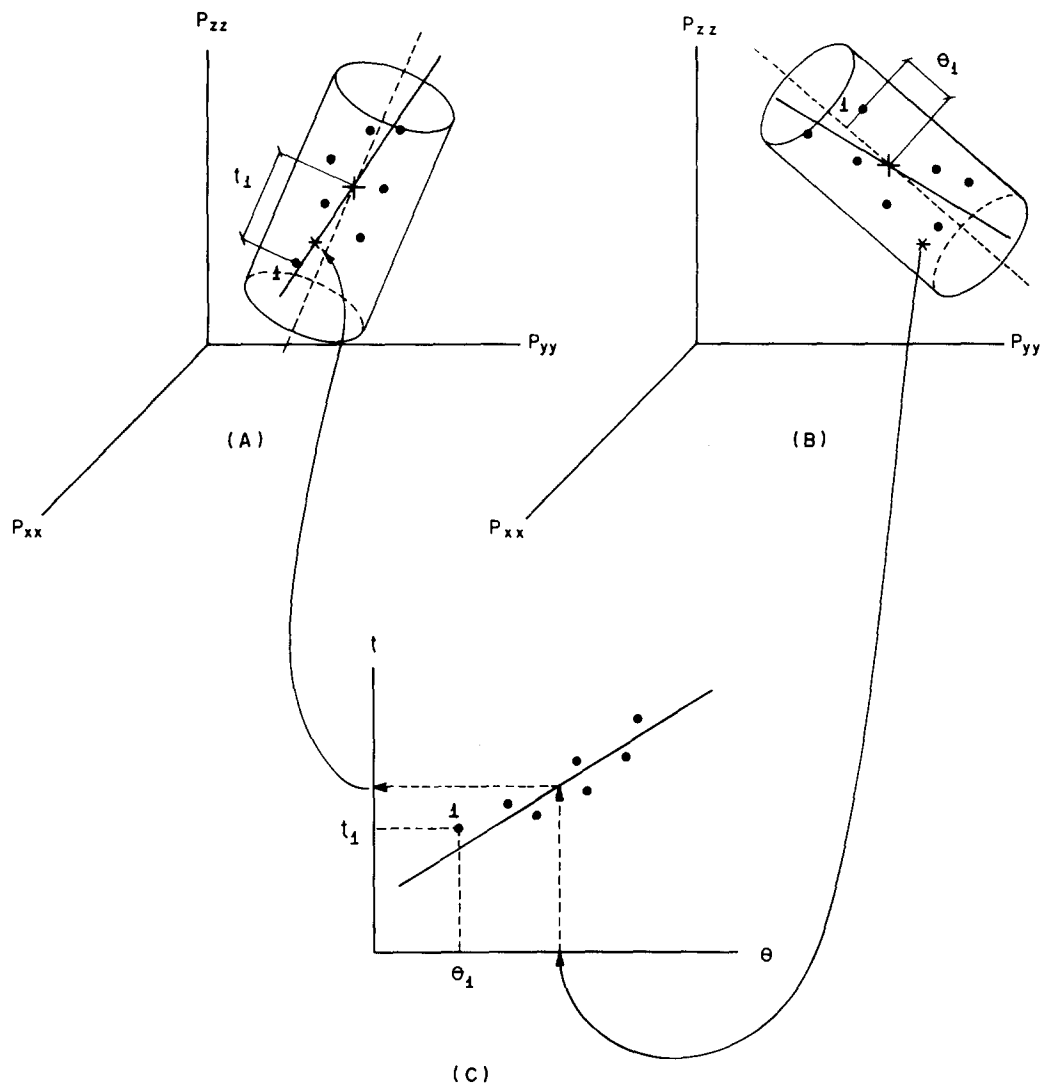


FIG. 2. One component PC models for diagonal polar tensors. (a) Experimental polar tensors; (b) theoretical polar tensors; (c) regression of the experimental scores on the theoretical scores.

Similarly, with the residuals of the tensors used for model building a "typical" variance can be defined

$$S_0^2 = \sum_k \sum_{\sigma, \nu} \frac{\epsilon_{k, \sigma \nu}^2}{(M-A)(n-A-1)}. \quad (9)$$

Assuming a reasonably normal distribution, both variances can be compared in an  $F$  test with  $(M-A)$  and  $(M-A)(n-A-1)$  degrees of freedom,

$$F = S_n^2 / S_0^2. \quad (10)$$

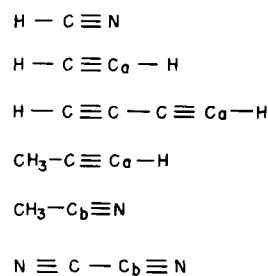
An  $F$  value showing that the residual variance of tensor  $n'$  is significantly larger than the variance of the model set is an indication that the predictions of the experimental polar tensor elements for atom  $n'$  are unreliable.

For the applications reported here and, at present, for other attempts to predict IR gas phase intensities using similarity modeling, not only is the  $F$  test scheme very approximate, but the calculated model itself is quite unstable and vulnerable to changes in the model-building set of polar tensors. This is the case because, for statistical purposes, IR intensities have been measured for very few molecules, and one is then forced to use half a dozen or less polar tensors to define the similarity model.

Although it is not possible to ascertain the stability of

the models calculated here, confidence in the predicted results is improved by the fact that several models result in similar values for the predicted polar tensors. Moreover, the alternative of direct transference, being a special case of PC modeling, suffers more acutely from the same shortcomings.

Reference:



Test: H - C<sub>a</sub>  $\equiv$  C - C<sub>b</sub>  $\equiv$  N

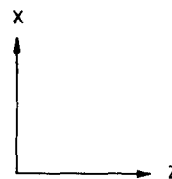


FIG. 3. Molecules used in the PLS prediction of the cyanoacetylene polar tensors.

TABLE I. Experimental and *ab initio* theoretical polar tensor reference values. Units of  $e$ ; coordinate system as in Fig. 3.

Molecule	Atom	Experimental		STO-3G		4-31G	
		$p_{xx} = p_{yy}$	$p_{zz}$	$p_{xx} = p_{yy}$	$p_{zz}$	$p_{xx} = p_{yy}$	$p_{zz}$
HCN <sup>a</sup>	H	0.237	0.218	0.208	0.393	0.321	0.282
	C	0.084	-0.292	0.042	-0.495	-0.032	-0.103
	N	-0.321	0.074	-0.250	0.102	-0.289	-0.179
C <sub>2</sub> H <sub>2</sub> <sup>b</sup>	H	0.205	0.183	0.115	0.358	0.253	0.216
	C <sub>α</sub>	-0.205	-0.183	-0.115	-0.358	-0.253	-0.216
C <sub>4</sub> H <sub>2</sub> <sup>c</sup>	H	0.200	0.238	0.122	0.469	0.257	0.281
	C <sub>α</sub>	-0.223	-0.394	-0.112	-0.276	-0.262	-0.403
CH <sub>3</sub> CCH <sup>d</sup>	H	0.211	0.180	0.119	0.372	0.256	0.203
	C <sub>α</sub>	-0.249	-0.465	-0.174	-0.480	-0.301	-0.781
CH <sub>3</sub> CN <sup>e</sup>	N	-0.331	-0.171	-0.287	0.018	-0.313	-0.464
	C <sub>β</sub>	0.060	0.078	0.118	-0.279	0.025	0.176
C <sub>2</sub> N <sub>2</sub> <sup>f</sup>	N	-0.245	0.123	-0.231	0.344	...	...
	C <sub>β</sub>	0.245	-0.123	0.231	-0.344	...	...

<sup>a</sup>K. Kim and W. T. King, *J. Chem. Phys.* **71**, 1967 (1979).

<sup>b</sup>K. Kim and W. T. King, *J. Mol. Struct.* **57**, 201 (1979).

<sup>c</sup>Calculated using the intensities reported in Th. Kooops, T. Visser, and W. M. A. Smit, *J. Mol. Struct.* **125**, 179 (1985).

<sup>d</sup>P. Jona, M. Gussoni, and G. Zerbi, *J. Phys. Chem.* **85**, 2210 (1981).

<sup>e</sup>Y. Koga, S. Kondo, S. Saeki, and W. B. Person, *J. Phys. Chem.* **88**, 3152 (1984).

<sup>f</sup>K. Kim and W. T. King, *J. Chem. Phys.* **80**, 974 (1984).

With increasing amounts of experimental data, similarity modeling of atomic polar tensors will rest on a more solid footing. Indeed, similarity modeling calculations can indicate which experimental measurements should be given highest priority to yield the most secure prediction of polar tensor values.

### AN APPLICATION: CYANOACETYLENE

Similarity models were constructed using the linear molecules shown in Fig. 3, the molecular axis being taken as the  $z$  axis. The reference (or training) set comprises the molecules employed in the model-building phase. The cyanoacetylene molecule, HC<sub>3</sub>N, for which predictions are made, constitutes the test set. The experimental polar tensors of the reference set molecules consistent with the Cartesian system of Fig. 3 are presented in Table I. Where necessary, the experimental polar tensors taken from the literature were rotated to conform to the orientations shown in Fig. 3.

Included in Table I are the *ab initio* STO-3G and 4-31G calculated polar tensor elements for all of the reference set molecules. The coordinate system is the same as the one shown in Fig. 3. The MO calculations were performed running the HONDO package<sup>14</sup> on a DEC PDP-10 mainframe, with standard basis sets as included in the computer program. It should be noted that the 4-31G values for the polar tensors of C<sub>2</sub>N<sub>2</sub> are not shown in Table I. The reason is that our 4-31G calculations for this molecule met with convergence difficulties (already experienced by King and co-workers<sup>15</sup>) that led to unacceptably large values for the dipole derivatives. For the STO-3G calculations no such troubles were encountered, and STO-3G values are therefore shown for the entire reference set. Experimental equilib-

rium geometries<sup>16</sup> were used throughout. The PLS calculations were carried out using the SIMCA-3B program<sup>17</sup> on a CPM 64 kbyte-8-bit microcomputer.

A critical stage in the similarity transference procedure is the definition of the similarity models to be used, which are characterized by the molecules included in the reference sets. Figure 3 contains all molecules having either a carbon-carbon or a carbon-nitrogen triple bond for which complete gas phase intensity data are available. Since the test molecule, cyanoacetylene, contains both types of triple bonds, the other molecules in Fig. 3 are qualified as members of the different categories or classes defining the various similarity models.

A class that suggests itself is that of the acid hydrogen atoms in HCN, C<sub>2</sub>H<sub>2</sub>, C<sub>4</sub>H<sub>2</sub>, and CH<sub>3</sub>CCH. All these molecules have a classical H-C≡ bond arrangement and one therefore anticipates similar electronic structures for these protons. These structures are not equivalent, of course, because of the different moieties in which the triple bond is embedded. On the other hand, since C-H bonds adjacent to triple bonds have protons with effective charges quite different from those of other protons,<sup>18</sup> other hydrogen polar tensors were not considered for inclusion in this class. Another group of similar atomic polar tensors may be defined by the nitrogen tensors from the HCN, CH<sub>3</sub>CN, and C<sub>2</sub>N<sub>2</sub> molecules, which display a -C≡N classical bond arrangement.

For the inner atoms the similarity models are a bit more difficult to define. Both the acid and nitrogen carbons have a -C≡ arrangement with differing groups participating in the single and triple bonds. For the carbon atom bonded to nitrogen, HCN, C<sub>2</sub>N<sub>2</sub>, and CH<sub>3</sub>CN provide reference atomic polar tensors. For the acid carbon the C<sub>2</sub>H<sub>2</sub>, HCN, C<sub>4</sub>H<sub>2</sub>,

and  $\text{CH}_3\text{CCH}$  molecules could form the reference class. However, separation of the carbon polar tensors into two different similarity models may not be warranted, since the carbon from HCN would be common to both models. This can be confirmed by inspection of Fig. 4, in which 4-31G, STO-3G, and experimental atomic polar tensors are represented graphically. In all three plots the hydrogen polar tensors of the reference class from a compact and well-separated cluster. For the STO-3G results three groupings can be distinguished, corresponding to the hydrogen, nitrogen, and carbon polar tensors. On the other hand, the 4-31G carbon and nitrogen tensors together form a large disperse group and probably only two similarity models should be formed, one for the hydrogen tensors and the other for the rest. In Table II the groups of polar tensors actually used to calculate similarity models for these molecules are summarized. One advantage of using a smaller number of models is evident from this table: more polar tensors can be used to calculate the parameters of each model, leading to statistically more stable solutions.

Once the similarity models are established the atomic polar tensors of cyanoacetylene can be calculated using the 4-31G and STO-3G molecular orbital estimates of the atomic polar tensors in Eqs. (5)–(7). In this way polar tensor values for all atoms but the central carbon of  $\text{HC}_3\text{N}$  are obtained. The remaining polar tensor is obtained from the charge conservation condition, i.e.,  $\sum_{\alpha} \mathbf{P}_X^{(\alpha)} = 0$ . The intensity values are proportional to the squares of the elements of the  $\mathbf{P}_Q$  tensor, which contains the dipole moment derivatives with respect to the normal coordinates. This tensor is given by<sup>19</sup>

$$\mathbf{P}_Q = \mathbf{P}_X \mathbf{AUL}' \quad (11)$$

where  $\mathbf{P}_X$  is an ordered juxtaposition of the atomic polar tensors and the matrix product  $\mathbf{AUL}'$  carries out the transformation from Cartesian to normal coordinates. The normal coordinates for  $\text{HC}_3\text{N}$  were calculated using the general valence force field of Uyemura *et al.*<sup>6</sup>

## RESULTS AND DISCUSSION

Table III contains the experimental and molecular orbital polar tensor values for cyanoacetylene along with the

partial least squares values using 2, 3, and 4 similarity models, as described in Table II. The root mean square errors for both the  $A_1$  and  $E$  symmetry species polar tensor elements, as well as the total root mean square errors, are also listed and can be used as a criterion of accuracy for the different types of calculations.

The total root mean square errors for all the PLS calculations are significantly less than the errors for the 4-31G and STO-3G molecular orbital results. Certainly the minimal basis set STO-3G approximation cannot be expected to provide a wave function which accurately reflects the electronic densities for small geometrical distortions from equilibrium, but the more sophisticated 4-31G basis set gives even poorer results, with an rms error of 0.199  $e$  which is larger than the error of 0.167  $e$  for the STO-3G calculations. King and co-workers<sup>15</sup> have previously commented on the deficiencies of the 4-31G wave functions in calculating electron densities for the  $\text{C}\equiv\text{N}$  bond. In our case this even led to the exclusion of 4-31G values for  $\text{C}_2\text{N}_2$ , as already mentioned.

On the other hand, the PLS estimates, which have significantly smaller rms errors, are based on crude similarity models, because of the dearth of experimental intensity data for appropriate reference molecules. For example, the 4-31G value of  $\partial p_z / \partial z_{C_b}$  for cyanoacetylene fits poorly any of the (4-31G)-based PLS models, and should lead to correspondingly poor predictions. This also occurs for the PLS calculation in which three similarity models were constructed for the STO-3G theoretical polar tensor elements. In spite of these deficiencies, the PLS rms errors are (with the exception of the 4-31G model with four classes) at most one-half of the errors encountered with the direct MO calculations.

More sophisticated wave functions can certainly be expected to provide results with lower errors than the ones obtained for the 4-31G and STO-3G calculations. Such wave functions, however, should also give a more faithful description of the similarities of the molecules involved, thus leading to improved similarity models and lower PLS errors.

The STO-3G theoretical values combined with the reference set experimental values result in better PLS models than do the 4-31G values. The best STO-3G PLS results occur for the calculations in which four similarity models,

TABLE II. Groups of atomic polar tensors used to construct similarity models. Entries in parentheses refer to the absence of  $\text{C}_2\text{N}_2$  from 4-31G models (see the text).

Number of models	Number of tensors in each model	Description of model	Origin of polar tensors
4	4	H	HCN, $\text{C}_2\text{H}_2$ , $\text{C}_4\text{H}_2$ , $\text{CH}_3\text{CCH}$
	3	$C_a$	$\text{C}_2\text{H}_2$ , $\text{C}_4\text{H}_2$ , $\text{CH}_3\text{CCH}$
	3(2)	$C_b$	HCN, $\text{CH}_3\text{CN}$ , ( $\text{C}_2\text{N}_2$ )
	3(2)	N	HCN, $\text{CH}_3\text{CN}$ , ( $\text{C}_2\text{N}_2$ )
3	4	H	HCN, $\text{C}_2\text{H}_2$ , $\text{C}_4\text{H}_2$ , $\text{CH}_3\text{CCH}$
	6(5)	$C_a + C_b$	HCN, $\text{C}_2\text{H}_2$ , $\text{CH}_3\text{CN}$ , $\text{C}_4\text{H}_2$ , $\text{CH}_3\text{CCH}$ , ( $\text{C}_2\text{N}_2$ )
	3(2)	N	HCN, $\text{CH}_3\text{CN}$ , ( $\text{C}_2\text{N}_2$ )
2	4	H	HCN, $\text{C}_2\text{H}_2$ , $\text{C}_4\text{H}_2$ , $\text{CH}_3\text{CCH}$
	9(7)	$C_a + C_b + \text{N}$	HCN, $\text{C}_2\text{H}_2$ , $\text{CH}_3\text{CN}$ , $\text{C}_4\text{H}_2$ , $\text{CH}_3\text{CCH}$ , ( $\text{C}_2\text{N}_2$ )

TABLE III. Experimental, molecular orbital, and partial least squares values for the polar tensors of cyanoacetylene ( $e$ ).

Symmetry species	Atom	Exptl <sup>a</sup>	4-31G				STO-3G			
			MO	PLS <sup>b</sup>	PLS <sup>c</sup>	PLS <sup>d</sup>	MO	PLS <sup>b</sup>	PLS <sup>c</sup>	PLS <sup>d</sup>
$(P_{zz}^A)$	H	0.267	0.306	0.263	0.263	0.263	0.516	0.261	0.261	0.261
	C <sub>a</sub>	-0.168	-0.150	-0.157	-0.189	0.076 <sup>e</sup>	-0.135	-0.078	-0.090	-0.163
	C <sub>b</sub>	0.168	0.583	0.146 <sup>e</sup>	0.187 <sup>e</sup>	-0.107 <sup>e</sup>	-0.050	0.155	0.149 <sup>e</sup>	0.152
	N	-0.085	-0.454	-0.240	-0.049	-0.049	0.183	0.017	-0.005	-0.005
	rms error <sup>f</sup>		0.278	0.078	0.023	0.185	0.214	0.068	0.056	0.041
$(P_{xx}^E = P_{yy})$	H	0.277	0.235	0.183	0.183	0.183	0.119	0.195	0.195	0.195
	C <sub>a</sub>	-0.175	-0.125	-0.075	-0.067	-0.152 <sup>e</sup>	-0.055	-0.091	-0.224	-0.184
	C <sub>b</sub>	0.217	0.216	0.312 <sup>e</sup>	0.298 <sup>e</sup>	0.072 <sup>e</sup>	0.167	0.178	0.034 <sup>e</sup>	0.230
	N	-0.260	-0.307	-0.294	-0.326	-0.326	-0.266	-0.321	-0.303	-0.303
	rms error <sup>f</sup>		0.040	0.085	0.088	0.093	0.102	0.069	0.105	0.047
	Total rms error <sup>f</sup>		0.199	0.082	0.065	0.146	0.167	0.069	0.085	0.044

<sup>a</sup>Data from Ref. 21.

<sup>b</sup>Two similarity models used, as described in Table II.

<sup>c</sup>Three similarity models used.

<sup>d</sup>Four similarity models used.

<sup>e</sup>The PLS residual values indicate that these are unreliable estimates.

<sup>f</sup>(rms error)<sup>2</sup> = (1/N)Σ<sub>i</sub> [(P<sub>ii</sub>)<sub>exptl</sub> - (P<sub>ii</sub>)<sub>calcd</sub>]<sup>2</sup>.

one for each atom of cyanoacetylene (except the central atom), are used. The total rms error in this case, 0.044  $e$ , is about one-fourth of the error resulting from direct use of the STO-3G polar tensor values.

In all the partial least squares calculations the hydrogens formed a class by themselves, as suggested by the plots in Fig. 4. For  $\partial p_z/\partial z_H$  the almost identical PLS values of 0.263 and 0.261  $e$  are obtained from the two kinds of similarity models, whereas direct 4-31G and STO-3G calculations yield quite different values, 0.306 and 0.516  $e$ , respectively. A similar situation exists for the  $\partial p_x/\partial x_H$  values: the PLS values, 0.183 and 0.195  $e$ , are very close, and the corresponding MO values, 0.235 and 0.119  $e$ , are not. The use of the reference set experimental values thus appears to have a stabilizing effect on the PLS calculated results.

The results for the  $\partial p_z/\partial z_N$  element illustrate the ability of the partial least squares method to provide accurate results even when poorly estimated molecular orbital values

are occasionally used. In this case the 4-31G and STO-3G values of  $-0.454$  and  $+0.183$   $e$  differ by more than 0.6  $e$  and predict opposite signs for the C $\equiv$ N stretching derivative. Such discrepancies have been observed for this derivative in other molecules, the best known example being HCN, for which only configuration interaction calculations were capable of resolving the  $\partial p/\partial r_{CN}$  sign ambiguity.<sup>20</sup> Regardless of whether 4-31G or STO-3G values are used for modeling, however, the PLS results are all close to zero, with the exception of the  $-0.240$   $e$  value obtained when C<sub>a</sub>, C<sub>b</sub>, and N were all included into a very diffuse 4-31G model. These results are in excellent agreement with the small negative experimental value, which predicts a  $\leftarrow \overset{+}{C} \equiv \overset{-}{N} \rightarrow$  change in charge distribution upon stretching. For the bending coordinate,  $\partial p_x/\partial x_N$ , the PLS, 4-31G and STO-3G values are very similar and in good agreement with the experimental result.

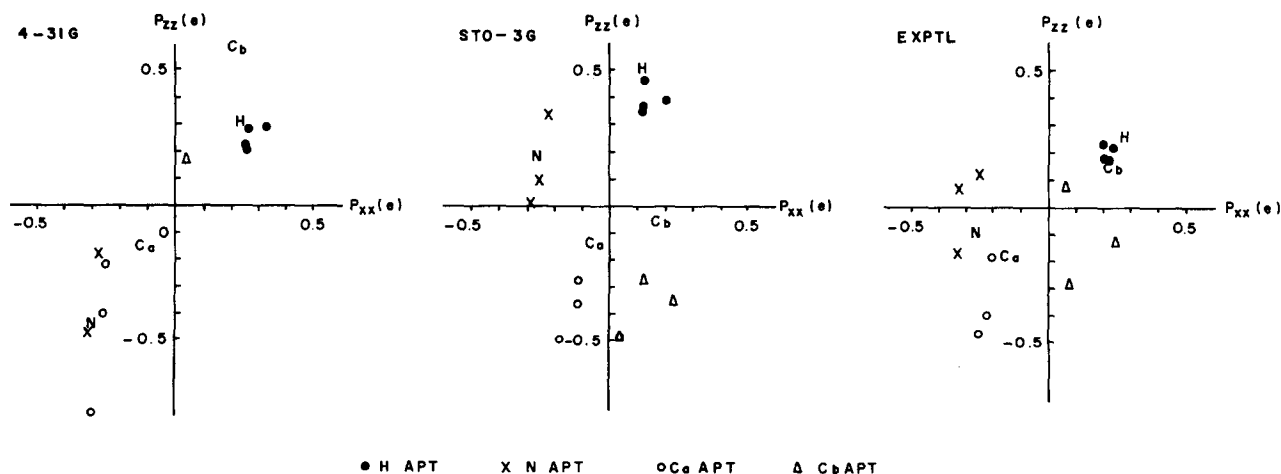


FIG. 4. Graphical representations of polar tensor elements of cyanoacetylene. (a) 4-31G values; (b) STO-3G values; (c) experimental values. The Latin capitals correspond to HC<sub>3</sub>N values.



TABLE IV. Experimental and calculated vibrational intensities of cyanoacetylene ( $\text{km mol}^{-1}$ ).

Symmetry species	$i$	$\nu_i$ ( $\text{cm}^{-1}$ )	$A_i$ (exptl)	$A_i$ (4-31G)	$A_i$ (STO-3G)	$A_i$ ( $9s5p1d/3s2p$ ) <sup>a</sup>	$A_i$ (PLS) <sup>b</sup>	$A_i$ (PLS) <sup>b</sup>
$A_1$	1	3327	60.4	81.6	225.4	104.8	61.4	55.7
	2	2272	9.92	59.2	11.0	63.8	8.2	12.0
	3	2077	1.93	0.1	50.2	0.2	4.6	10.2
	4	863	0.06	0.0	0.1	0.0	0.6	0.7
$E$	5	664	68.4	53.4	14.4	...	41.3	36.5
	6	499	8.0	4.1	1.4	...	1.7	2.1
	7	223	0.18	1.0	2.2	...	3.5	1.6
			rms error <sup>c</sup>	21.1	68.2	34.9	10.7	12.8

<sup>a</sup> Value from Ref. 21.

<sup>b</sup> Based on STO-3G results using four and two similarity models, respectively.

<sup>c</sup>  $(\text{rms error})^2 = (1/N) \sum_i [(A_i)_{\text{exptl}} - (A_i)_{\text{calcd}}]^2$ .

For the  $\partial p_z / \partial z_{C_2}$  and  $\partial p_x / \partial x_{C_2}$  derivatives all the PLS and molecular orbital values have negative signs, with the exception of the 4-31G PLS estimate indicated to be unreliable by its large residual variance. These values are in reasonable agreement with the  $-0.168$  and  $-0.175 e$  experimental elements.

Comparisons are more limited for the values of the  $C_b$  polar tensor. For  $\partial p_z / \partial z_{C_b}$  only two of the PLS results based on the STO-3G MO values are statistically reliable. These values are  $0.152$  and  $0.155 e$ , in excellent agreement with the  $0.168 e$  experimental value. For  $\partial p_x / \partial x_{C_b}$  again only two reliable PLS values are calculated:  $0.178$  and  $0.230 e$ , using STO-3G values with 2 and 4 similarity models, which compare well with the experimental value of  $0.217 e$ .

In Table IV vibrational intensities calculated for cyanoacetylene using the polar tensors in Table III and Eq. (11) are presented. Intensities corresponding to polar tensor sets for which the PLS method indicated at least one value to be unreliable have not been included in this table.

The root mean square errors for the PLS estimates of the intensities are about half the size of the errors resulting from direct calculation with 4-31G polar tensors. Although the most accurate PLS models were obtained with STO-3G polar tensors, the STO-3G wave function itself leads to poorer intensity values for  $\text{HC}_3\text{N}$  than does the 4-31G wave function. The 4-31G intensity values are in good agreement with the experimental intensities, except for the  $\nu_2$  band, for which the 4-31G result,  $59.2 \text{ km mol}^{-1}$ , is six times the  $9.92 \text{ km mol}^{-1}$  experimental value. Since this band is assigned to the CN stretching normal mode, this discrepancy is easily attributed to the too large absolute magnitudes calculated for  $\partial p_z / \partial z_{\text{N}}$  and  $\partial p_z / \partial z_{C_b}$  of cyanoacetylene. All the other polar tensor elements calculated from the 4-31G wave function are in excellent agreement with the experimental values, leading to the accuracy observed for the remaining 4-31G intensities.

Included for comparison are intensity values for the  $A_1$  symmetry species calculated from a  $9s5p1d/3s2p$  wave function.<sup>21</sup> The  $\nu_2$  intensity is also badly in error in this case. On the other hand, the PLS estimates of this intensity value are extremely good, differing by only  $2 \text{ km mol}^{-1}$  ( $\sim 20\%$ ) from the experimental value. The relative values of the other

PLS estimated intensities are also very encouraging, although the  $\nu_5$  intensity is calculated to be about one-half as strong as the  $\nu_1$  band, whereas in practice  $A_5$  is a little higher than  $A_1$ .

## FINAL CONSIDERATIONS

Direct transference of atomic polar tensors can be considered as a special case of the transference procedure based on similarity model calculations. This similarity transference procedure can be extended to other molecular parameters such as force constants, bond lengths and angles, etc. Especially intriguing is the possibility of employing different types of parameters (for example, force constants and polar tensors or electro-optical parameters) in the same calculation for simultaneously predicting fundamentally different but related properties (the frequencies and intensities of the fundamental bands of the gas phase infrared spectra).

The PLS calculations are a valuable complement to the molecular orbital calculations of molecular parameters. As a consequence, a library of theoretical and experimental parameters is created and can be continually updated. Indeed the modeling procedure can be used to establish priorities about parameter measurements necessary for more precise predictions. As the molecular orbital estimates of the parameters improve one can also expect more accurate PLS determinations. Since the PLS method is a variable reduction technique, less data are necessary for reliable estimation than in the case of multiple regression, which could also be used to relate experimental and theoretical reference data. With sufficient data, parametric statistical approaches could be included in the similarity transference method.

## ACKNOWLEDGMENTS

The authors gratefully acknowledge partial financial support from the Conselho Nacional de Desenvolvimento Científico e Tecnológico (M. N. R. and B. B. N.), the Fundação de Amparo à Pesquisa do Estado de São Paulo (R. E. B.) and the Financiadora de Estudos e Projetos (R. E. B.).

<sup>1</sup>J. Applequist, *Acc. Chem. Res.* **10**, 79 (1977); J. Applequist and C. O. Quicksall, *J. Chem. Phys.* **66**, 3455 (1977).

- <sup>2</sup>L. Pauling, *The Nature of the Chemical Bond*, 3rd ed. (Cornell University, Ithaca, 1960).
- <sup>3</sup>I. M. Mills, in *Infrared Spectroscopy and Molecular Structure*, edited by M. Davies (Elsevier, Amsterdam, 1963), Chap. 5.
- <sup>4</sup>K. V. Mardia, J. T. Kent, and J. M. Bibby, *Multivariate Analysis* (Academic, London, 1979), Chap. 8.
- <sup>5</sup>H. Wold, in *Systems Under Indirect Observation*, edited by K. G. Jöreskog and H. Wold (North-Holland, Amsterdam, 1982), Vol. II, Chap. 1.
- <sup>6</sup>(a) M. Uyemura and S. Maeda, *Bull. Chem. Soc. Jpn.* **47**, 2930 (1974); (b) M. Uyemura, S. Deguchi, Y. Nakada, and T. Onaka, *ibid.* **55**, 384 (1982).
- <sup>7</sup>S. Wold, C. Albano, W. J. Dunn III, U. Eudland, D. Esbensen, P. Geladi, S. Hellberg, E. Johansson, W. Lindberg, and M. Sjöström, in *Chemometrics: Mathematics and Statistics in Chemistry*, edited by B. R. Kowalski (Reidel, Dordrecht, 1984), Chap. 2.
- <sup>8</sup>S. Wold, *Pattern Recognition* **8**, 127 (1976).
- <sup>9</sup>W. B. Person and J. Overend, *J. Chem. Phys.* **66**, 1443 (1977); B. J. Krohn, W. B. Person, and J. Overend, *ibid.* **65**, 969 (1976).
- <sup>10</sup>S. Wold, *Technometrics* **20**, 397 (1978).
- <sup>11</sup>T. M. Cover and P. E. Hart, *IEEE Trans. Info. Theory* **IT-13**, 21 (1967).
- <sup>12</sup>W. S. Meisel, *Computer Oriented Approach to Pattern Recognition* (Academic, London, 1972).
- <sup>13</sup>M. Sjöström, S. Wold, W. Lindberg, J. A. Persson, and H. Martens, *Anal. Chim. Acta* **150**, 61 (1983).
- <sup>14</sup>M. Dupuis, J. Rys, and H. F. King, *Quantum Chemistry Program Exchange* **403** (1978).
- <sup>15</sup>K. Kim and W. T. King, *J. Chem. Phys.* **80**, 974 (1984); See also, Y. Koga, S. Kondo, S. Saeki, and W. B. Person, *J. Phys. Chem.* **88**, 3152 (1984).
- <sup>16</sup>I. Suzuki, M. A. Pariseau, and J. Overend, *J. Chem. Phys.* **44**, 3561 (1966); G. Strey and I. M. Mills, *J. Mol. Spectrosc.* **59**, 103 (1976); L. Halonen and I. M. Mills, *ibid.* **73**, 494 (1978); M. Tamimoto, K. Kuchitsu, and Y. Morino, *Bull. Chem. Soc. Jpn.* **44**, 386 (1971); J. L. Duncan, D. C. McKean, P. D. Mallinson, and R. McCulloch, *J. Mol. Spectrosc.* **46**, 232 (1973); M. Tamimoto, K. Kuchitsu, and Y. Morino, *Bull. Chem. Soc. Jpn.* **42**, 2519 (1969); for C<sub>2</sub>N<sub>2</sub>, L. E. Sutton, *Supplement to Tables of Interatomic Distances* (Chemical Society, London, 1963); C. C. Constain, *J. Chem. Phys.* **29**, 864 (1968).
- <sup>17</sup>Principal Data Components, 2505 Shepard Blvd., Columbia, MO; Se-panov AB, Ostrand 14, Enskede, Sweden.
- <sup>18</sup>W. T. King and G. B. Mast, *J. Phys. Chem.* **80**, 2521 (1976); R. E. Bruns and B. B. Neto, *J. Chem. Phys.* **68**, 847 (1978).
- <sup>19</sup>W. B. Person and J. H. Newton, *J. Chem. Phys.* **61**, 1040 (1974).
- <sup>20</sup>(a) B. Lin, K. M. Sando, C. S. North, H. B. Friedrich, and D. M. Chipman, *J. Chem. Phys.* **69**, 1425 (1978); (b) J. E. Gready, G. B. Backsay, and N. S. Hush, *Chem. Phys.* **31**, 467 (1978); (c) A. B. M. S. Bassi, *J. Chem. Phys.* **68**, 5667 (1978); (d) D. F. Hornig, *ibid.* **68**, 5668 (1978); (e) R. E. Bruns and W. B. Person, *ibid.* **53**, 1413 (1970).
- <sup>21</sup>M. N. Ramos, B. B. Neto, R. E. Bruns, and O. M. Herrera, *J. Mol. Struct.* **142**, 209 (1986).