in most practical situations, the first method is to be preferred.

## III. RESULTS AND DISCUSSION

The results of the third phase consist of two printouts. The first printout contains one page for every OMA peak (see Figure 4). This page contains each standard wavelength and subcategory that was found to match the given OMA peak. The second printout contains, for each subcategory, the number of OMA peaks that matched (see Figure 5).

Analysis of this information leads to conclusions regarding the identities of the species which are contained within a laser-induced plasma. Since the small dense plasma core is dominated by high-energy short-lived ionic species, and the region outside of the core is dominated by radical and excited-state species of transient existence, their identities cannot be determined by using the steady-state forms of spectroscopy. By use of the optical multichannel analyzer, one is able to acquire data over a wide wavelength range simultaneously with or delayed by a predetermined amount from the initiating laser pulse. The large number of species leads to a complex spectrum which precludes identification on a line by line basis. By use of our technique, however, most species may be identified by comparison with known spectroscopic constants of all potential species.

Although the technique presented here was described in terms of the analysis of relatively complex spectra obtained by means of an optical multichannel analyzer, it may readily be applied to the deconvolution of multicomponent spectra using more traditional forms of spectroscopy.

## REFERENCES AND NOTES

(1) C. P. Robinson, *Ann. N.Y. Acad. Sci.*, **267**, 81, (1976).
(2) A. M. Ronn, *Sci. Am.*, **240** (5), 114 (1979), and references cited therein.
(3) E. Fuss and T. P. Colter, *Appl. Phys.* **12**, 265 and references cited therein (1977).
(4) P. D. Maker, R. W. Terhune, and C. M. Savage, *Quantum Electron.*, *Proc. Int. Congr., 3rd, 1963*, 155 (1964).
(5) A. M. Ronn, *Chem. Phys. Lett.*, **42** (2), 207 (1976).
(6) S. T. Lin and A. M. Ronn, *Chem. Phys. Lett.*, **56** (3), 414 (1978).
(7) S. T. Lin and A. M. Ronn, *Chem. Phys. Lett.*, **49**, 255 (1977).
(8) Y. Langsam and A. M. Ronn, *Chem. Phys.*, in press.
(9) E. Yablonovitch and H. S. Kwok, *Opt. Commun.*, **18**, (1), 103 (1976).
(10) W. F. Meggers, C. H. Corliss, and B. F. Scribner, "Tables of Spectral Line Intensities Part I—Arranged by Elements", *NBS Monogr.* (U.S.) No. **145** (1975).
(11) W. L. Wiese, M. W. Smith, and B. M. Glennon, "Atom Transition Probabilities—Hydrogen Through Neon" *Natl. Stand. Ref. Data Ser.* (*U.S. Natl. Bur. Stand.*), **NSRDS-NBS 4** (1966).
(12) M. W. Smith and W. L. Wiese, *J. Phys. Chem. Ref. Data*, **2**, 85 (1973).
(13) P. H. Krupenie, "The Band Spectrum of Carbon Monoxide", *Natl. Stand. Ref. Data Ser.* (*U.S. Natl. Bur. Stand.*), **NSRDS-NBS 5** (1966).
(14) P. H. Krupenie, "The Spectrum of Molecular Oxygen", *J. Phys. Chem. Ref. Data*, **1** (2), 423 (1972).
(15) R. W. B. Pearse and A. G. Gaydon, "The Identification of Molecular Spectra", Wiley, New York, 1963.

# Theory of Correlation Tables. 1

T. VESZPRÉMI* and G. CSONKA

Department of Inorganic Chemistry, Budapest Technical University, 1521-H, Budapest, Hungary

A possible mathematical model is presented for correlation tables used in spectroscopy. A process based on information theory is demonstrated through an example for the optimum construction of correlation tables. In this example the construction of an ${}^1$H NMR correlation table is investigated; the method is general, can be used for other spectroscopies, and is suitable for the construction of correlation tables used in computerized evaluation of spectra.

The accumulation of spectroscopic experience led to the recognition that some fragments of molecules can be observed irrespective of their chemical environment. The structural elements cause absorption signals observable whenever the fragments are present in the molecules. These rules can be expressed in two simple forms, i.e., in correlation tables and in so-called additivity rules. Both are successfully used in practical spectroscopy.

The horizontal axis of a correlation table represents the observed characteristic (e.g., frequency), and in the vertical column different structural elements can be found. As correlation tables are used primarily for gaining immediate information, their setup largely depends on the researchers themselves.

Shift limits are not exact, and researchers' definitions of fragments vary as well. As often as not, conditions (e.g., solvent) essential for using the tables successfully are excluded.

Under the given circumstances, the use of correlation tables in computational technology is far less effective than it could be.

In what follows, a mathematical model of correlation tables is specified, its structural principles are explained, and a procedure for the optimum use of correlation tables is described. Although the results are of a general nature, simple correlation tables related to ${}^1$H NMR spectroscopy are taken as examples for better understanding.

## CORRELATION TABLES IN THE LITERATURE

As an introductory illustration some ${}^1$H NMR correlation tables have been selected. Table I contains data of some structural elements whose ${}^1$H NMR shifts can all be clearly interpreted.

The methoxy group lucidly exemplifies the nature of the problem. The first investigated correlation table[1] makes a distinction between methyl ester and methyl ether, although the given intervals greatly overlap. Sasaki and his co-workers[2] consider the aromatic and aliphatic methyl ethers entirely distinguishable on the basis of chemical shifts. According to their correlation table the chemical shifts of methyl esters and aromatic methoxy groups are indistinguishable. In the correlation table[3] given in the third column of Table I, the methoxy group is not subdivided. On the other hand, Bible[4] gives different chemical shifts for the three methoxy groups as he regards methoxy groups as distinguishable when found in three differing environments. According to the last column

**Table I.** Comparison of ¹H NMR Correlation Tables

| FRAGMENTS | CHEMICAL SHIFT LIMITS IN $\delta$ (ppm) | | | | |
|---|---|---|---|---|---|
| | BEECH[1] | SASAKI[2] | BIBLE[3] | SOHÁR[4] | KOLONITS[5] |
| CH₃-C≡C | 1.6-2.6 | 1.6-2.2 | 1.8-2.2 | - | - |
| CH₃ Ar | 1.9-2.8 | 2.0-2.8 | 2.1-2.8 | 2.2-2.8 | 2.33 |
| CH₃COO Ar | 1.9-2.6 | 1.8-2.5 | (2.3-2.7) | - | - |
| CH₃COO R | | | (2.0-2.5) | 2.0 | 1.95-2.05 |
| CH₃COR | 1.7-2.7 | | 2.0-2.8 | 2.1 | 2.1 |
| CH₃COAr | | | (2.5-2.9) | 2.6 | 2.62 |
| CH₃OCO | 3.5-4.0 | 3.5-4.1 | 3.2-4.3 | 3.8-4.0 | 3.71-4.10 |
| CH₃OAr | 3.2-4.2 | 3.5-4.1 | | 3.5-3.8 | 3.73 |
| CH₃OR | | 3.1-3.5 | | 3.3 | 3.3 |
| CH₃S | 1.8-2.8 | — | 2.0-2.6 | 2.1-2.6 | 2.06-2.35 |
| CH₃C | 0.1-2.2 | 0.5-1.8 | 0.7-1.9 | 0.9-1.9 | — |
| (CH₃)₃ C-O | | 1.0-1.4 | 1.1-1.4 | | 1.18-1.27 |
| -Ar | | 1.2-1.6 | 1.05-1.7 | | 1.25 |
| -CO | (0.1-2.2) | 1.0-1.6 | 1.05-1.7 | — | 1.08-1.15 |
| -O=C | | 0.9-1.6 | - | | - |
| -R | | 0.6-1.1 | 0.65-1.05 | | 0.9 |
| -C=C | | 0.9 1.6 | 1.05 1.7 | | 1.0 |
| HCO -R | 9.3-10.6 | 9.0-10.0 | 9.4-10.0 | 9.6-10.0 | |
| -Ar | | 9.0-10.2 | 9.6-10.5 | | — |
| -O | 7.7-8.6 | 7.8-8.4 | 8.0-8.3 | — | |
| HO-CO- | 10.0-13.5 | 6.0-10.0 | 9.5-13.0 | 10.4-10,8 | |

**Table II.** Chemical Shifts of Methylene Group

| FRAGMENTS | BEECH[1] | SASAKI[2] | SOHÁR[3] | SHOOLERY |
|---|---|---|---|---|
| RO CH₂OR | 4.5 - 6.1 | 4.2 - 5.0 | 4.3 - 5.0 | 4.85 |
| ArO CH₂OAr | | | | 6.05 |
| -CO-O-CH₂O-CO | | | | 6.65 |
| (R)S CH₂ Ar | 3.2 - 4.0 | | 3.3 - 4.1 | 3.65 |
| -CH₂(CO-OH)₂ | 3.2 - 3.5 | | | 2.65 |
| CH₂(CO-OR)₂ | | | | 3.25 |
| CH₂(CONR)₂ | 3.3 - 3.9 | 2.7 - 4.0 | 3.3 - 3.6 | 3.45 |
| CH₂(COR)₂ | | | | 3.65 |
| O CH₂ CO | 3.7 - 4.2 | 4.0 - 5.4 | 4.3 - 4.8 | 4.15 - 4.85 |

of Table I, the aliphatic and aromatic ethers are distinguishable, but aromatic ethers and esters overlap.[5]

Another example: the correlation tables are in comparative agreement concerning the chemical shift of the aromatic methyl group. On the basis of examining the chemical shifts of methyl groups on 172 different compounds containing methyl-substituted phenyl rings, it was found that the correct figure is 1.80–2.50 ppm, and this departs from all the other data. The data of Table II reveal that chemical shifts calculated by the Shoolery constants do not fall within the given chemical shift limits in every case.

## INFORMATIONAL CHANNEL IN SPECTROSCOPY

Figure 1 shows the information flow at the time of evaluating the spectra. Noise I of the figure marks the natural line width and the disturbance caused by the magnetic inhomo-
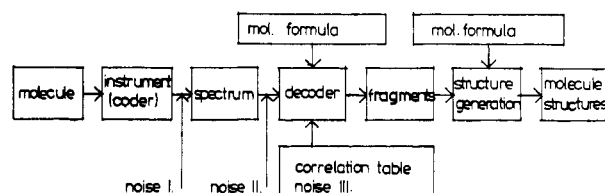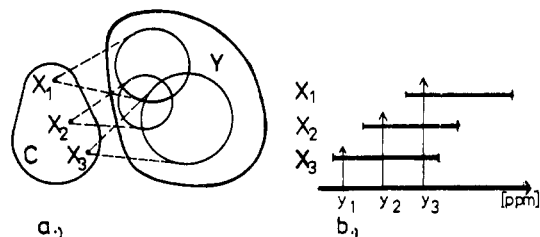


**Figure 1.** The information channel.



**Figure 2.** (a) Projection between the set of fragments (C) and signals (Y). (b) Application of measured frequencies for the selection of groups.

**Table III.** Effect of the Increasing Consideration of the Environment

| | -CH₃ | | | -CH₂- | | | -CH- | | |
|---|---|---|---|---|---|---|---|---|---|
| | — | $\alpha$ | $\beta$ | — | $\alpha$ | $\beta$ | — | $\alpha$ | $\beta$ |
| Band width (ppm) | 4.5 | ~1 | ~0.5 | 5 | ~0.8 | ~0.4 | 6.5 | ~0.8 | ~0.4 |
| Number of possible fragments | 1 | 20 | 400 | 1 | 400 | $16 \cdot 10^3$ | 1 | $8 \cdot 10^3$ | $64 \cdot 10^6$ |

geneity in case of ¹H NMR spectra. Noise II is the error in reading of the spectrum. The latter is usually larger. Noise II is 0.01 ppm; relevant literature defines ¹H NMR spectra to this accuracy.

Noise III is of greatest significance here because it is caused by the use of correlation tables for interpretation, and the consequent uncertainty is over an order of magnitude higher than the previous ones. Using correlation tables, the other two sources of noise cause negligible uncertainty.

If the influence of the setting of fragments is partially ignored, it may cause the widening of the given intervals in the correlation tables, because the signal of a less specifically defined group can be found within a wider interval.

A simplified correlation table can be seen in Figure 2. It can be noted that the given intervals of different fragments $(x_1, x_2, x_3)$ are overlapping. This is the reason it is impossible to say which group caused certain observed signals. Figure 2b shows such a case. It is apparent that if we observe the signal $y_3$, all of the three groups are possible.

If the fragments are made more specific, the observable signals will be concentrated in a smaller interval, but in this case the number of groups will greatly increase. This can be seen in some examples in Table III. The sign "—" in the table marks the case when the chemical shift limits are given irrespective of the environment, $\alpha$ marks the consideration of the substituents in the $\alpha$ position, and $\beta$ marks that of the substituents in the $\alpha$ and $\beta$ positions. Possible fragments have been calculated for 20 different substituents.

There are two consequences of the more comprehensive consideration of environment: (a) the correlation table requires more memory in the case of computer processing; (b) the number of compounds containing the fragment decreases. Our goal was to establish criteria for deciding which grouping would yield a correlation table giving maximum information. This can be done with the help of information theory.
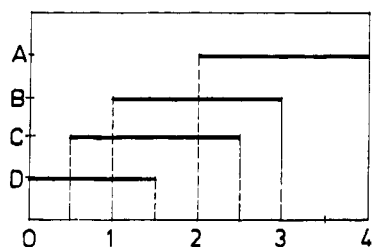
**Figure 3.** Model for simple correlation table.

**Table IV.** Processing of Correlation Tables

| Number of fragments (n) | Interval | Probability (p) | $H_v$ |
|---|---|---|---|
| 1 | 0-0. 5; 3-4 | $p_1$- 0.375 | 0 |
| 2 | 0. 5-1; 1.5-2 2.5-3 | $p_2$- 0.375 | 1 |
| 3 | 1. 0-1. 5; 2-2.5 | $p_3$- 0.25 | 1.585 |
| 4 | 0 | $p_4$- 0 | 2 |

Recently, information theory has been applied to quantitative and qualitative analysis, e.g., for the identification of unknown compounds with thin-layer chromatographic[6] and gas chromatographic retention values[7,8] and with binary coded mass spectra.[9] It has also been applied to retrieval with coded infrared[10,11] spectra.

Rotter and Varmuza[12] have discussed, from an information theoretical viewpoint, a variety of possible criteria for the evaluation of binary pattern classifiers. The procedure and its modified version have been applied to binary mass[13] and [13]C NMR spectra.[14,15]

Recently, Milne and his co-workers[16] have investigated the spectra–structure relationships in [13]C NMR data basis.

## OPTIMIZATION OF CORRELATION TABLES

The value of a correlation table, i.e., its information content, depends on the extent of overlap in the interval of signals. This, in turn, depends on the kind of groups chosen for the correlation table. It is immediately apparent that the more minute details the correlation table contains, the more significant overlaps arise. Therefore, considering details beyond a certain extent is useless. In the following, an example will be enlarged upon in order to see in what kind of grouping would give the correlation table maximum information. For our calculations we use the mathematical operations of information theory.[17] The mathematical basis of the information theory is summarized in the Appendix.

Let 4 units be the width of the available zone and let it include four equally probable groups A, B, C, and D. The initial entropy is

$$H_0 = \log_2 4 = 2$$

The width of the interval belonging to A, B, and C is 2 units and to D is 1.5 units. Construct the correlation table in Figure 3. A signal can appear anywhere in the interval of 0–4 as the result of a measurement. The probability that it will appear in the interval of 3–4 is the following:

$$p = (4 - 3)/4 = 0.25$$

Divide the whole domain on the basis of the number of groups in the interval. The results are summarized in Table IV. $p_1$, $p_2$, $p_3$, and $p_4$ geometrical probabilities show whether the number of possibilities is 1, 3, 4, or 4. The entropy after the evaluation of the signal is marked by $H_v$:

$$H_{vn} = \log_2 n \qquad (n = 1, 2, 3, 4)$$

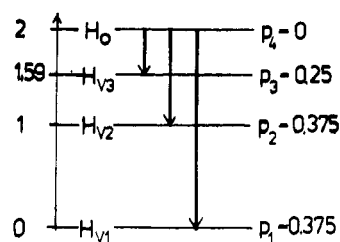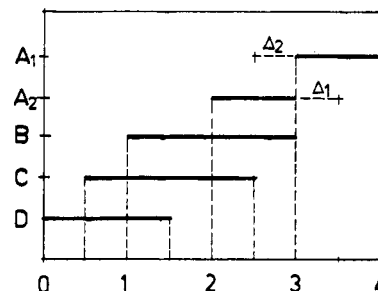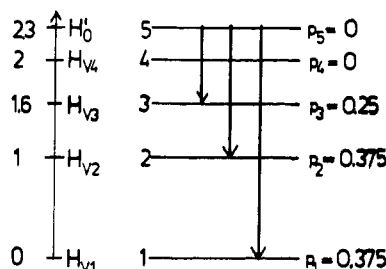Figure 4 shows the levels of entropy, and the tendency of



**Figure 4.** Process of entropy decrease.



**Figure 5.** Substitution of fragment A by fragments $A_1$ and $A_2$.



**Figure 6.** Entropy decrease under changed initial conditions. $\Delta_1 = \Delta_2 = 0$.

entropy decrease is marked by the arrows. (The final entropy levels are in the figure too.) The expected value of entropy decrease is the information ($I$), which in this case is (see Appendix eq 5)

$$I = \sum_{n=1}^{4} p_n \Delta H_n = \sum_{n=1}^{4} p_n(H_0 - H_{vn}) = 0.375 \cdot 2 + 0.375 +$$

$$0.25 \cdot 0.415 = 1.229$$

Let us assume that it seems more practical to use groups $A_1$ and $A_2$ instead of group A. Hereby, the basic interval would fall into two parts, yet the two intervals are not disjunct. Let us mark the stretch of the interval assigned to $A_1$ by $\Delta_2$ and to $A_2$ by $\Delta_1$ (Figure 5).

**I.** Let us calculate what will happen if $A_1$ and $A_2$ do not overlap, that is, $\Delta_1 = \Delta_2 = 0$. As now five groups are equally possible, the initial entropy is $H_0' = \log_2 5$; the probability of different final states does not change because the intervals of the five groups do not overlap. The new situation is shown in Figure 6.
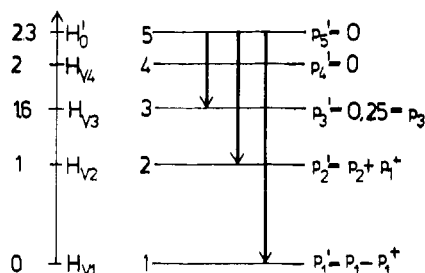
The information in the altered situation is

$$I' = \sum_{n=1}^{5} p_n \Delta H_n' = \sum_{n=1}^{5} p_n(H_0' - H_{vn})$$

The information increase is

$$\Delta I = I' - I = \sum_{n=1}^{3} p_n(H_0' - H_{vn}) - \sum_{n=1}^{3} p_n(H_0 - H_{vn}) =$$

$$H_0' - H_0 \quad (1)$$

In eq 1 the summing must be done only until 3 because the probabilities, $p_n$ ($n = 1, 2, 3, 4$), of certain levels of entropy do not change and $\sum p_n = 1$, so $p_5 = 0$.

THEORY OF CORRELATION TABLES

*J. Chem. Inf. Comput. Sci., Vol. 20, No. 4, 1980* **237**



**Figure 7.** Entropy decrease under changed probabilities. $\Delta_1 > 0$.



**Figure 8.** Entropy decrease under changed probabilities. $\Delta_2 > 0$.

The result is totally general. If there is a way to divide a group without overlap of the signals of the new groups (subgroups), the expected value of entropy decrease (the information) would increase. If we have $N$ groups in the correlation table, and we can choose one from which we can form $K$ new groups in a way that their signal intervals do not overlap, then

$$\Delta I = I' - I = H_0' - H_0 = \log_2 \frac{N + K - 1}{N} > 0 \quad (2)$$

The introduction of this kind of modification always involves information change. If $N$ increases, this information decreases and its limit value is zero.

**II.** Let us calculate what will happen if $\Delta_2 = 0$ and $\Delta_1 > 0$. In this case, besides the change of initial entropy, the probabilities would change at the expense of each other because $\sum p_n = \sum p_n' = 1$ ($p'$ is the changed probability). $p_1^+$ denotes the geometrical probability belonging to $\Delta_1$. The probability $p_1^+$ can be calculated with the equation

$$p_1^+ = \Delta_1 / \Delta_m$$

where $\Delta_m$ is the width of the total measured interval. Taking $0 < \Delta_1 \le 1$ and $0 < p_1^+ \le {}^1/_4 \Delta_1$ into consideration causes some change in our example:

$$p_2' = p_2 + p_1^+ \quad (3)$$

$$p_1' = p_1 - p_1^+$$

The change is graphed in Figure 7. The entropy levels of the final states are unchanged: $H_{vn}' = H_{vn}$ ($n = 1, ..., 5$). On the basis of the equation $\Delta I = I'' - I$, we can again calculate the change of information:

$$I'' = H_0'' - (p_1' H_{v1} + p_2' H_{v2} + p_3' H_{v3}) \quad (4)$$

$$I'' = H_0' - (p_1 H_{v1} + p_2 H_{v2} + p_3 H_{v3}) - p_1^+(H_{v2} - H_{v1}) \quad (5)$$

$$I = H_0 - (p_1 H_{v1} - p_2 H_{v2} + p_3 H_{v3}) \quad (6)$$

and

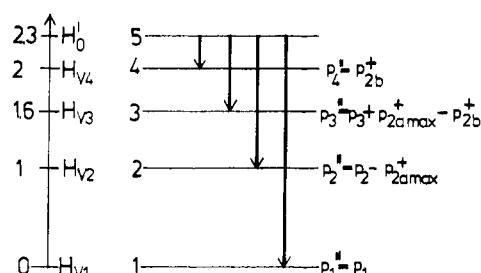$$\Delta I = H_0' - H_0 - p_1^+(H_{v2} - H_{v1}) \quad (7)$$

where $H_{v2} = 1$ and $H_{v1} = 0$.

$$\Delta I = H_0' - H_0 - p_1^+ = \log_2 \frac{N + K - 1}{N} - p_1^+ \quad (8)$$

In our example, this can also be calculated by the substitution of $N = 4$ and $K = 2$:

$$\Delta I = \log_2 \frac{5}{4} - p_1^+ = 0.32 - p_1^+$$

This change is worthwhile only if $\Delta I > 0$ is fulfilled and the larger $\Delta I$ is, the more effective the change. As this example has $0 < p_1^+ \le 0.25$, it is advisable to perform the change regardless of the value of $\Delta_1$. In the case of $N = 10$ and $K = 2$, the decision is not so certain because $\log_2 11/10 = 0.138$; if $p_1^+ < 0.138$ ($\Delta_1 > 0.55$), it is worth enlarging the correlation table, but if this fails to happen, the new table will be worse than the old one.

**III.** The next case is when $\Delta_1 = 0$ and $\Delta_2 > 0$. The geometrical probability belonging to $\Delta_2$ is $p_2^+$: $0 < \Delta_2 \le 1$; therefore, $0 < p_2^+ \le 0.25$. Now the result of the increase of the interval is different from the previous case. There are two different cases.

(a) $0 < \Delta_2 \le 0.5$, then $0 < p_{2a}^+ \le 0.125$.
(b) $0.5 < \Delta_2 \le 1$, then $0 < p_{2b}^+ \le 0.125$.

$p_{2b}^+$ appears only if $p_{2a}^+$ reached its maximum value ($p_{2a}^+{}_{max}$) and at this point $p_2^+ = p_{2a}^+{}_{max} + p_{2b}^+$. The changed probabilities are:

(a)
$$p_2' = p_2 - p_{2a}^+$$
$$p_3' = p_3 - p_{2a}^+$$

(b)
$$p_2'' = p_2 - p_{2a}'{}_{max}$$
$$p_3'' = p_3 + p_{2a}^+{}_{max} - p_{2b}^+$$
$$p_4'' = 0 + p_{2a}^+$$

In case a, the probability of the final state belonging to $H_{v2}$ decreases in favor of the entropy $H_{v3}$. In case b, the probability of $H_{v2}$ reached the minimum of case a, and it stays at this constant value, while the probability of the occurrence of $H_{v4}$ increases to the detriment of $H_{v3}$ (Figure 8).

On the basis of the correlations conducted earlier, the change in information can be calculated:

(a) $$\Delta I_a = I' - I = H_0' - H_0 - p_{2a}^+(H_{v3} - H_{v2}) \quad (9)$$

$$\Delta I_a = \log_2 \frac{5}{4} - p_{2a}^+ \log \frac{3}{22} = 0.32 - p_{2a}^+ = 0.585$$

Because $p_{2a} \le 0.125$, $\Delta I_a$ is always positive so it is reasonable to perform the change.

(b)
$$\Delta I_b = H_0' - H_0 - p_{2a}^+{}_{max}(H_{v3} - H_{v2}) + p_{2b}^+(H_{v4} - H_{v3})$$
$$\Delta I_b = 0.247 - p_{2b}^+ = 0.415 \quad (10)$$

The conclusions can be drawn in a similar way: $\Delta I_b$ is positive in the given condition so it is useful to carry out the changes.

From the very different results obtained on the basis of conditions II and III, the conclusion is obvious: only after analyzing the correlation table is it possible to decide whether the subdivision of an interval is useful or not. The information change resulting from the introduction of new fragments is significantly influenced by the existing structure of the correlation table.

The approach used in the example above considers the different fragments as equally probable. This is justifiable if nothing is known about the frequency of the occurrence of individual fragments. That is, optimization is performed on the basis of a given initial correlation table. According to this approach if a fragment is broken up into several parts, the sum of the probabilities of the new fragments will be greater than the initial probability. For example, let us assume that there are $n$ fragments. In this case the probability of the presence of fragment $k$ is $p_k = 1/n$. If this fragment is divided into three parts, the probabilities of the new fragments are $p_{k_1} = p_{k_2} = p_{k_3} = 1/(n + 2)$. This means that by specializing a group

**238**  *J. Chem. Inf. Comput. Sci., Vol. 20, No. 4, 1980*

Veszprémi and Csonka

A into fragments we neglect the fact that each of the new fragments derives from A: the individual fragments are totally independent and equal in rank. Simple correlation tables of this type are often used in computer systems, so this optimization process is advisable there.

## APPENDIX. INFORMATION THEORY

Let us define two random variables. $A$ is a random variable defined by the set of possible fragments $\{x_1, ..., x_k, ..., x_n\}$ where $n$ is the number of fragments; $B$ is defined by the set of possible signals $\{y_1, ..., y_j, ..., y_m\}$ where $m$ is the number of signals. Events $A_k$ and $B_j$ denote that $A = x_k$ and $B = y_j$, respectively.

**Specific Information.** The specific information $I_k$ of an event $A_k$ (the presence of a particular fragment) is the logarithm to base 2 of the reciprocal of the probability $p_k$ of its occurrence before analysis.

$$I_k = -\log_2 p_k \tag{1A}$$

**Entropy.** The entropy $H(A)$ (also called uncertainty or average information) is the weighted average of the values of the specific information for each event:

$$H(A) = -\sum_{k=1}^{n} p_k \log_2 p_k \tag{2A}$$

It is evident that $\sum_{k=1}^{n} p_k = 1$. $H$ is expressed in bits. From the point of view of analytical procedure, $H(A)$ is the uncertainty about the identity of the unknown fragment before analysis.

**Conditional Entropy.** The conditional entropy $H_B(A)$ is the uncertainty remaining after analysis, i.e., after measuring a signal $y_j$, and can be expressed by

$$H(A|B_j) = -\sum_{k=1}^{n} p(A_k|B_j) \log_2 p(A_k|B_j) \tag{3A}$$

where $p(A_k|B_j)$ is the conditional probability of identity $x_k$ provided a signal $y_j$ has been measured. Also $\sum_{k=1}^{n} p(A_k|B_j) = 1$ and $p(A_k|B_j) = p_{kj}/p_j$ where $p_{kj}$ is the joint probability of two mentioned events, $p_j$ is the probability of measuring a signal $y_j$ and $\sum_{y=1}^{m} p_j = 1$. The general eq 3a contains more than one term $p(A_k|B_j)$ because of the presence of errors and the signal not being unique.

The expected value of $H(A|B_j)$ is the uncertainty remaining after analysis and it is obtained by averaging over all possible signals.

$$H_B(A) = \sum_{j=1}^{m} p_j H(A|B_j) \tag{4A}$$

Furthermore:

$$H_B(A) = \sum_{j=1}^{m} p_j \left( -\sum_{k=1}^{n} \frac{p_{kj}}{p_j} \log_2 \frac{p_{kj}}{p_j} \right) = -\sum_{j=1}^{m} \sum_{k=1}^{n} p_{kj} \log_2 \frac{p_{kj}}{p_j} \tag{5A}$$

**Information Content.** The decrease of uncertainty as a result of the analysis (information obtained) can be written (provided a signal $y_j$ has been measured) as

$$\Delta H(A|B_j) = H(A) - H(A|B_j) \tag{6A}$$

The expected value of this entropy decrease is the information content of the analytical procedure, and it is obtained by averaging over all possible signals:

$$I(A,B) = \sum_{j=1}^{m} p_j \Delta H(A|B_j) = H(A) - \sum_{j=1}^{m} p_j H(A|B_j) \tag{7A}$$

and application of eq 4 to 7 leads to the following expression:

$$I(A,B) = H(A) - H_B(A) \tag{8A}$$

As $p_k = \sum_j p_{kj}$

$$H(A) = -\sum_k \sum_j p_{kj} \log_2 p_k \tag{9A}$$

and application of eq 5 and 9 to 8 leads to the following expression for the information content:

$$I(A,B) = \sum_k \sum_j p_{kj} \log_2 \frac{p_{kj}}{p_k p_j} \tag{10A}$$

The following can be established.

(a) $I(A,B) \geq 0$ because of $p_{kj} \geq p_k p_j$ for all $k$'s and $j$'s. The equality applies if and only if $A$ and $B$ are statistically independent ($p_{kj} = p_k p_j$ for all $k$'s and $j$'s).

(b) $I(A,B) \leq H(A)$. The equality applies only if $A$ and $B$ are unambiguously dependent.

(c) $I(A,B) = I(B,A)$.

## REFERENCES AND NOTES

(1) Beech, G.; Jones, R. T.; Miller, K. *Anal. Chem.* **1974**, *46*, 714.
(2) Sasaki, S.; Kudo, Y.; Ochai, S.; Abe, H. *Mikrochim. Acta* **1971**, 726.
(3) Sohár, P. "Mágneses magrezonancia spektroszkópia"; Akadémiai Kiadó: Budapest, 1976.
(4) Bible, R. H. "Guide to Empirical Method"; Plenum Press: New York, 1967.
(5) Kolonits, P. Szerves Kémiai Praktikum, Tankönyvkiadó, Budapest, 1973.
(6) Massart, D. L. *J. Chromatogr.* **1973**, *79*, 157.
(7) Dupois, P. F.; Dijkstra, A. *Anal. Chem.* **1975**, *47*, 379.
(8) Eskes, A.; Dupois, P. F.; Dijkstra, A.; De Clercq, H.; Massart, D. L. *Anal. Chem.* **1975**, *47*, 2168.
(9) Van Marlen, G.; Dijkstra, A. *Anal. Chem.* **1976**, *48*, 595.
(10) Wyandotte-ASTM (Kuentzel) Punched Cards Index; American Society for Testing and Materials, Philadelphia, Pa., 1969.
(11) Dupois, P. F.; Dijkstra, A.; van der Maas, J. H. *Fresenius Z. Anal. Chem.* **1978**, *291*, 27.
(12) Rotter, H.; Varmuza, K. *Org. Mass Spectrom.* **1975**, *10*, 874.
(13) Soltzberg, L. J.; Wilkins, C. L.; Kaberline, S. L.; Fai Lam, T.; Brunner, T. R. *J. Am. Chem. Soc.* **1976**, *98*, 7139.
(14) Woodruff, H. B.; Snelling, C. R.; Shelley, C. A.; Munk, M. E. *Anal. Chem.* **1977**, *49*, 2075.
(15) Wilkins, C. L.; Brunner, T. R. *Anal. Chem.* **1977**, *49*, 2136.
(16) Milne, G. W. A.; Zupan, J.; Heller, S. R.; Miller, J. A. *Org. Magn. Reson.* **1979**, *No. 5*, 289.
(17) Reza, F. M. "An Introduction to Information Theory"; McGraw-Hill (Hungarian translation): Müszaki Kñyvkiadó, Budapest, 1966.