

**NUMERICAL ANALYSES OF SOME COMMON ERRORS IN
CHEMOSYSTEMATICS**

BY ROBERT P. ADAMS

Made in United States of America
Reprinted from BRITTONIA
Vol. 24, No. 1, January-March, 1972
pp. 9-21

NUMERICAL ANALYSES OF SOME COMMON ERRORS IN CHEMOSYSTEMATICS

ROBERT P. ADAMS

Adams, Robert P. (Colorado State University, Fort Collins, Colorado 80521). Numerical analyses of some common errors in chemosystematics. *Brittonia* **24**: 9-21. 1972.—The effects of random and nonrandom errors were analyzed. Four types of nonrandom errors were simulated and examined in detail: 1. misidentification of compounds; 2. unresolved components; 3. errors in the quantification of individual compounds; and 4. errors in the quantification of a series of analyses. These cases were compared using data from studies in *Juniperus*. Misidentification of 1 of 40 components had only a slight effect on the resulting classification even when that error accounted for 20% of the total character weights. Unresolved components also had only a slight effect. Nonrandom errors in the quantification of individual compounds (throughout a study) were shown to have no effect on the similarity measure used in this study. Errors in the quantification of a series of analyses (such as encountered in seasonal variations) appeared to be the most serious source of error and were shown to be of considerable importance in chemosystematic studies. Several preventative procedures are suggested.

INTRODUCTION

Several types of errors are rather common in biochemical systematic studies, but their effects usually have been either discounted or magnified as impossible obstacles. Chemosystematists are often either organic chemists whose primary interests are in the identification of plant products and the elucidation of biochemical pathways, or classically trained taxonomists whose primary interest (in chemosystematics) is in the use of biochemical markers to aid in the classification of the taxa being studied and to gain insight into the evolutionary relationships of these taxa. With this dichotomy of interests it is not surprising that error analysis has been largely ignored in chemosystematics.

Random errors are certainly as common to biochemical data as they are to morphological data. Some of the more common random errors in chemosystematics are: variations in the amount of plant material to be extracted; variations in the amounts of solvents used in extraction procedures; differences in extraction times; normal voltage and current variations which influence the accuracy of equipment utilized in the detection and quantification of compounds. Most chemosystematists are at least aware of these errors and usually take some precautions to minimize them. If the data are analyzed statistically, these random errors should have little effect on most analyses.

On the other hand, nonrandom (or systematic) errors pose an entirely different problem. They can not be accounted for in analysis of variance procedures and are generally not detectable in any statistical analyses of the data. Thus, what may appear to be evidence of clinal variation or introgression, may in fact be due to seasonal differences in the collection of samples, or to the effect of gradual changes (e.g., deterioration) in a column used in gas/liquid chromatography.

The purpose of this paper is to examine the effects of four common nonrandom errors encountered in the utilization of chemical data.

I. Misidentification of a compound. The identity of a compound can be determined by several methods. These include: retention time, infrared spectrophotometry, mass spectral analysis, nuclear magnetic resonance analysis, ultraviolet spectrophotometry (especially in the presence of various reagents, as in flavonoid identification) and synthesis of the compound.

Some of these methods are more useful than others. Identification based upon retention times is usually quite error prone unless one is working with organisms that are closely related (e.g., at the infraspecific level). Even infrared analysis is useless in the identification of D and L isomers. Since the identification of a rare or unknown compound may require several months, most chemosystematists compromise by identifying only the major components of a complex mixture. The minor components are often tentatively identified based on retention times, etc.

II. Unresolved components. Most chemosystematists usually attempt to resolve complex mixtures into various components. Most techniques separate the mixtures into bands (electrophoresis, serology), spots (paper and thin layer chromatography), or peaks (gas/liquid and liquid/liquid chromatography). In complex mixtures involving 50 to 100 compounds, it is extremely difficult to resolve all components completely. When these constituents are present in small quantities (as is often the case), they may be obscured by components present in larger quantities. Frequently, identification of the major compound will not reveal the presence of the minor constituent. This type of error differs somewhat from a misidentification in that the measurement of the quantity of the major component is increased by an unknown amount.

III. Errors in the quantification of individual compounds. When mixtures are separated by means of physical-chemical characteristics, some compounds react in a manner that interferes with accurate quantification (e.g., certain compounds tend to streak when paper chromatographed; others partially decompose during analysis; certain components exhibit severe tailing in gas/liquid chromatography). In addition, the sensitivity of detection may vary with different compounds according to such factors as: number of carbon atoms per molecule, degree of saturation, presence of halogens, pH of a buffer solution, or presence of nitrogen.

IV. Errors in the quantification of a series of analyses. Seasonal variations in the chemical constituents of plant materials may be considerable; this variation must consequently be taken into consideration in planning research programs. Samples collected early may differ significantly from those collected later (Adams, 1970). This is particularly true in studies of geographic variation since apparent regional trends may only reflect seasonal variation due to the times of collection. Another example of this type of error is that encountered by the change in extraction and analysis procedures through time. These include variation in solvent purity from different suppliers, changes in columns (due to aging, bleeding, contamination, etc.), changes in extraction procedures (e.g., length of distillation time, use of different amounts of solvents), and changes in the size or shape of chambers used in paper chromatography. The latter two examples are nonrandom errors that would result, for example, from using a 2-hour distillation in the first part of a study, and subsequently switching to a 2.5-hour distillation.

MATERIALS AND METHODS

Thirteen populations of *Juniperus ashei* Buch. were sampled in December and January 1967-68. Population samples consisted of terminal branches from 5 trees at each site. Samples from individual trees consisted of 6 to 8 inches of fresh foliage from 4 or 5 branches. The fresh foliage was steam-distilled for one hour to remove the volatile terpenoids (Adams, 1970). The major components were identified by infrared spectroscopy, and the minor components were tentatively identified by their retention times (see Adams, 1970; Adams & Turner, 1970). The relative percentage of each compound was determined by an electronic digital integrator with automatic printed output, attached to the gas/liquid chromatograph.

The volatile oil of *J. ashei* contains approximately 90 terpenoid components (von Rudloff, 1968), many of which are present in only trace amounts (less than .1% of the total oil). In the first step in the analysis of infraspecific variation in *J. ashei*, a one-way analysis of variance (ANOVA) was performed on the data for each of the terpenoids for all 13 populations (65 samples). The among/within F ratios were calculated for each of the terpenoids, of which 39 had F ratios greater than 1.0 (and values greater than trace amounts in some population). These 39 compounds were selected for use in computing similarity measures.

The similarity measure used is basically a weighted matching coefficient (Sokal & Sneath, 1963). The character weight is the F ratio (determined by ANOVA) which weights those characters most heavily that show large variation among OTUs (populations) and relatively smaller variation within OTUs (populations), following the ideas of Farris (1966) and Flake & Turner (1968).

$$1.1 \quad R_{d_{xy}} = \frac{\sum_i^n F_i |x_i - y_i| / R_{g_i}}{\sum_i^n F_i}$$

$$1.2 \quad S_{r_{xy}} = 1 - R_{d_{xy}}; 0 \leq S_{r_{xy}} \leq 1$$

The clustering method used was the single linkage method of Sneath (1957). Table I shows the F weights used for the 39 terpenoid characters used in the previous study (Adams & Turner, 1970). The identities of most of these components are given in Adams (1970).

Figure 1 shows the results of clustering these 13 populations of *J. ashei*. Notice that

Let:

- x_i = value of character i in OTU (population) x
- y_i = value of character i in OTU (population) y
- $R_{d_{xy}}$ = relative dissimilarity between OTUs x and y
- $S_{r_{xy}}$ = relative similarity between OTUs x and y
- $\quad = 1 - R_{d_{xy}}$
- F_i = F ratio for character i
- R_{g_i} = range of character i encountered in all population averages
- n = number of comparisons between OTUs x and y (excluding negative matches)

Then:

populations 33, 40, and 30 appear to be somewhat different from the rest of the populations, with population 35 being rather loosely connected to the main group. The similarity matrix (not shown) and this phenogram will serve as references for the first three cases below.

RESULTS

I. Misidentification of a compound. In order to simulate a misidentification, 2A(α -pinene) was recorded as compound 2B(α -thujene) in populations 42 and 44. Since compound 2B is normally absent in all populations of *J. ashei*, it was anticipated that this would represent a rather drastic error. It should be noted that generally only quantitative variations were observed in the terpenoids within this species. Thus, when F ratios were computed, compound 2A appeared in varying amounts in all populations except populations 42 and 44, where it was recorded as 0.0. Likewise compound 2B was

TABLE I
THIRTY-NINE COMPOUNDS OF *J. ashei* WITH THE F RATIOS USED AS WEIGHTING FACTORS.
 $F_{.01} = 2.61$, $df = 12/51$. $\Sigma F = 147.79$

Cpd. #	F ratio	Cpd. #	F ratio	Cpd. #	F ratio
1	1.449	2A	3.536	3	2.507
7	14.327	9	10.112	11	8.985
13	4.554	14	5.033	17	2.683
18	1.908	21	4.201	24	1.447
26	1.907	27	2.356	29	1.232
30	5.498	32	2.875	32A	2.524
33	2.497	34	1.247	35	1.158
36A	4.148	37	1.800	37A	3.481
39	2.086	40A	1.093	43	14.254
44	2.799	47	1.593	48	1.381
50	5.450	51	5.878	62	2.852
64	1.312	66	5.480	67	1.473
69	5.140	74	1.265	75	4.268

recorded as having varying amounts in populations 42 and 44 and 0.0 in the other 11 populations.

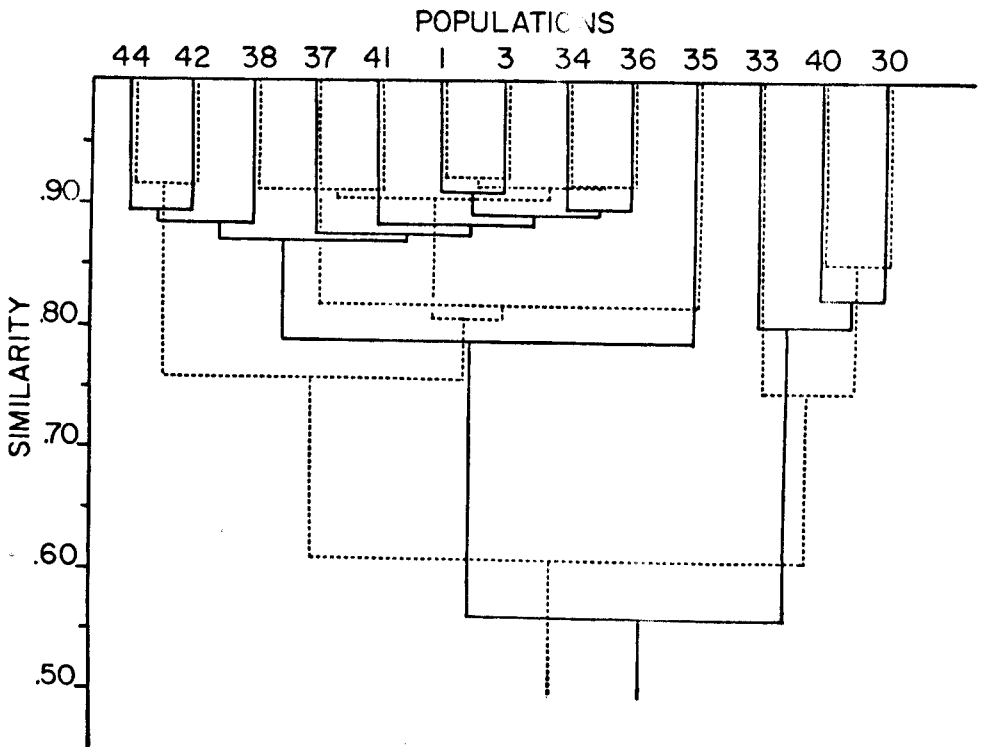
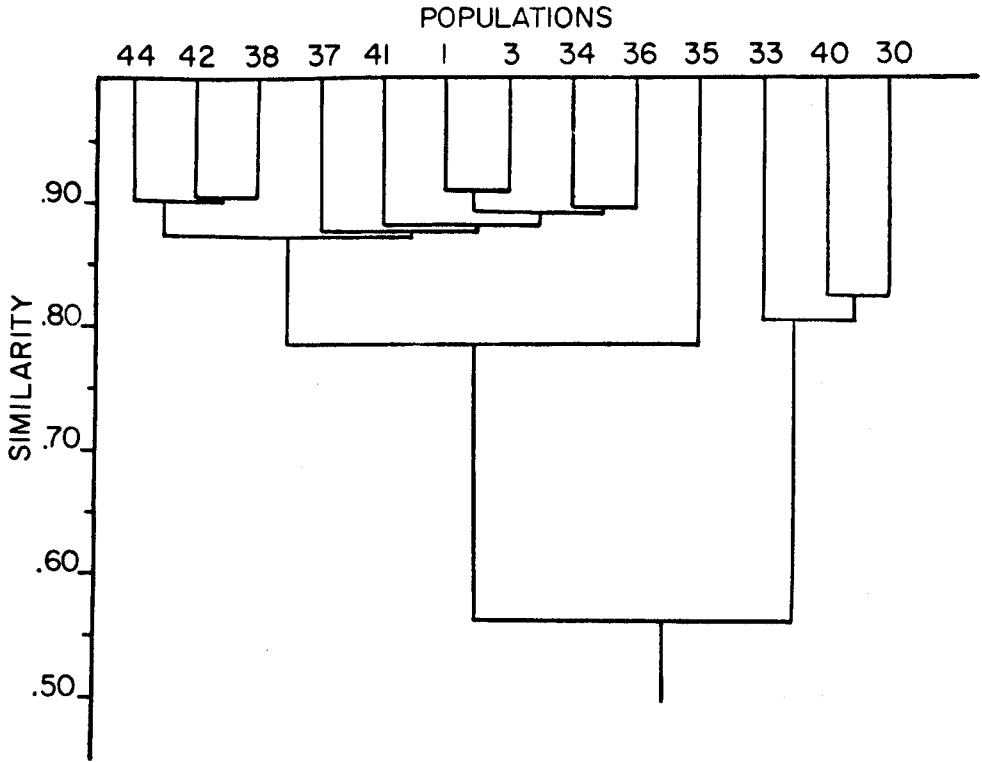
Whereas compound 2A previously had an F ratio of 3.54, the new F ratio was 5.70, and compound 2B (which was previously absent) now had an F ratio of 4.76 ($F_{.01} = 2.61$, $df = 12/51$). Inclusion of this compound now makes a character set of 40.

Since the characters were weighted, it is misleading to think that this represents only a 1/40 change of the similarity measures. Note that the sum of the F weights (Table I) is 147.79. In the original analysis the weight of compound 2A was 3.54, or $3.54/147.79 = 2.4\%$ of the total character weighting. When the new character was created, the combined weight of characters 2A and 2B was 10.46, or $10.46/(144.25 + 10.46) = 6.8\%$ of the total character weighting. Although 6.8% represents a significant portion of the similarity measure, it would be interesting to know how an even greater contribution of this pair (2A/2B) would affect the classification. Therefore, the modified data was reanalyzed using the following F values as character weights for both 2A and 2B: 8.0 each or $16.0/(144.25 + 16.0) = 10\%$ of the total weights; 18.0 each or $36.0/(144.25 + 36.0) = 20\%$ of the total weights; 31.0 each or $62.0/(144.25 + 62.0) = 30\%$ of the total weights; and 72.0 each or $144/(144.25 + 144.0) = 50\%$ of the total weights. It should be noted that none of the original F weights for the other 38 characters was changed from the original analysis.

These five similarity matrices were computed along with the mean and standard deviation. In addition, the correlation of each of these matrices with the original similarity matrix was determined (Sokal & Sneath, 1963) as well as the 95% confidence limits of the correlation coefficient (Steel & Torrie, 1960). Figure 2 shows the phenogram (in solid lines) obtained by single linkage clustering of first modified data (F weights of 2A + 2B = 6.8% of total); the phenogram in dotted lines (Figure 2) is that of the fifth modification (F weights of 2A + 2B = 50% of total weights). When the misidentification

FIG. 1. (Upper). A phenographic representation of single linkage clustering of 13 populations of *J. ashei* using 39 terpenoid characters, F weighted.

FIG. 2. (Lower). Solid lines indicate the phenogram resulting from single linkage clustering using a simulated misidentification of compound 2A in populations (OTUs) 42 and 44, with the F weight of cpd. 2A plus the F weight of cpd. 2B being 6.8% of the total character weights. The dotted lines show the phenogram using the same data with the F weight of cpd. 2A plus the F weight of cpd. 2B equal to 50% of the total character weights. See text for further explanation.



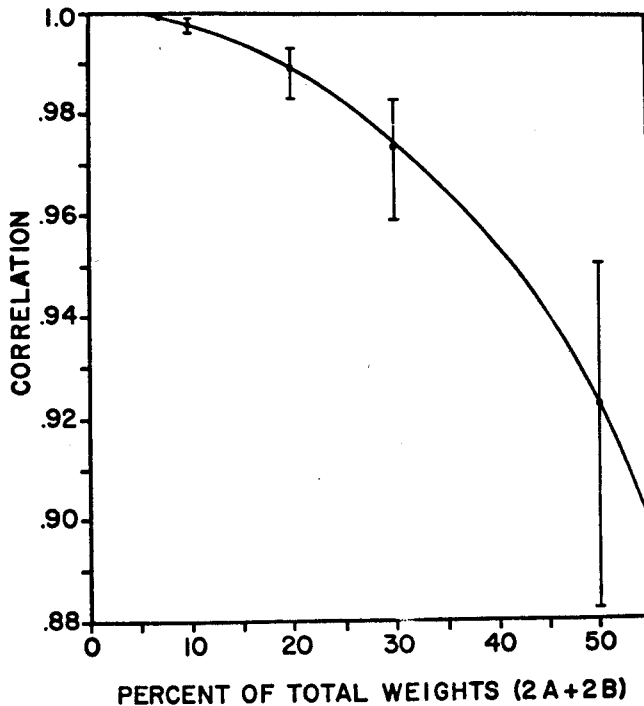


FIG. 3. Correlation of the original similarity matrix with subsequent modifications of the data involving weighting factors of 6.8, 10.0, 20.0, 30.0, and 50 percent of the total character weights ascribable to characters 2A and 2B. Vertical bars indicate the 95% confidence limits for each of the correlation coefficients.

tion represents only 6.8% of the total weights there is very little influence on the phenogram. Notice that OTUs 42 and 44 (which now share compound 2B) did pull away from the group a little, but otherwise the classification is practically unaffected. This is evident from Figure 3, which shows the correlation of the original data with the subsequent modifications. As the proportion of the total weighting increases, we see an approximately logarithmic decrease in the correlation coefficient. Even though there is a fairly high correlation between the original data and the fifth modification (.923), Figure 2 (dotted lines) shows that the phenogram has lost much of its structure and one is approaching a 2-character classification. Figure 4 shows the change in average similarity ratios. It appears that there is a linear decrease in the average similarity ratio which reflects the increased dichotomy within the similarity matrix. The standard deviations of the means of the similarity matrices varied from .025 to .026.

Thus it would appear from this limited analysis that one might tolerate misidentification which contributed up to 20% of the total character weighting without severely distorting the classification. Further analysis is needed to determine if this is true when several characters might be misidentified with a cumulative weight of 20% of the total weights.

II. Unresolved components. Previous work (Adams & Turner, 1970) showed that compounds 7 (myrcene) and 43 (carvone) had very high F ratios, 14.33 and 14.25, which is about $28.58 \div 147.79 = 19.3\%$ of the total character weights. Furthermore, these compounds were inversely correlated within populations. Thus it seemed that the combination of these two compounds (in each plant analysis) would represent a good example of important, unresolved components. The original data deck was modified

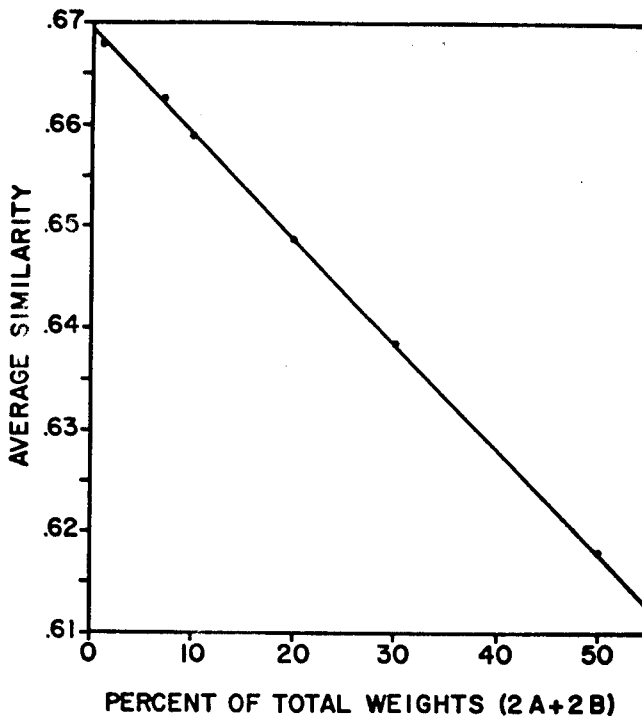


FIG. 4. Average similarity measures for the original and five modifications of the data where the weights of character 2A plus 2B were 6.8, 10.0, 20.0, 30.0, and 50.0 percent of the total character weights. The standard deviation of the means varied from .025 to .026.

such that for each plant, component 7 = compounds 7 + 43, and compound 43 = 0.0. When the modified data was analyzed, component 7 (#7 + #43) had a new F ratio of 0.74 and was therefore not used in computing the similarity measures since the F ratio was less than 1.0. Thus, when compounds 7 and 43 were not resolved into two components, they essentially canceled each other. The remaining 37 characters were used to obtain the phenogram in Figure 5.

Notice that in comparison with Figure 1, the chief difference is that OTU 37 has failed to enter the group 41, 1, 3, 34, 36, but joins a super-set at a lower similarity value. In addition, OTU 34 joined first with OTUs 1 and 3 rather than with OTU 36. The correlation between the original similarity matrix and that using 37 characters was 0.9965, which fortifies the comparisons between the phenograms.

The use of F ratios as weighting factors is exceptionally well suited in this case since the loss of significant information due to unresolved components is reflected in the decrease in the F ratio.

To investigate this case further, the original data deck was modified by the consolidation of not 2 but 4 compounds into component 7. These were compounds 7 (myrcene), 26 (camphor), 27 (linalool), and 43 (carvone). The original values of the last three compounds were replaced by zero for each plant.

The F ratio of the new component 7 (compounds 7 + 26 + 27 + 43) was 2.53. The similarity measures were therefore computed using 36 characters (compounds 26, 27, 43 having been submerged with compound 7 into component 7). The result of clustering is shown in Figure 6. The merging of these compounds seems to have had little effect in comparison with Figure 5. The correlation between the original similarity

matrix and the matrix obtained using the modified data was 0.9962, indicating very little distortion. Examination of the F weights of characters 26 ($F = 1.91$) and 27 ($F = 2.36$) reveals that these two characters were not very heavily weighted (cf. F ratios in Table I). Thus, it is not too surprising that little additional distortion was produced when compounds 26 and 27 were added to the composite component 7 (#7 + #43). It should be noted, however, that although compounds 7, 26, 27, and 43 account for over 20% of the total F weights, there is still good agreement in the phenograms and similarity matrices when their combined weights are dropped to less than 2% of the total. This seems to indicate considerable robustness in the numerical procedures.

III. Errors in the quantification of individual compounds. Quantitative measurements of individual compounds are generally in error by some constant percent for each compound. For instance, a flame ionization detector might consistently underestimate the quantity of myrcene in a sample as 94% of the true value. Since chemosystematists often do not apply correction factors to their data, it is instructive to examine the effect of this type of error.

- Let: n = number of comparisons between OTUs x and $y = 1$
 X_i = true value of compound i in the sample representing OTU x
 Y_i = true value of compound i in the sample representing OTU y
 α = relative error in measuring the true value of compound i (e.g., 95% of true value, 103% of true value)
 Rg_i = true value of the range of compound i in all observations
= maximum X_i - minimum X_i
 F_i = variance, S_a^2 , of compound i among populations (OTUs)/variance, S_w^2 , of compound i within populations (OTUs)

Then:

- 2.1 measured value of compound $i = \alpha \cdot X_i$
2.2 measured value of the range of compound $i = \alpha(\max_i - \min_i) = \alpha \cdot Rg_i$

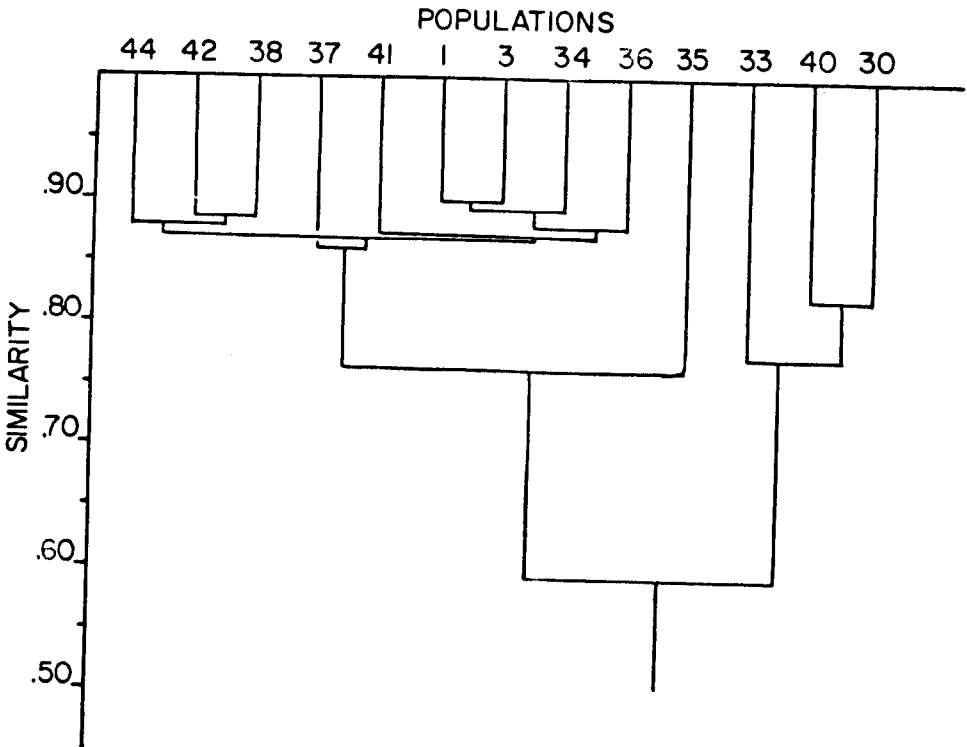
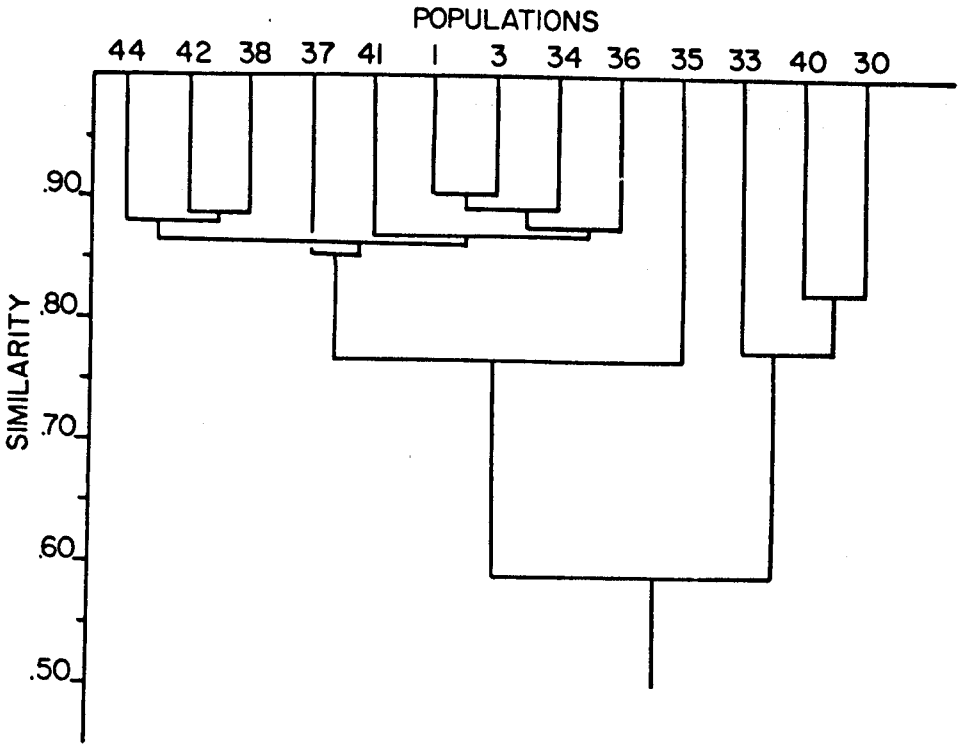
Since the observed values of X_i are multiplied by a common factor, α , the variance of $\alpha \cdot X_i = \alpha \cdot S^2$, (Sokal & Rohlf, 1969). Thus the weighting factor, F_i , is unaffected by this systematic error since, $F_i = \alpha \cdot S_a^2 / \alpha \cdot S_w^2 = S_a^2 / S_w^2$. If the relative error, α , in the quantification of compound i does not change from OTU x to OTU y , then the relative similarity between OTUs x and y (using the measured values of compound i) is:

$$2.3 \quad S_{r_{xy}} = 1 - R_{d_{xy}} = 1 - \frac{\sum_i F_i |\alpha \cdot X_i - \alpha \cdot Y_i| / \alpha \cdot Rg_i}{\sum_i F_i}$$

$$2.4 \quad S_{r_{xy}} = 1 - \frac{\sum_i F_i \cdot \frac{\alpha}{\alpha} \cdot |X_i - Y_i| / Rg_i}{\sum_i F_i}$$

FIG. 5. (Upper). A phenogram showing the results of clustering using unresolved components (cpd. 7 = cpd. 7 + cpd. 43).

FIG. 6. (Lower). The phenogram resulting from clustering using 36 compounds in which 4 are unresolved (cpd. 7 = cpd. 7 + cpd. 43 + cpd. 26 + cpd. 27).



$$2.5 \quad S_{r_{xy}} = 1 - \frac{\sum R_i |X_i - Y_i| / R_{g_i}}{\sum F_i}$$

Thus, equation 2.5 reduces to the original similarity measure, equations 1.1 and 1.2. For those readers who are not mathematically inclined, it will suffice to report that when compound 7 is replaced throughout the study by 105% of compound 7, and compound 43 is replaced by 95% of compound 43, the resulting similarity matrix is identical to the original.

It is obvious that if the relative error, α , is not changed by the addition of other compounds, the total similarity measure will be unaffected by this type of systematic error. It should be noted, however, that the effects of this type of error might be quite significant when other types of similarity measures are used.

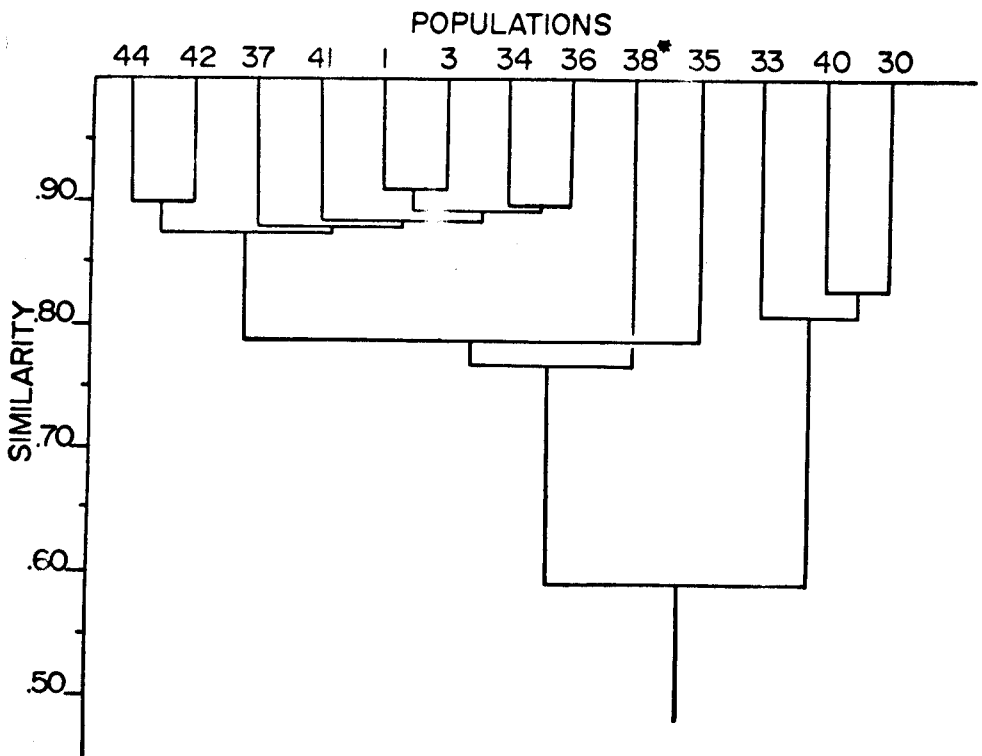
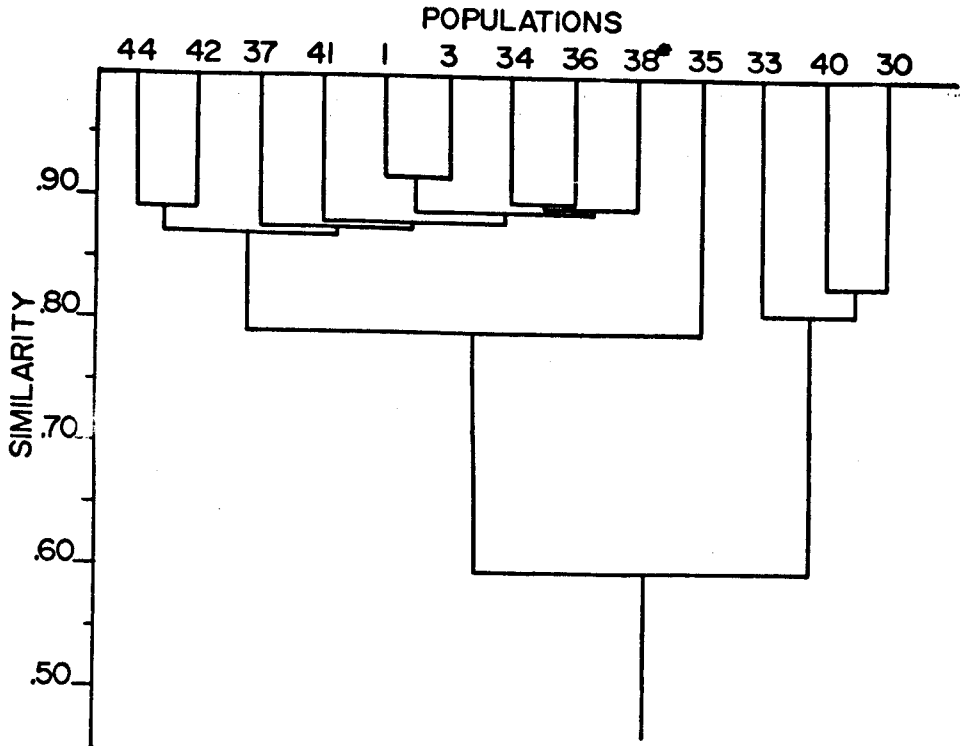
IV. Errors in the quantification of a series of analyses. The effects of this type of error are usually extremely difficult to anticipate. It is therefore with considerable benefit from hindsight that the effects of seasonal variation will be examined. Since data was available from 4 trees sampled in both July and January in the vicinity of Austin, Texas (population 38), it was decided to replace the original 5 tree samples of population 38 (cf. Figure 1) with these 4 samples. Figure 7 shows the new phenogram utilizing 4 different winter samples for population 38* and the original January samples of the other 12 populations. Note that the use of these 4 samples (instead of the original 5) had some effect on the phenogram in that population 38* now clusters more closely with populations 34 and 36 rather than with population 42 (cf. Figure 1). Otherwise the phenograms are very similar. When the similarity measures are recalculated using the same four trees sampled in July, 1967, the resulting F weights were slightly changed but the phenogram was drastically modified (Figure 6)! In comparison with Figure 5 (winter sampling of all populations), the phenogram in Figure 6 (summer sampling of population 38*, winter sampling of the other 12 populations) depicts population 38* to be rather dissimilar to all of the Central Texas populations (44, 42, 37, 41, 1, 3, 34, 36). The correlation between the similarity matrices of the summer and winter data (for population 38 only) was 0.997, which indicates good general agreement. In this case, when only 1 OTU has been changed, the correlation coefficient is quite inappropriate in showing the changes in the resulting classification.

Although this error appears very significant in a sensitive study (such as geographical variation), it would generally be of much less importance in a study of differences between species or genera. Nevertheless, it is certainly necessary to minimize these effects in a chemosystematic study. Some of the methods of lessening these errors in the quantification of a series of analyses are:

1. Randomize the order of analyses in the laboratory. This tends to randomize the effects of column aging and bleeding, changes in solvent purity, etc.
2. Utilize more time in the development of extraction and analysis procedures. More thorough planning of laboratory procedure will usually eliminate the need to change techniques in the middle of a study.

→
FIG. 7. (Upper). The results of clustering using 4 different winter samples for population 38* (instead of the original 5 trees) and the original winter samples for the other 12 populations.

FIG. 8. (Lower). The phenogram obtained using 4 summer samples for population 38* and the original winter samples of the other 12 populations.



3. Attempt to collect all of the material needed in the study at the same stage of development and as nearly as possible under the same environmental conditions.
4. Use a standard reference which is run at regular intervals to insure that if systematic errors do occur, correction factors can be applied to the data.

CONCLUSION

Crovello (1969) has found that different subsets of characters yielded somewhat different classifications of *Salix*. In general his work seems to confirm the "matches asymptotes" hypothesis of Sokal & Sneath (1963), but sheds essentially no light on the problems which plague most chemosystematists. The work by Fisher & Rohlf (1969), although primarily a validation of the distance coefficient (versus the correlation coefficient) as a measure of similarity, did indicate that one might tolerate random noise in perhaps as many as 6 of 74 characters without producing severe distortion of the classification. In fact they obtained a correlation of $0.95 \pm$ between the original similarity matrix and a subsequent similarity matrix when 48 of 74 characters were scrambled!

In general the robustness of numerical taxonomic methods was borne out in the present study in that: errors in misidentification were shown to be small even when given 20% of the total character weights; unresolved components may result in some loss of detail due to lower weighting factors (F ratios), at least when discordant variation occurs, but the use of statistical procedures in assigning weights tends to eliminate this "noise"; uniform errors in the quantification of an individual component throughout a study was shown to have no effect on the similarity measures or resulting classification.

The most serious error analyzed in this study was that due to seasonal variation in the quantities of terpenoids. In order to minimize this kind of error, I would urge chemosystematists to: use standardized collection and analysis procedures; use a standard reference at regular intervals during the analyses; randomize the order of analyses in the laboratory.

ACKNOWLEDGMENTS

The author is grateful to Drs. Amy Jean Gilmartin and T. J. Crovello for critical readings of this manuscript. The first part of this research was supported by NSF Grant GB-5548X and the latter part was supported by NSF Grant GB-24320. Computer time was furnished by Colorado State University.

LITERATURE CITED

- Adams, R. P. 1970. Seasonal variation of terpenoid constituents in natural populations of *Juniperus pinchotii* Sudw. *Phytochemistry* **9**: 397-402.
- Adams, R. P. & B. L. Turner 1970. Chemosystematic and numerical studies of natural populations of *Juniperus ashei* Buch. *Taxon* **19**(5): 728-752.
- Crovello, T. J. 1969. Effects of change of characters and of number of characters in numerical taxonomy. *Amer. Midl. Naturalist* **81**(1): 68-86.
- Farris, J. S. 1966. Estimation of conservatism of characters by constancy within biological populations. *Evolution* **20**(4): 587-591.
- Fisher, D. R. & F. J. Rohlf 1969. Robustness of numerical taxonomic methods and errors in homology. *Syst. Zool.* **18**(1): 33-36.
- Flake, R. H. & B. L. Turner 1968. Numerical classification for taxonomic problems. *J. Theor. Biol.* **20**: 260-270.
- Rudloff, E. von 1968. The volatile oil of the leaves of *Juniperus ashei* Buchholz. *Canad. J. Chem.* **46**: 679-683.

- Sneath, P. N. A.** 1957. The application of computers to taxonomy. *J. Gen. Microbiol.* **17**: 201-226.
- Sokal, R. R. & F. J. Rohlf** 1969. *Biometry*. W. H. Freeman and Co., San Francisco and London.
- Sokal, R. R. & P. N. A. Sneath** 1963. *Principles of Numerical Taxonomy*. W. H. Freeman and Co., San Francisco and London.
- Steel, R. G. D. & J. H. Torrie** 1960. *Principles and Procedures of Statistics*. McGraw-Hill Book Co., New York.