

SCALING PROCEDURE*

DAVID KLAHR

UNIVERSITY OF CHICAGO

Recent advances in computer based psychometric techniques have added a collection of powerful tools for analyzing nonmetric data. These tools, although primarily well suited to the nonmetric case, have general potential pitfalls. Among other things, there is no statistical test

of the statistical significance of results yielded by Kruskal's nonmetric multidimensional scaling. The estimates of the relative frequency with which randomly generated sets of data reveal the relative frequency with which apparent structure is erroneously found in unstructured data. For a small number of points (i.e., six or seven) it is very likely that a good fit will be obtained in two or more dimensions when in fact the data are generated by a random process. The estimates presented here can be used as a benchmark to which to analyze the significance of the results obtained from empirically based nonmetric multidimensional scaling.

1. Introduction

Recent developments in psychological scaling [Shepard, 1962a, 1962b; Kruskal, 1964a, 1964b; Lingoes, 1965] have yielded a collection of computer

procedures which, when used in conjunction with these more traditional methods, provide powerful tools for the analysis of such data. They are being used with increasing frequency in a diverse range of applications, e.g., color vision, Morse Code perception [Shepard, 1962] tests [Russell & Gregson, 1966] college admissions [Klahr, 1968], marketing [Green, Carmone & Robinson, 1968].

However, the great strength of these new procedures must be employed with careful application. Torgerson [1955, p. 361], in reviewing some of the advances in multidimensional scaling, says:

The new procedures would . . . seem to offer nothing but advantages over the old; to require very little and to yield very much. Yet there are many problems connected with their use which . . . have not been at all obvious . . . It is like doing a factor analysis and the factor analysis methods always yield an answer. But it can be even more difficult to fully comprehend the meaning of that answer.

*A preliminary version of this paper was presented at the International Federation for Information Processing Congress 68 in Edinburgh, Scotland, August 5-10, 1968.

Torgerson goes on to discuss several of the difficulties. One which he treats most extensively in his paper is the problem of the *range* of similarity itself. He shows that under certain circumstances the implicit models underlying the scaling techniques are quite inappropriate. In particular, the procedure of Green [1966] is inappropriate for the scaling of items, i.e., upon the other relevant items in the collection of things to be judged. The scaling models do not allow for this kind of variability in similarity judgments.

A second difficulty, the problem of multiple solutions (i.e., non-uniqueness) is also discussed. For 15 or more points we can be very sure that the solution is unique, and for 15 or more points we can be virtually certain that the best fit is the one which best fits the data.

Green [1966] suggests that "these . . . procedures need Monte Carlo and other computer runs to determine their properties. . . ." This paper deals

with the statistical significance of the results. We investigate, through Monte Carlo simulation, the following question: How likely are we to falsely reject the

null hypothesis? There are at present no statistical methods for testing the significance of the results generated by the scaling procedures. As a first step towards filling this

gap, a Monte Carlo simulation is run on randomly generated data. The results to be presented are based upon a series of attempts to scale randomly generated data. They can be used as a bench mark against which to assess the significance of empirical results.

2. *Nonmetric Multidimensional Scaling*

All multidimensional scaling procedures assume that there is an underlying structure—a spatial configuration of items in which inter-item distances are a function of the relative similarity of those items. The goal of the procedures is to construct such a configuration from the inter-item proximity measures. For example, if the proximity measures are judgments of the relative similarity of pairs of stimuli, the scaling procedure would attempt to construct a spatial configuration in which the items are close together.

The three most widely used procedures—Shepard's [1962], Kruskal's [1964a, 1964b], and Lingoes' [1966]—are designed for the same purposes. That is, they yield essentially similar outputs, given the same input. Although the procedures achieve their desired result in different ways, for the purposes of our analysis it will suffice to discuss only one of them. Therefore, the re-

remainder of this paper will be based upon Kruskal's nonmetric multidimensional scaling procedure [1964a, 1964b].

The goal of the procedure is to find the spatial configuration of a set of points in which the rank order of the interpoint distances is maximally similar to the rank order of the corresponding interitem similarity measures. It starts with an arbitrary configuration of points and iteratively adjusts the interpoint distances to be the opposite of the rank order of the data perfectly. As the dimensionality of the space is reduced and the solution becomes more highly constrained, we are apt to get some departures from perfect fit. Some of the distances may be "out of order." A measure of

developed by Kruskal [1964a]: it is quite similar to a residual sum of squares.

Kruskal suggests that departures from perfect fit (stress) be interpreted as follows: "If the stress is small, the solution is a good one; if it is large, the solution is poor."

We expect minimum stress to increase as the dimensionality decreases, starting in $n - 1$ space with zero stress.

The decision as to which configuration is the most appropriate representation of items rests upon scientific judgments and is not a direct output of the scaling technique. In most applications the decision is based upon the configuration itself. Kruskal's stress measure is not a statistical test for the significance of the stress, although it is appropriate to use it as a measure of fit. It is possible to obtain estimates of significance through the use of the techniques described in the next section.⁴

⁴Since this work was done, a similar study by Stenson and Knoll (1969) has appeared covering a different range of parameters: $n = 1$ to 10 in steps of 1, and $n = 10$ to 60. For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

⁵For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

⁶For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

⁷For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

⁸For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

⁹For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

¹⁰For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

¹¹For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

¹²For large n , there is much less variance in the final stress, and the averages based on 3 data sets are sufficient. For the smaller studies in our paper, there is a need for more samples.

3. Procedure

In all applications of multidimensional scaling techniques the input is a matrix of interitem similarity measures. Whatever the metric for the raw data, the scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves.

Whatever the metric for the raw data, the scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves.

Whatever the metric for the raw data, the scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves.

Whatever the metric for the raw data, the scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves.

Whatever the metric for the raw data, the scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves.

Whatever the metric for the raw data, the scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves.

Whatever the metric for the raw data, the scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves. The scaling procedure makes use of the similarity measures themselves.

randomly assigned values from a uniform distribution on the open interval from 0 to 1. One hundred such sets of random proximities were generated for each value of $n = 6, 7, 8$ and 10; fifty sets each were generated for $n = 12$ and 16. Every one of the 500 sets was scaled by Kruskal's nonmetric multidimensional scaling program by minimizing from 0 to 1000 iterations the values that control Kruskal's program were: minimum stress sought—zero; maximum number of iterations—75; type of proximity measure—dissimilarities; spatial metric—Euclidean.

4. Results

The stress of each final configuration, for it is the behavior of stress under conditions of pure noise that this study seeks to portray.

For each value of n we have plotted the cumulative distribution of the final stress values from one, two and three dimensional configurations in Figures 1 through 3 respectively. If an analytic expression could be developed for the statistical behavior of the stress, then it would presumably be possible of generating a family of curves quite similar to these empirical estimates.

Table 1 presents some summary statistics for each value of n . The number of solutions less than or equal to Kruskal's "good" (.05) and "excellent" (.01) stress levels are included along with the maximum, minimum, average and standard deviation of final stress values. Table 2 contains some selected plots of the average final stress shown in Table 1.

5. Discussion

These data may provide some assistance in detecting and avoiding

For example "good" solutions (i.e. stress $< .05$) are often attainable for

levelly generated. (For 6 points, 66 out of 100; for 7 points, 74 out of 100;

for 8 points, 83 out of 100.) A small increase in the number of points (e.g.,

$n > 10$) provides a substantial reduction in the likelihood of this kind of

for $n > 10$

for a given number of points stress increases with decreasing dimensionality.

One of the criteria that Kruskal suggests for selecting the appropriate

dimensionality is an "elbow" in the plot of stress vs. dimensionality. In

major decrease in the number of points is observed by decreasing dimensionality.

The lack of distinctive elbows in the plots in Figure 4 arises from the

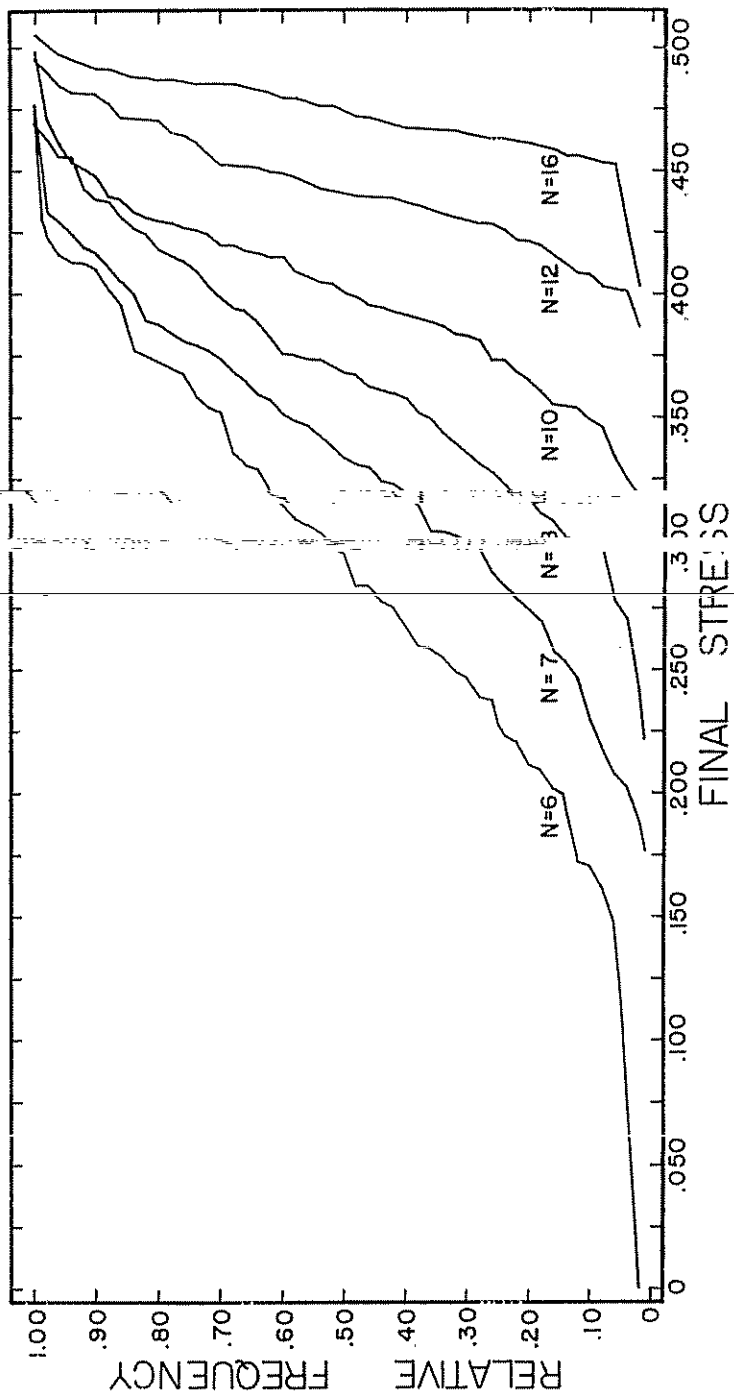


Fig. 1.--Cumulative distributions of relative frequency of final stress in a 1 dimensional configuration for 6, 7, 10, 12, and 16 points.

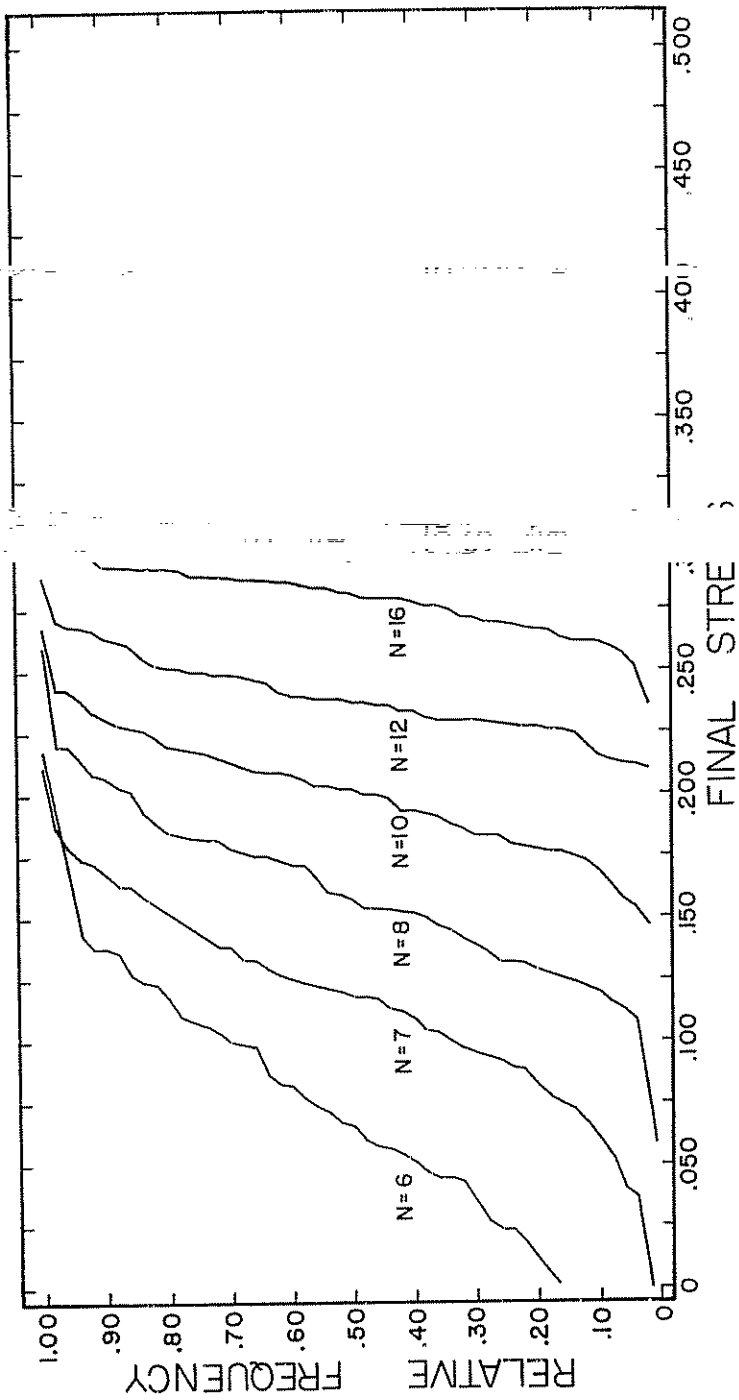


Fig. 2.--Cumulative distributions of relative frequency of final stress in a 2 dimensional configuration for 6, 7, 10, 12, and 16 points.

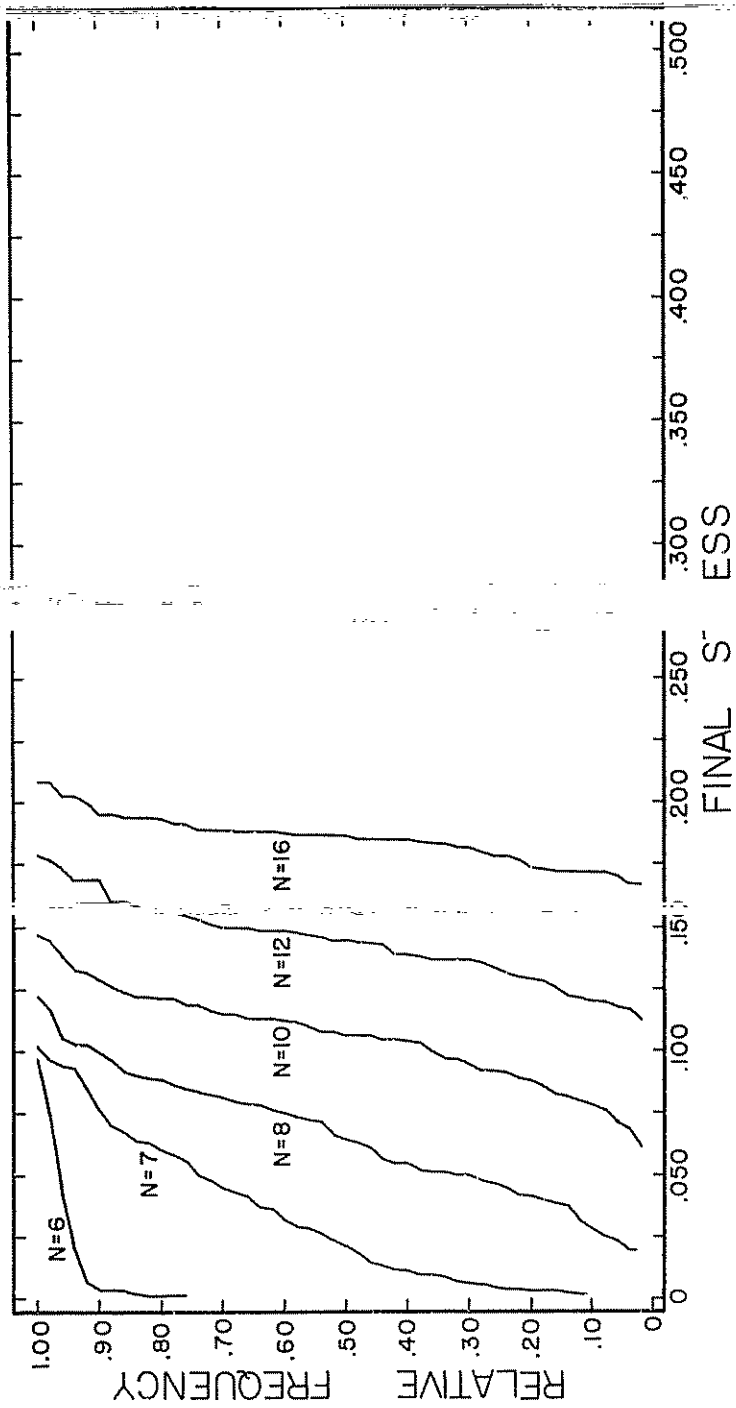


Fig. 3.--Cumulative distributions of relative frequency of final stress in a 3 dimensional configuration for 7, 8, 10, 12, and 16 points.

TABLE 1

Summary Statistics From Monte Carlo Runs: Average, Standard Deviation, Max, Min, Good and Excellent Stress

	1	2	3	4	5
Number of "Good" Final Stresses	4	39	95	100	100
Number of "Excellent" Final Stresses	4	27	95	100	100
Average Final Stress	0.288	0.010	0.010	0.010	0.010
Standard Deviation of Final Stresses	0.097	0.003	0.003	0.003	0.003
Maximum Final Stress	0.471	0.008	0.008	0.001	0.001
Minimum Final Stress	0.	0.	0.	0.	0.
Number of "Good" Final Stresses	0	7	74	100	100
Number of "Excellent" Final Stresses	0	3	53	100	100
Average Final Stress	0.332	0.004	0.004	0.002	0.002
Standard Deviation of Final Stresses	0.066	0.003	0.003	0.005	0.005
Maximum Final Stress	0.477	0.002	0.002	0.025	0.025
Minimum Final Stress	0.177	0.	0.	0.	0.
Number of "Good" Final Stresses	0	0	53	53	53
Number of "Excellent" Final Stresses	0	0	8	73	73
Average Final Stress	0.368	0.010	0.010	0.016	0.016
Standard Deviation of Final Stresses	0.055	0.004	0.004	0.018	0.018
Maximum Final Stress	0.470	0.000	0.000	0.079	0.079
Minimum Final Stress	0.222	0.009	0.009	0.	0.
Number of "Good" Final Stresses	0	0	0	3	3
Number of "Excellent" Final Stresses	0	0	0	0	0
Average Final Stress	0.400	0.011	0.011	0.055	0.055
Standard Deviation of Final Stresses	0.038	0.005	0.005	0.017	0.017
Maximum Final Stress	0.498	0.007	0.007	0.096	0.096
Minimum Final Stress	0.314	0.012	0.012	0.008	0.008

6 points
100 sets

7 points
100 sets

8 Points
100 sets

10 points
100 sets

Dimensions

4
1
0
0.0
0.0
0.1
0.0

12 poi
50 set

TABLE 1--Contin

	Nu								
Number of "Good" Final Stresses	0	0	0	0	0	0	0	0	0
Number of "Excellent" Final Stresses	0	0	0	0	0	0	0	0	0
Average Final Stress	0.4	0.240	0.185	0.1	0.057	0.057	0.057	0.057	0.057
Standard Deviation of Final Stresses	0.0	0.017	0.010	0.0	0.015	0.015	0.015	0.015	0.015
Maximum Final Stress	0.4	0.288	0.297	0.1	0.089	0.089	0.089	0.089	0.089
Minimum Final Stress	0.3	0.210	0.166	0.1	0.027	0.027	0.027	0.027	0.027
Number of "Good" Final Stresses	0	0	0	0	0	0	0	0	0
Number of "Excellent" Final Stresses	0	0	0	0	0	0	0	0	0
Average Final Stresses	0.4	0.279	0.185	0.1	0.096	0.096	0.096	0.096	0.096
Standard Deviation of Final Stresses	0.0	0.014	0.010	0.0	0.011	0.011	0.011	0.011	0.011
Maximum Final Stress	0.5	0.300	0.178	0.1	0.127	0.127	0.127	0.127	0.127
Minimum Final Stress	0.4	0.237	0.112	0.1	0.074	0.074	0.074	0.074	0.074

TABLE 2

Selected Percentile Points From Cumulative
Distributions of Final Stress

		Percentile							
		5	10	25	50	75	90	95	
Dimensions	1	.121	.171	.227	.296	.363	.410	.415	6 points 100 sets
	2	.000	.000	.023	.065	.107	.138	.148	
	3	.000	.000	.000	.000	.000	.004	.020	
	4	.000	.000	.000	.000	.000	.000	.000	
Dimensions	1	.200	.230	.269	.334	.391	.447	.467	7 points 100 sets
	2	.000	.000	.000	.000	.000	.000	.000	
	3	.000	.000	.004	.021	.056	.077	.094	
	4	.000	.000	.000	.000	.000	.006	.009	
	5	.000	.000	.000	.000	.000	.000	.000	
Dimensions	1	.276	.298	.326	.396	.412	.439	.456	8 points 100 sets
	2	.111	.120	.132	.157	.183	.206	.215	
	3	.022	.028	.046	.064	.084	.100	.104	
	4	.000	.000	.002	.010	.023	.042	.050	
	5	.000	.000	.000	.000	.001	.005	.007	
Dimensions	1	.333	.349	.373	.401	.427	.440	.461	100 sets
	2	.153	.160	.180	.202	.216	.222	.246	
	3	.000	.070	.092	.107	.119	.129	.130	
	4	.027	.035	.044	.054	.067	.077	.079	
	5	.004	.009	.018	.024	.032	.040	.050	
Dimensions	1	.402	.408	.429	.441	.465	.481	.485	50 sets
	2	.000	.000	.000	.000	.000	.000	.000	
	3	.000	.000	.000	.000	.000	.000	.000	
	4	.071	.072	.075	.086	.094	.109	.118	
	5	.032	.037	.047	.057	.063	.078	.085	
Dimensions	1	.496	.497	.500	.500	.500	.500	.500	16 points
	2	.257	.261	.270	.282	.289	.293	.298	
	5	.077	.083	.089	.096	.102	.112	.119	

averaging of many curves, some with quite distinctive elbows at different dimensionalities. Thus Figure 4 should not obscure the fact that pure noise

may give the spurious appearance of a true dimensionality.

The importance of these findings to a user of Kruskal's program rests

used to test *a priori* hypotheses about the dimensionality or spatial arrange-

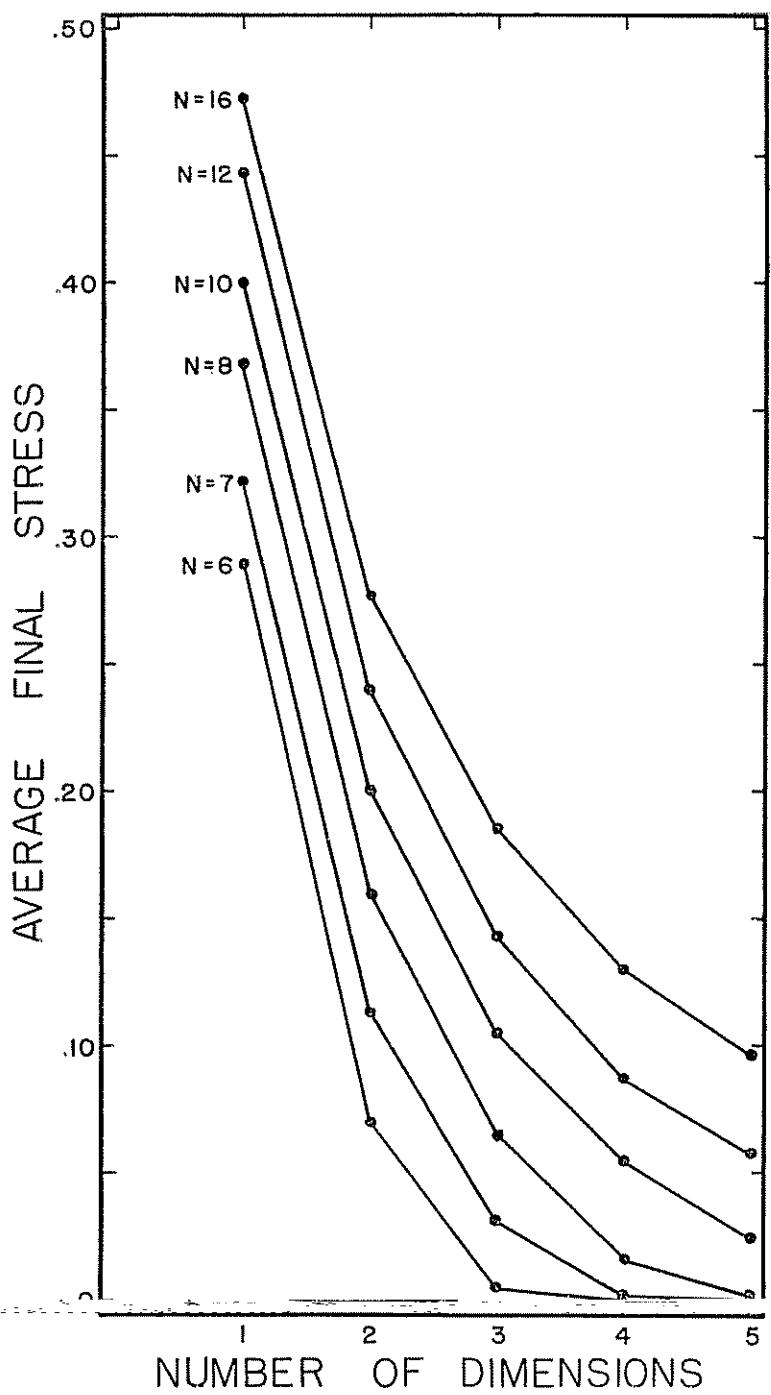


Fig. 4.--Average final stress vs. number of dimensions for 6, 7, 8, 10, 12, and 16 points.

ment of a stimulus set, then the significance criterion is but one of many pieces of evidence that can be used in interpreting results. However, if the scaling procedure is used in an exploratory manner, where there are no preconceptions about the configuration, then any results for only 8 or 9 points in 2 or 3 dimensions cannot be considered very convincing evidence for the existence of structure. For example, from Figure 3 it is evident that with eight points there are about 2 chances in 3 that pure noise could be accounted for in 3 dimensions with a stress less than .075. If the experimenter is at the design stage, then he can use these results as a lower bound on the number of stimuli and will need a significant number of stimulus points to obtain a satisfactory

6. Conclusion

The scaling procedure discussed in this paper is one of a class of new techniques that are so powerful and convenient that they are being used as exploratory devices to see if any structure exists in a set of proximity measures. It is a straightforward procedure and is suitable for a wide variety of multidimensional scaling applications.

REFERENCES

- overview. *Wharton Quarterly*, 1968, 3.
- Huber, D. D. *Dimensionality reduction in psychology*. (In press, 1968, coming).
- Kruskal, J. B. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 1964, 29, 1-27. (a)
- Kruskal, J. B. Nonmetric multidimensional scaling: A numerical method. *Psychometrika*, 1964, 29, 28-42. (b)
- Lings, J. C. An IBM 7090 program for Guttman-Lingoes smallest space analysis. *Behavioral Science*, 1967, 12, 200-202.
- Russell, P. N., & Gregson, P. A. A comparison of intermodal and intramodal methods in multidimensional scaling of three-component taste mixtures. *Australian Journal of Psychology*, 1967, 18, 244-254.
- Shepard, R. N. The analysis of proximities: Multidimensional scaling with an unknown distance function, I. *Psychometrika*, 1962, 27, 125-140. (a)
- Shepard, R. N. The analysis of proximities: Multidimensional scaling with an unknown distance function, II. *Psychometrika*, 1962, 27, 219-264. (b)
- Shepard, R. N. Analysis of proximities as a technique for the study of information processing in man. *Human Factors*, 1963, 5, 33-48.
- Shepard, R. N. Metric structures in ordinal data. *Journal of Mathematical Psychology*, 1966, 3, 287-315.
- Shepard, R. N., & Kelly, R. J. Goodness of fit for nonmetric scaling in Kruskal's nonmetric scaling procedure. *Psychological Bulletin*, 1969, 71, 122-126.
- Torgerson, W. S. Multidimensional scaling of similarity. *Psychometrika*, 1965, 30, 370-303.

Manuscript received: 1/1/68

Revised manuscript received: 9/15/68