



Faculty Working Papers

INTERTECHNIQUE CROSS-VALIDATION
IN CLUSTER ANALYSIS

A. Marvin Roscoe, Jr.
Jagdish N. Sheth and
Welling Howell

#208

College of Commerce and Business Administration
University of Illinois at Urbana-Champaign

FACULTY WORKING PAPERS

College of Commerce and Business Administration

University of Illinois at Urbana-Champaign

September 25, 1974

INTERTECHNIQUE CROSS-VALIDATION
IN CLUSTER ANALYSIS

A. Marvin Roscoe, Jr.
Jagdish N. Sheth and
Welling Howell

#208

Intertechnique Cross-Validation in Cluster Analysis

A. MARVIN ROSCOE, JR.


JAGDISH N. SHETH,

and

WELLING HOWELL *

Clustering methods are often used in marketing research to define homogeneous market segments, and it should be determined in these studies that the derived clusters represent actual clusters. However, replication or external validation is not always practical. An alternative procedure, cross-validation using intertechnique comparisons, is described in a study of geographical market heterogeneity for the telephone industry.

* A. Marvin Roscoe, Jr. and Welling Howell are Marketing Supervisors in the Market Research Section of the Marketing Department of the AT&T Company. Jagdish N. Sheth is I.B.A. Distinguished Professor and Research Professor at the University of Illinois, Urbana - Champaign.



Digitized by the Internet Archive
in 2012 with funding from
University of Illinois Urbana-Champaign

<http://www.archive.org/details/intertechniquecr208rosc>

Cross-validation among techniques seems essential in cluster analysis because most clustering methods tend to be heuristic algorithms instead of analytically optimal solutions. (See Joyce and Channon [6] and Frank and Green [2] for a review of the numerous clustering methods available today). As heuristic algorithms, they have no sampling theory for statistical inferences about the size and the number of clusters. Also, there are no external validation procedures to ensure that the clusters derived from a specific cluster analysis are in reality the true invariant clusters. The potential statistical problem of obtaining artifacts as clusters is further compounded in some procedures which require a priori assumptions about the size and the number of clusters. Although a number of clustering methods perform statistical tests such as the F ratio or Wilks' Lambda based on analysis of variance principles to guard against obtaining random solutions, no procedure exists which will increase the assurance that a nonrandom cluster solution is in fact the true cluster solution.

Because clustering methods are used in marketing research to identify homogeneous market segments for selective marketing efforts, it is critical that the clusters derived from a heuristic algorithm are the true clusters. One procedure to ensure cluster invariance is replication which, however, is not always practical. Another procedure is the common practice in psychometrics of cross-validating the results by external validation. Surprisingly, there are very few studies in which cross-validation has been utilized to insure that the derived clusters are indeed invariant. Although several studies have pointed out the dramatic changes in the cluster structures as a function of data input [4,8] there seems to be only one published study to our knowledge which has examined the question of intertechnique validation of clusters [3].

The objective of this paper is to describe a cross-validation procedure which utilizes intertechnique comparisons of the clustering results. Although the actual study entailed applications of five different clustering techniques, our discussion is limited to two techniques in this paper due to space limitations. A brief description of the large scale research project is provided in which the clustering results were essential to formulating an experimental design for a field experiment.

DESCRIPTION OF THE STUDY

The major research study consisted of a three factorial-64 cell experimentation on survey research methods. The three factors were: first, two different lengths of the questionnaire; second, four different follow-up procedures; and, third, the market heterogeneity of geographical areas of the United States with respect to consumer telephone behavior and socioeconomic-demographic characteristics (see [9]). The levels of the first two factors were predetermined based on theory, prior research and practical implications for the ongoing research on a longitudinal national panel of telephone customers. For the third factor, it was necessary to determine the heterogeneity of the markets by empirical research which utilized clustering methods.

To define the market heterogeneity, profile data on 30,000 residential telephone customers were used for clustering. These customers are part of a longitudinal consumer panel called the Marketing Research Information System which is maintained for the Bell System by AT&T. The panel members are selected based on a multistaged stratified sample in which the first stage of the sampling procedure consists of 100 Revenue Accounting Offices (RAOs) representing the entire Bell System. The profile consists of essentially three types of information about each panel member:

(a) his socioeconomic - demographic status and housing characteristics determined by a survey conducted in early 1970 and matched with the 1970 Census, (b) his monthly telephone behavior broken down into several categories as determined by the industry practice, and (c) an inventory of his telephone equipment including number and types of telephones, and additional services.

Since it was required to empirically investigate the geographical heterogeneity of the markets, an average profile of the residential telephone customers was determined for each of the 86 RAOs for which detailed and complete information was available.

A total of 65 customer descriptors were used to represent the total profile of customers. A list of the variables is shown in Table 1. A factor analysis (principal components) solution with orthogonal Varimax rotation was performed on the data for the following reasons: (a) to reduce the multicollinearity among variables so that the profile consisted of orthogonal factor scores which are geometrically essential to calculate Euclidian distances, (b) to equalize the relative weights of each of the underlying dimensions which could otherwise be easily changed by arbitrary dropping or adding of profile variables, and (c) to standardize the diverse scales of measurement common across the socioeconomic, demographic and telephone information [7]. Ten significant factors were extracted from the analysis which summarized 92 percent of the total variance. A brief description of the factors is provided in Table 2.

The number of significant factors was determined using several criteria, both statistical and judgmental, following the recommendations of Rummel [10]. In addition, the stability of the factor structure was also determined by comparing the results with other data analyses to ensure the invariance of the fundamental dimensionality and structure of the profile data.

The standardized rotated factor scores for each RAO were then utilized to compute Euclidian distances between all combinations of RAOs. The resultant 86 X 86 distance matrix became the input to the clustering procedures.

Due to the following distinct advantages, Johnson's Hierarchical Clustering method [5] was chosen as the primary clustering technique for determining the market heterogeneity. First, it is strictly empirical; second, no prior assumptions are required on the part of the researcher; and third, a hierarchical display is provided of the clusters being formed based on a function minimizing the pairwise distances among entities. While the size of the distance matrix is a limitation of the technique, it was not a problem in our case because of the relatively small number of RAOs to be clustered. Due to the structure of the distance matrix and the presumption of the "ultrametric inequality", [5, p. 248-9] the diameter method was chosen instead of the connectedness method in the BE-HICLUST solutions. The results are diagramed in Figure 1.

While the hierarchical clusters from HICLUST were meaningful and had strong face validity, it was necessary to cross-validate the results by at least one other technique which was essentially similar in its input requirements, analytic strategies and the output format. For this we chose the cluster analysis program developed as part of the BMDP Series which is also a hierarchical clustering routine based on sum of squares distances and the amalgamation principle [1]. In short, BMDP2M amalgamates entities based on the criterion of the smallest distance. Once a cluster is formed, consisting of at least two entities, it calculates the average profile of the cluster and treats it as if it were a new entity which is then clustered with other entities or clusters based on the principle of smallest distances. The process continues until all entities and clusters are hierarchically linked at different levels of distances. The results of the BMDP2M analysis are diagramed in Figure 2.

As can be seen, the two hierarchical clusters are similar in their structure and hierarchy suggesting that there is a good cross-validation between the two analyses. In order to quantitatively assess the degree of congruence between the two hierarchical clusters, two distinct statistical procedures were utilized. The first procedure consisted of calculating the correlation coefficient for the two distributions of distances at which linkages were made between entities or clusters in each hierarchical analysis. Since the number of linkages is not likely to be identical, we have selected the maximum number of links of one technique and the corresponding number of the other technique. The correlation coefficient between the sequential linkage distances is 0.994 which is highly positive indicating extreme closeness of the hierarchical structure of the two cluster analyses.

Another procedure for cross-validation consisted of examining the clusters developed at some specific levels of distances. Based on the plotting of distances at which linkages were made, for the BE-HICLUST results a distance of 5.0 was indicated as a cutoff point due to the natural break in the curve suggesting a clear truncation.

The linkage for the BMDP2M results were also plotted and the natural break in the linkages occurred at 3.1. This was at the point where all the clusters had been formed. After this point the BMDP2M analysis indicated 15 unique entities that were not identified with any of the defined clusters. In order to produce comparable results, the cutoff point for the BE-HICLUST diagram was moved to 3.5 for the cross-validation. The clusters could be identified by their geographical orientation and have been labeled Eastern, Southern, Central and Western. Metropolitan has been used for large urban areas not specifically associated with regional areas. The clusters derived from the two techniques are marked in Figures 1 and 2 and are listed in Table 3.

A total of 17 clusters are displayed in Table 4, consisting of 13 regional clusters (Eastern, Southern, Central and Western), three metropolitan city clusters and the last one representing all the unique RAOs which could not be clustered due to their extreme distances from other RAOs. The cross-tabulation between HICLUST and BMDP2M clustering results indicates that 62 out of 86 RAOs fell on the diagonal of the crosstab matrix which represents a hit of 72 percent correct classifications in terms of intertechnique results. Furthermore, most of the off-diagonal elements generally fall across clusters within the same geographical region. In Table 5, a cross-tabulation at the regional level is provided which shows that 75 out of 86 RAOs could be correctly classified on an intertechnique basis. This represents a hit of 87 percent.

While the two results are quite comparable, there are differences in the example worth noting. The BE-HICLUST algorithm appears to provide a more logical structure to the clusters which are grouped by region as indicated in Figure 2. In addition, the BE-HICLUST method seems to work better where large distances are involved, associating 8 of the 14 unique entities with meaningful clusters. Such differences reinforce the need to use several techniques and to understand the advantages of each especially where the researcher's judgement plays such an important role.

SUMMARY AND CONCLUSIONS

We have pointed out the need for intertechnique cross-validation in cluster analysis due to the heuristic nature of most clustering procedures and the subjective judgements required to interpret the results. In this paper, we have also presented a concrete application of two statistical procedures which enable the researcher to quantitatively measure the congruence of structure and content of clusters across techniques. The first consists of a correlation coefficient index calculated on the distributions of distances at which sequential linkages are made among entities or clusters or both. The second consists of a cross-tabulation of specific clusters derived across two different solutions.

In this paper, the intertechnique cross-validation procedures have been applied with respect to two hierarchical clustering procedures in which the problem was the determination of geographical heterogeneity of markets for the telephone industry. This application considered the general housing and population characteristics along with a complete profile of telephone behavior. However, other uses of the intertechnique cross-validation procedure have been made by the authors for a variety of telephone behavior and markets.

Table 1

LIST OF VARIABLES

Housing	Family
1. Own-rent home	9. Income
2. Type of residence	10. Number in family
3. Number of rooms	11. Average Age
	12. Life cycle
	13. SES status
Mobility	Telephone Service and Equipment
4. Length of residence	
Head of Household	
5. Sex	14. Class of service
6. Age	15. Grade of service
7. Education	16. Number of telephones
8. Occupation	17. Number of vertical services

Billing Items 12 months

18-29	Local service
30-41	Local message
42-53	Intrastate long distance
54-65	Interstate long distance

Table 2

FACTOR DIMENSION LABELS

1. Local service billing	6. Life cycle
2. Local message billing	7. Service and equipment
3. Intrastate long distance	8. Interstate long distance 1 *
4. Family - housing	9. Interstate long distance 2 *
5. Interstate long distance	10. Socioeconomic characteristics

* The two factors for interstate long distance represent different seasonal patterns of calling across geographical areas.

Figure 1

HIERARCHICAL STRUCTURE OF TELEPHONE AREAS - BE-HICLUST

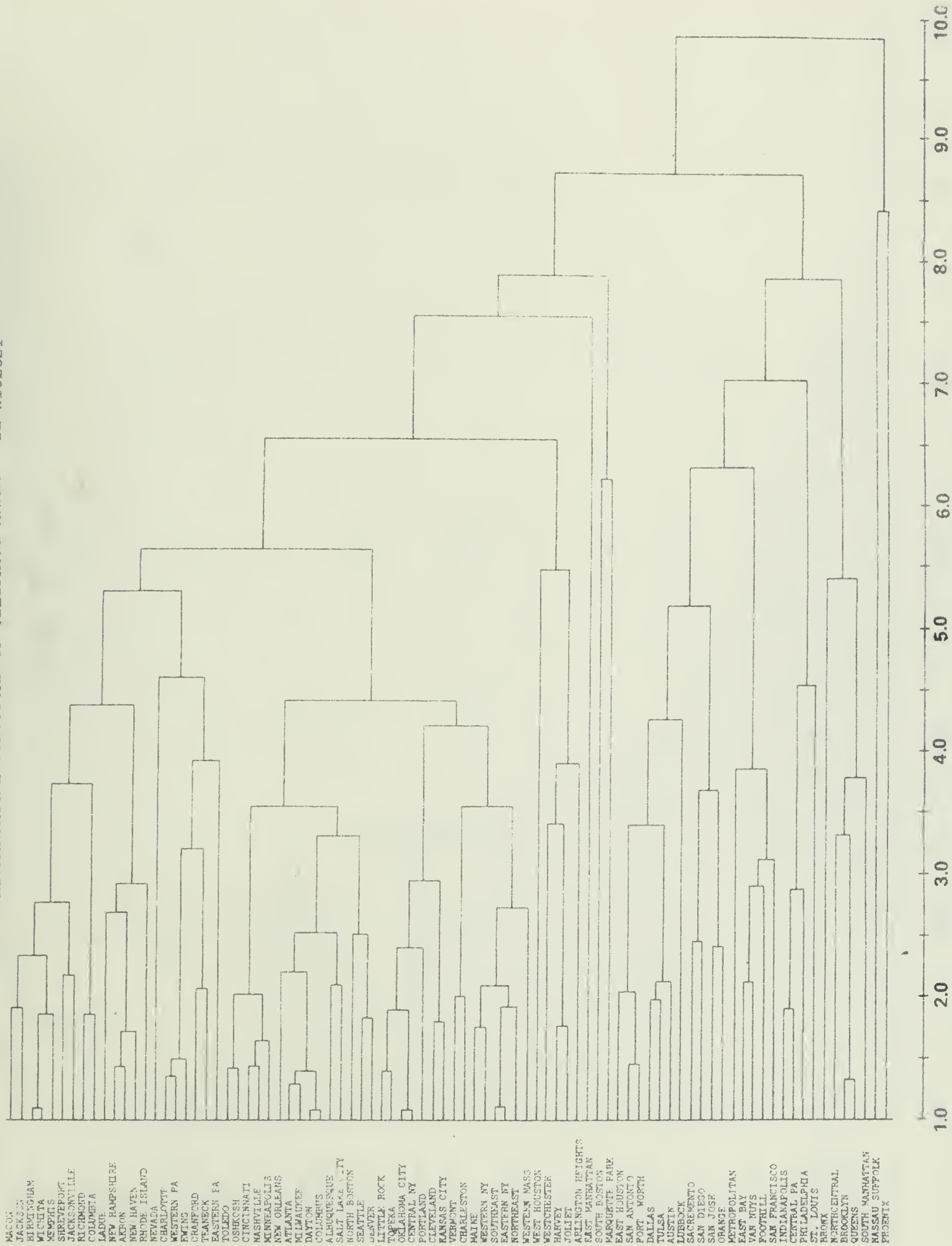
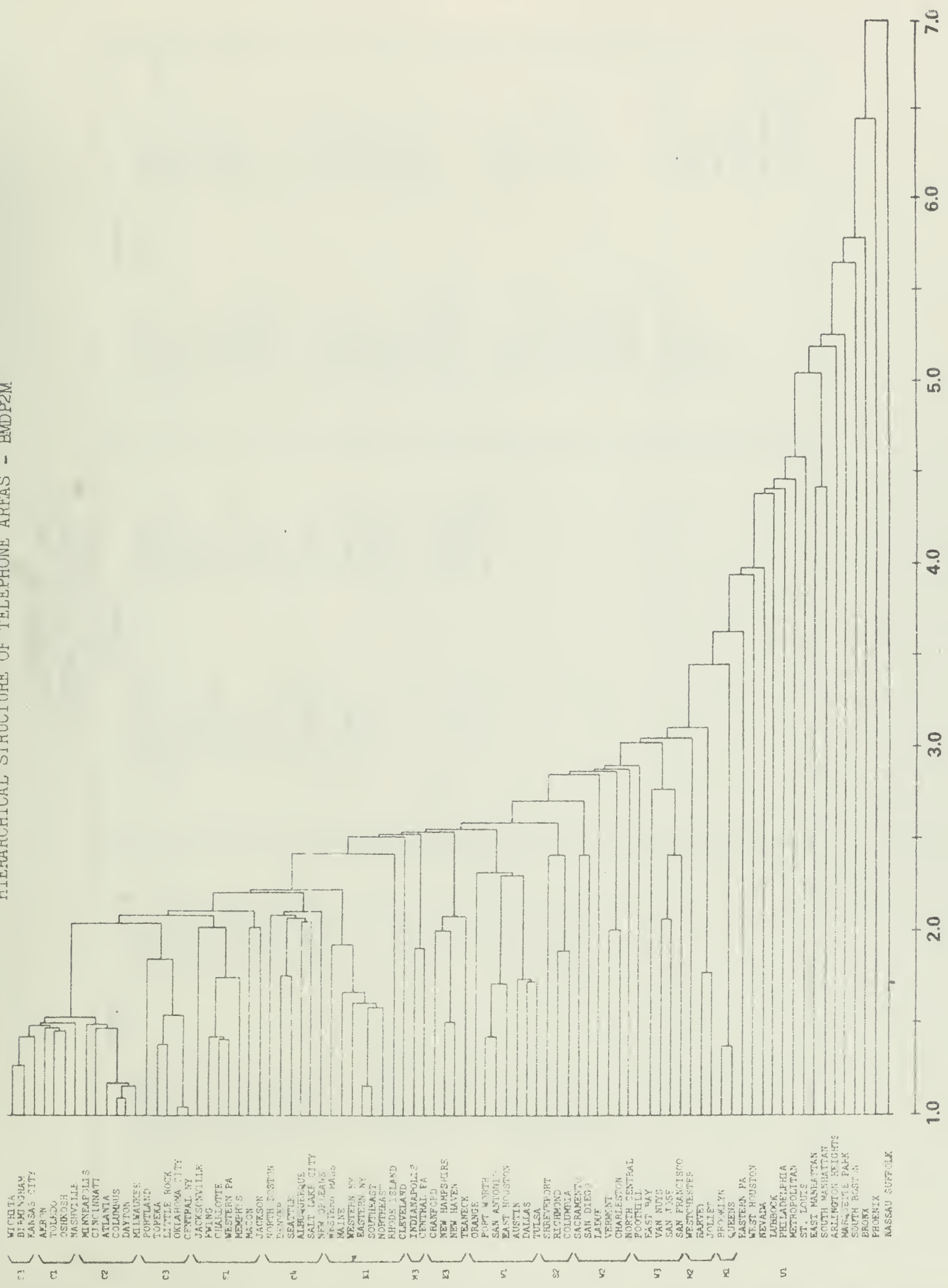


Figure 2

HIERARCHICAL STRUCTURE OF TELEPHONE AREAS - BMDP2M



- WICHITA
- BIRMINGHAM
- KANSAS CITY
- AKRON
- TOLEDO
- OSHKOSH
- NASHVILLE
- MINNEAPOLIS
- CINCINNATI
- ATLANTA
- COLUMBUS
- DAYTON
- MILWAUKEE
- PORTLAND
- TOPEKA
- LITTLE ROCK
- OKLAHOMA CITY
- CENTRAL NY
- JACKSONVILLE
- FWINING
- CHARLOTTE
- WESTERN PA
- MEMPHIS
- MACON
- JACKSON
- NORTH DOSTON
- DENVER
- SEATTLE
- ALBUQUERQUE
- SALT LAKE CITY
- NEW ORLEANS
- WASHINGTON
- MAINE
- WESTERN NY
- EASTERN NY
- SOUTHEAST
- NORTHEAST
- RHODE ISLAND
- CLEVELAND
- INDIANAPOLIS
- CENTRAL PA
- CRANFORD
- NEW HAMPSHIRE
- NEW HAVEN
- TEANECK
- ORANGE
- PORT WORTH
- SAN ANTONIO
- EAST HOUSTON
- AUSTIN
- DALLAS
- TULSA
- SREVEPORT
- RICHMOND
- COLUMBIA
- SACRAMENTO
- SAN DIEGO
- IANUF
- VERMONT
- CHARLESTON
- NORTH CENTRAL
- FORT HILL
- EAST WAY
- VAN NUYS
- SAN JOSE
- SAN FRANCISCO
- WESTCHESTER
- RARITY
- JOLIET
- BROOKLYN
- QUEENS
- EASTERN PA
- WEST HOUSTON
- REVADA
- LURECK
- PHILADELPHIA
- METROPOLITAN
- ST. LOUIS
- EAST MANHATTAN
- SOUTH MANHATTAN
- ARLINGTON HEIGHTS
- MARJEMIE PARK
- SOUTH BOSTON
- BROWK
- PHOENIX
- MASSAU SUFFOLK

1.0 2.0 3.0 4.0 5.0 6.0 7.0

Table 3

LISTING OF CLUSTERS

<u>Eastern</u>		<u>Southern</u>		<u>Central</u>	
BE-HICLUST	BMDP2M	BE-HICLUST	BMDP2M	BE-HICLUST	BMDP2M
11	Western Mass	S1	Jacksonville	C1	Akron
	Maine	2	Knox		Toledo
	Western Va	3	Charlottesville		Columbus
	Western Va	4	Richmond		Cincinnati
	Western Va	5	Richmond		Dayton
	Western Va	6	Richmond		Dayton
	Western Va	7	Richmond		Dayton
	Western Va	8	Richmond		Dayton
	Western Va	9	Richmond		Dayton
	Western Va	10	Richmond		Dayton
	Western Va	11	Richmond		Dayton
	Western Va	12	Richmond		Dayton
	Western Va	13	Richmond		Dayton
	Western Va	14	Richmond		Dayton
	Western Va	15	Richmond		Dayton
	Western Va	16	Richmond		Dayton
	Western Va	17	Richmond		Dayton
	Western Va	18	Richmond		Dayton
	Western Va	19	Richmond		Dayton
	Western Va	20	Richmond		Dayton
	Western Va	21	Richmond		Dayton
	Western Va	22	Richmond		Dayton
	Western Va	23	Richmond		Dayton
	Western Va	24	Richmond		Dayton
	Western Va	25	Richmond		Dayton
	Western Va	26	Richmond		Dayton
	Western Va	27	Richmond		Dayton
	Western Va	28	Richmond		Dayton
	Western Va	29	Richmond		Dayton
	Western Va	30	Richmond		Dayton
	Western Va	31	Richmond		Dayton
	Western Va	32	Richmond		Dayton
	Western Va	33	Richmond		Dayton
	Western Va	34	Richmond		Dayton
	Western Va	35	Richmond		Dayton
	Western Va	36	Richmond		Dayton
	Western Va	37	Richmond		Dayton
	Western Va	38	Richmond		Dayton
	Western Va	39	Richmond		Dayton
	Western Va	40	Richmond		Dayton
	Western Va	41	Richmond		Dayton
	Western Va	42	Richmond		Dayton
	Western Va	43	Richmond		Dayton
	Western Va	44	Richmond		Dayton
	Western Va	45	Richmond		Dayton
	Western Va	46	Richmond		Dayton
	Western Va	47	Richmond		Dayton
	Western Va	48	Richmond		Dayton
	Western Va	49	Richmond		Dayton
	Western Va	50	Richmond		Dayton
	Western Va	51	Richmond		Dayton
	Western Va	52	Richmond		Dayton
	Western Va	53	Richmond		Dayton
	Western Va	54	Richmond		Dayton
	Western Va	55	Richmond		Dayton
	Western Va	56	Richmond		Dayton
	Western Va	57	Richmond		Dayton
	Western Va	58	Richmond		Dayton
	Western Va	59	Richmond		Dayton
	Western Va	60	Richmond		Dayton
	Western Va	61	Richmond		Dayton
	Western Va	62	Richmond		Dayton
	Western Va	63	Richmond		Dayton
	Western Va	64	Richmond		Dayton
	Western Va	65	Richmond		Dayton
	Western Va	66	Richmond		Dayton
	Western Va	67	Richmond		Dayton
	Western Va	68	Richmond		Dayton
	Western Va	69	Richmond		Dayton
	Western Va	70	Richmond		Dayton
	Western Va	71	Richmond		Dayton
	Western Va	72	Richmond		Dayton
	Western Va	73	Richmond		Dayton
	Western Va	74	Richmond		Dayton
	Western Va	75	Richmond		Dayton
	Western Va	76	Richmond		Dayton
	Western Va	77	Richmond		Dayton
	Western Va	78	Richmond		Dayton
	Western Va	79	Richmond		Dayton
	Western Va	80	Richmond		Dayton
	Western Va	81	Richmond		Dayton
	Western Va	82	Richmond		Dayton
	Western Va	83	Richmond		Dayton
	Western Va	84	Richmond		Dayton
	Western Va	85	Richmond		Dayton
	Western Va	86	Richmond		Dayton
	Western Va	87	Richmond		Dayton
	Western Va	88	Richmond		Dayton
	Western Va	89	Richmond		Dayton
	Western Va	90	Richmond		Dayton
	Western Va	91	Richmond		Dayton
	Western Va	92	Richmond		Dayton
	Western Va	93	Richmond		Dayton
	Western Va	94	Richmond		Dayton
	Western Va	95	Richmond		Dayton
	Western Va	96	Richmond		Dayton
	Western Va	97	Richmond		Dayton
	Western Va	98	Richmond		Dayton
	Western Va	99	Richmond		Dayton
	Western Va	100	Richmond		Dayton

Table 3 - Continued

LISTING OF CLUSTERS

<u>Western</u>		<u>Metropolitan</u>		<u>Unique</u>	
BE-HICLUST	BMDP2M	BE-HICLUST	BMDP2M	BE-HICLUST	BMDP2M
MI East	Orange	MI North Central	Brooklyn	Nevada	Eastern Pa
San Antonio	Fort Worth	Brooklyn	Queens	Eastern Pa	West Houston
Fort Worth	San Antonio	Queens		West Houston	Nevada
Dallas	East Houston			Arlington Heights	Lubbock
Dallas	Austin			East Manhattan	Philadelphia
Atlanta	Dallas			South Boston	Metropolitan
Atlanta	Las Vegas	Westchester	Westchester	Marquette Park	St. Louis
Atlanta	Phoenix	Harvey	Harvey	Lubbock	East Manhattan
Atlanta	San Diego	Joliet	Joliet	Metropolitan	South Manhattan
Atlanta	San Diego			St. Louis	Arlington Heights
Atlanta	Las Vegas			Bronx	Marquette Park
Atlanta	Meritt	Indianapolis	Indianapolis	South Manhattan	South Boston
Atlanta	Manhattan	Central Pa	Central Pa	Nassau Suffolk	Bronx
Atlanta	North Central	Philadelphia		Phoenix	Phoenix
Atlanta	Atlanta			Nassau Suffolk	Nassau Suffolk
Atlanta	East Bay				
Atlanta	San Diego				
Atlanta	San Diego				
Atlanta	San Francisco				

REFERENCES

1. Dixon, W.J. "BMD P Series Documentation," Health Sciences Computing Facility. Los Angeles: University of California, 1971.
2. Frank, Ronald B. and Paul E. Green. "Numerical Taxonomy in Marketing Analysis: A Review Article," Journal of Marketing Research, 5 (February 1968), 83-98.
3. Golob, Thomas F., Eugene T. Canty, and Richard L. Gustafson. "Classification of Metropolitan Areas for the Study of New Systems of Arterial Transportation," Paper presented at the 1972 Annual Meeting of the Transportation Research Forum, Denver, Colorado, November 8-10, 1972.
4. Green, Paul E., Ronald E. Frank, and Patrick J. Robinson. "Cluster Analysis in Test Market Selection," Management Science, 13 (April 1967), 387-400.
5. Johnson, Stephen C. "Hierarchical Clustering Schemes," Psychometrika, 32 (September 1967), 241-54.
6. Joyce, Timothy and C. Channon. "Classifying Market Segment Respondents," Applied Statistics, 15 (November 1966), 191-215.
7. Morrison, Donald G. "Measurement Problems in Cluster Analysis" Management Science, 13 (August 1967), B775-80.
8. Neidell, Lester. "Comments on Typology and Cluster Analysis," Paper presented at the AMA Workshop on Multivariate Methods in Marketing, Chicago, Illinois, January 1970.
9. Roscoe, A. Marvin, Dorothy Lang, and Jagdish N. Sheth. "Experimental Effects of Follow-up Methods, Questionnaire Length, and Market Heterogeneity in Mail Surveys," Manuscript submitted for publication, 1974.
10. Rummel, R.J. Applied Factor Analysis. Evanston: Northwestern University Press, 1970, Chapter 15.

UNIVERSITY OF ILLINOIS-URBANA



3 0112 060296743