

AUTHOR COCITATION ANALYSIS USING CUSTOM BIBLIOGRAPHIC DATABASES: AN EXPLORATORY TOOL FOR DIGGING UP REFERENCE DISCIPLINES

Sean, Eom, Department of Accounting & MIS, Southeast Missouri State University, 1 University Plaza, Cape Girardeau, MO, USA, sbeom@semo.edu

Abstract

Researchers in any academic discipline build on each other's and their own previous work. Definitions, topics and concepts are shared. It is necessary to continuously follow up on interesting lines of inquiry. It is also necessary to identify, examine, and trace the intellectual linkage to each other in a given academic field as a basis of assessing the current state of its field to guide future development. Over the past 80 years, the way we count and analyze the intellectual linkage dramatically changed from the early manual transcribing and statistical computation of citation data to computer-based citation data creation and its manipulation. Most citation and cocitation analyses rely on commercial citation databases such as Social Science Citation Index.

This paper introduces an alternative approach to conducting author cocitation analysis (ACA) without relying on commercial citation databases, based on custom bibliographic database and cocitation matrix generation systems specifically developed to use the custom database. The alternative approach overcomes several weaknesses of commercial online data-based ACA research. This guide to an alternative approach to ACA will encourage other researchers to explore the intellectual structures of various MIS fields and guide the future development as well as revealing their reference disciplines.

Keywords: intellectual linkage, author cocitation analysis, reference disciplines, intellectual structure, factor analysis, cluster analysis, multidimensional scaling.

1 INTRODUCTION

A huge body of knowledge existing today is the result of a cumulative research tradition. Researchers build on each other's and their own previous work. Definitions, topics and concepts are shared. It is necessary to continuously follow up on interesting lines of inquiry. It is also necessary to identify, examine, and trace the intellectual linkage to each other in a given academic field as a basis of assessing the current state of its field to guide future development. The intellectual linkages are established through the process of referencing and citation. These intellectual linkages can be systematically examined by means of counting and analyzing the various facets of intellectual activity outputs in the form of written communications.

Over the past 80 years, the way we count and analyze the intellectual linkage dramatically changed from the early manual transcribing and statistical computation of citation data to computer-based citation data creation and its manipulation. The term statistical bibliography was coined by Hulme in 1922 (Hulme 1923) as a research tool for examining the intellectual development and structure of an academic discipline. Since then, we have seen continuous development in the field of bibliometrics. The principal method of bibliometrics is citation analysis through counting and analyzing the citation frequencies. The most important milestone in the development of citation analysis was established by Garfield. He presented an idea for the management of scientific information using a comprehensive citation index in 1955 and three years later founded the Institute for Scientific Information (ISI)(Garfield 1955). For a detailed description of theory and application of citation indexing, see (Garfield 1979). As of November 2003, the company carried about 70 different products including [Web of Science](#), [Science Citation Index](#), etc. A citation index is a listing of all referenced or cited source items published in a given time span associated with the citing articles. These citation index files are online bibliographic databases accessible only through online-based information service companies such as Dialog®, Profound®, DataStar™, [Questel/Orbit](#), etc.

This paper introduces an alternative approach to conducting author cocitation analysis (ACA) without relying on commercial citation databases, based on custom bibliographic database and cocitation matrix generation systems specifically developed to use the custom database.

The next section briefly discusses the basic concepts of bibliometrics and author cocitation analysis along with its assumptions, purposes, benefits, limitations, and criticisms. The following section overviews ACA steps. They include selection of authors, retrieval/generation of paired author cocitation frequencies, preparing inputs to the SAS system, multivariate statistical analyses of author cocitation matrix, and validations/interpretation of SAS outputs. The input to the SAS system is cocitation frequency matrix. The data matrix is processed by three multivariate procedures (factor, cluster, and multidimensional scaling).

2 BASIC CONCEPTS OF AUTHOR COCITATION ANALYSIS

The term statistical bibliography was coined by E. Wyndham Hulme in 1922 (Hulme 1923). The purposes of statistical bibliography are:

1. to shed light on the processes of written communication and of the nature and course of development of a discipline (in so far as this is displayed through written communication), by means of counting and analyzing the various facets of written communications (Prichard 1969).
2. the assembling and interpretation of statistics relating to books and periodicals ... to demonstrate historical movements, to determine the national or universal research use of books and journals, and to ascertain in many local situations the general use of books and journals (Raisig 1962).

Citation analysis is often used to determine the most influential scholars, publications, or universities in a particular discipline by counting the frequency of citations received by *individual* units of analysis (authors, publications, etc.) over a period of time from a particular set of citing documents. However, citation analysis cannot establish relationships among units of analysis. ACA is the principal bibliometric tool to establish *relationships* among authors in an academic field and thus identify subspecialties of a field and how closely each subgroup is related to each of the other subgroups. By establishing relationships among authors, ACA provides a basis of revealing the intellectual structure of literature and defining the principal subject (major area of subspecialties in an academic discipline and their contributing disciplines) through the empirical consensus of numerous authors in an academic discipline.

2.1 Author Cocitation Analysis

There are two primary types of cocitation analysis to map the intellectual structure of an academic field: document cocitation analysis and author cocitation analysis (ACA). Document cocitation analysis involves the analysis of a set of selected documents (e.g., journal articles, books, proceedings, etc.) in terms of which pairs of documents are cited together. Author cocitation analysis, introduced in 1981, is a more general approach to identify, examine, and trace the intellectual structure of an academic discipline by counting the frequency with which any work of an author is cited to any work by another author in the references of citing documents (Bayer et al. 1990).

The cocitation of authors occurs when a citing paper cites any work of authors in reference lists. Many information scientists and author cocitation analysis researchers define an author as "a body of writings by a person" or "a body of contributions by a person." The term "contributions" may be better since it can include any type of contribution that can be cited as a reference such as speeches delivered at professional meetings, personal communications including conversation and letters, and other media. These different uses of terms are related to citation databases used in the study. Most commercial citation databases and software access only the first author, regardless of the number of multiple authors, when retrieving author cocitation counts. This has been the critical weakness of using the commercial citation databases and software. However, this paper is based on the bibliographic database I have created, which includes all contributions such as speeches delivered at various meetings and software we have developed that can access all multiple authors. With custom-built bibliographic databases, and the bottom-up approach of the selection of author sets, ACA becomes an exploratory tool for digging up the roots (reference disciplines), locating the trunk (foundations), and sifting through the branches (subspecialties) of a tree (an academic discipline). The critical element that makes ACA an exploratory tool is the custom bibliographic databases and the author selection method of screening the entire databases to finalize the author set for ACA analysis.

3 OVERVIEW OF AUTHOR COCITATION ANALYSIS STEP

ACA consists of the assembling and interpretation of bibliographical statistics taken from the cited references which are taken from the selected citing articles. (See, Figure 1).

3.1 Selection of Authors/Cocitation Thresholds

An important purpose of ACA is an overall examination of the intellectual structure of an academic discipline. Therefore, it is critical to establish a diversified list of authors. Basically, two approaches are available: starting a predetermined list of authors in a given field (the subjective approach), or compiling a set of authors from scratch (the objective approach).

The first approach starts with a pre-determined list of authors, which is compiled by personal knowledge, consultation with researchers in the area to be studied, conducting surveys, using

directories, organizational membership, etc. Compiling a predetermined list of authors inevitably involves subjective judgments. This approach can be efficient in that no lengthy time is spent to finalize a list of authors for further analysis. A critical weakness of this approach is that it may often fail to identify emerging scholars in a given area of an academic discipline. The majority of previous

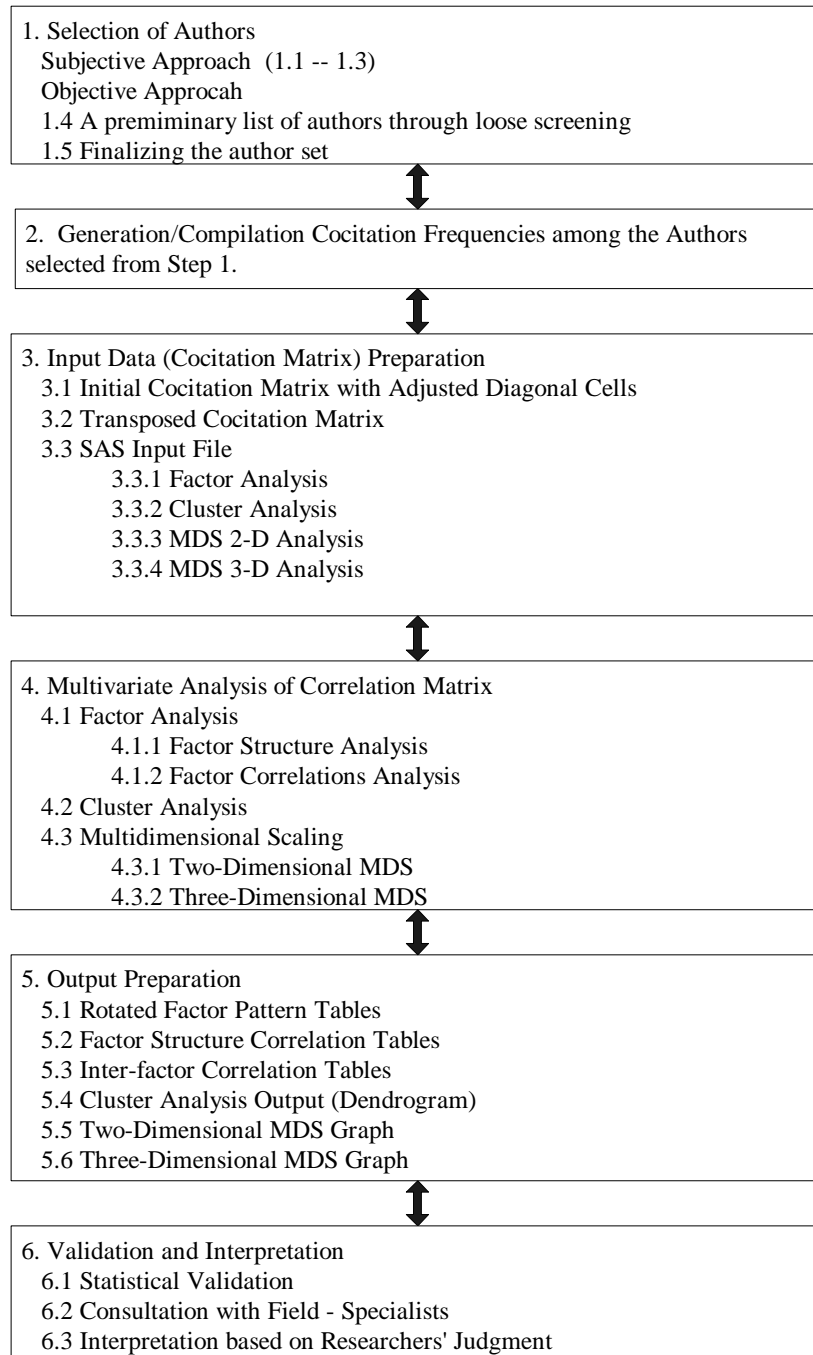


Figure.1 Author Cocitation Analysis Steps

research in this area has used this subjective approach (Culnan 1986; McCain 1986). The reason for doing so was not because this approach is superior but because most of the previous research used commercial online bibliographic databases to retrieve the cocitation frequency matrix.

The objective approach does not start with a predefined list of authors at all. It is relatively manageable to screen the whole custom-built databases to select a list of authors for further analysis. ACA, using custom databases, allows researchers to unobtrusively select authors for the study. Emerging scholars are more likely to be included in the selection process, unlike the subjective approach discussed earlier.

The other important advantage of using custom-built databases includes the identification of reference disciplines. *Citation* counts of the individual authors in the preliminary list are used to finalize the list of authors. This step filters the preliminary list compiled from the previous step further into a set of authors based on author *cocitation* counts (frequencies). Although it rarely happens, theoretically it is possible for any authors with higher citation counts to have very low cocitation counts with other authors. In this case, these authors may appear in the final author set. However, they will carry lower factor loadings that can be interpreted as insignificant authors to the formation of the intellectual structure of the academic discipline under study. To accurately and objectively examine the intellectual structure of a discipline, personal judgment must be avoided in selecting authors by objectively counting the frequency of each name from the data bases.

Using in-house databases, author selection criteria must be established in regard to the number of citation received. A higher number of threshold means a smaller number of authors for further analysis. The optimal number of authors is primarily dependent on the number of cited database records. In addition, there may be some differences in the citation behavior across the different academic disciplines. There are no quantitative tools that can be blindly applied in deciding the number of authors. A study of author cocitation of a journal in consumer research over the 15 year period used 4 citation as selection criteria to compile the list of authors (Hoffman et al. 1993). Other studies of ACA to map the intellectual structures of decision support systems, the adjusted diagonal cell values of 25 or more were used (Eom 2003). Once the threshold values are used to select the author list, it is important to apply the same criteria consistently to the subsequent studies to be followed to trace the changing structures of a discipline. Due to the possible instability of small cocitation counts, author cocitation analysis researchers introduced several ad hoc criteria for further screening a large pool of candidate authors to finalize a list of authors (See (Eom 2003)).

Regardless of the nature of bibliographic databases (commercial vs. custom-built), determining the threshold cocitation rate is not the result of a structured process; rather, it is an unstructured process requiring the investigator's personal judgments. An exact quantitative basis for deciding the threshold cocitation rate has not been developed. Lowering the threshold in general increases the number of authors to be included in a study, which in turn may or may not change the number of meaningful factors in the study. Also, it is important to point out that cocitation thresholds themselves, as sole connection criteria, are suspect in a highly multidisciplinary area. One should look at the overall connectedness and the focused cocitation counts as well.

3.2 Generation Of Cocited Author Counts

Cocitation counts can be either retrieved from commercial online bibliographical databases such as Science Citation Index, Social Sciences Citation Index, and Arts and Humanities Citation Indexes, or generated from the custom-built bibliographical databases. If researchers use the commercial on-line databases, cocitation frequencies will be retrieved using the query command. Several examples of those commands are given in McCain (1990).

This paper introduces an alternative approach to ACA research -- the generation of cocited author counts using a custom-built bibliographic database and in-house cocitation count generation systems.

FoxBase database management systems are used to enter the bibliographic records. It computes author cocitation frequencies between any pair of all (primary and non-primary) authors under study. The author cocitation frequency generation system enables the users to overcome the problem with the Institute for Scientific Information (ISI) databases which code only the first author of a cited work. The cocitation matrix generation system we developed gives access to cited coauthors as well as first authors.

3.3 Conversion of Raw Cocitation Matrix

In ACA, cocitation frequency (cocited author counts) is the prime input data. The cocitation count generation system produces a raw author cocitation frequency table (Table 2). It is not suitable to use the raw citation matrix in triangle shape an input to the SAS system for ACA study. The next step is adjusting diagonal cell values in the raw cocitation matrix created. Raw cocitation frequencies of row and column authors fill the off-diagonal cells. The off-diagonal cell value is the total number of the cocitation count between these two authors. The term "author" in author cocitation analysis is neither an individual nor individuals. It refers to a body of writings by a person.

| Ref. #1 | Ref. #2 | Ref. #3 |
|----------|-----------|----------|
| Ackoff | Ackoff | Ackoff |
| Bonczek | Ackoff | Ackoff |
| Bonczek | Applegate | Blanning |
| Blanning | Applegate | |
| Blanning | Whinston | |
| Blanning | | |
| Whinston | | |

Table 1 Reference of Citing Papers

| | Ackoff | Applegate | Bonczek | Blanning | Whinston |
|-----------|--------|-----------|---------|----------|----------|
| Ackoff | 3 | | | | |
| Applegate | 1 | 1 | | | |
| Bonczek | 1 | 0 | 1 | | |
| Blanning | 2 | 0 | 2 | 2 | |
| Whinston | 2 | 1 | 1 | 2 | 2 |

Table 2 Sample Cocitation Matrix

The values of the diagonal cells are computed using the adjusted value approach, taking the three highest intersections for each author and dividing the sum by two. The value produced by the cocitation count generation system is to be replaced by the adjusted diagonal cell value. The rationale for using this adjusted value can be found in McCain (1990). In addition to the adjusted diagonal cell value approach, McCain (1990 p.435) discusses the second and third approaches--substituting the diagonal value with the highest off-diagonal cocitation counts for each author and treating the diagonal cell values as missing data. Her initial results indicate little difference between these two. We also found little difference between the approach used here and the second approach of using the highest off-diagonal cell value. In order to prepare a SAS input file, a diagonal cell value adjusted cocitation matrix (Table 3) needs further transformation (Table 4).

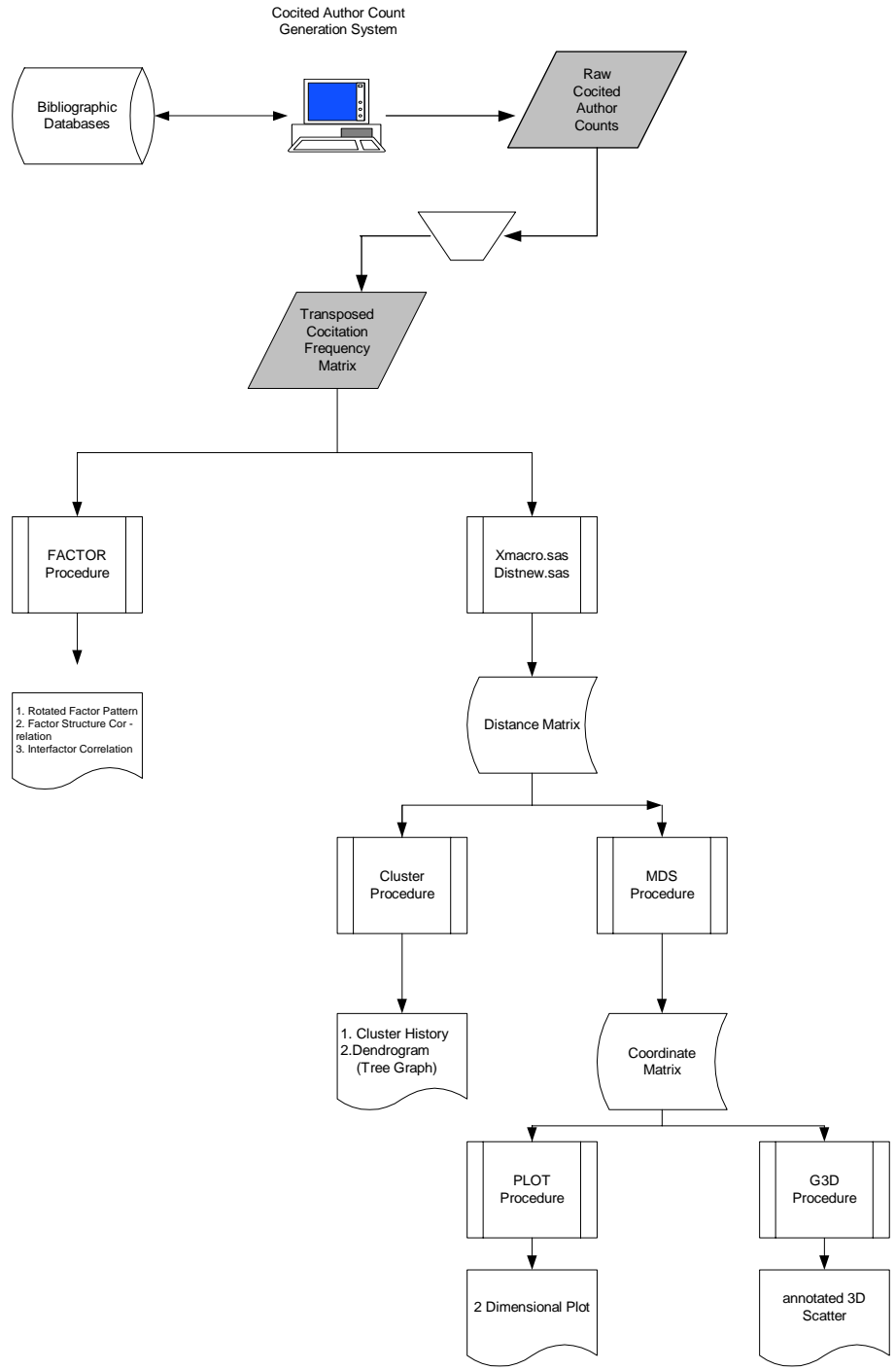


Figure2 Cocitation Frequency Matrix Preparation and SAS Steps in ACA

| | Ackoff | Applegate | Bonczek | Blanning | Whinston |
|----------------------|--------|-----------|---------|----------|----------|
| Ackoff | 2.5 | | | | |
| Applegate | 1 | 1 | | | |
| Bonczek | 1 | 0 | 2 | | |
| Blanning | 2 | 0 | 2 | 3 | |
| Whinston | 2 | 1 | 1 | 2 | 2.5 |
| largest value | 2 | 1 | 2 | 2 | 2 |
| Second largest value | 2 | 1 | 1 | 2 | 2 |
| Third largest value | 1 | 0 | 1 | 2 | 1 |
| Sum | 5 | 2 | 4 | 6 | 5 |
| Adj. diagonal value | 2.5 | 1 | 2 | 3 | 2.5 |

Table 3 Sample Cocitation Matrix (Diagonal Value Adjusted)

| | Ackoff | Applegate | Bonczek | Blanning | Whinston |
|-----------|--------|-----------|---------|----------|----------|
| Ackoff | 2.5 | 1 | 1 | 2 | 2 |
| Applegate | 1 | 1 | 0 | 0 | 1 |
| Bonczek | 1 | 0 | 2 | 2 | 1 |
| Blanning | 2 | 0 | 2 | 3 | 2 |
| Whinston | 2 | 1 | 1 | 2 | 2.5 |

Table 4 ransposed Cocitation Matrix

4 OVERVIEW OF INPUT, PROCEDURES, AND OUTPUTS

As Figure 2 shows, ACA data analysis requires the 7 different procedures of the SAS system. Each procedure requires inputs to produce outputs. It is extremely important for ACA researchers to prepare input files to these procedures in a most efficient manner. There is one necessary input (cocitation frequency matrix), seven procedures (Factor, Xmacro, distnew, cluster, MDS, PLOT, and G3D), and several outputs from 3 multivariate statistical procedures. All three techniques used in the ACA aim at grouping/classifying all variables into several subgroups with common underlying hidden structures, characteristics and/or attributes. The hidden structures/characteristics/attributes are given different terms: factors in factor analysis, clusters in cluster analysis, and dimensions in multidimensional scaling. Although all three techniques seek to summarize/simplify a large number of variables, there are some distinctive differences among these techniques. Factor procedures produce three outputs--rotated factor pattern, factor structure correlation, and interfactor correlation. Cluster procedures generate cluster history and dendrogram (tree graph). Multidimensional scaling procedures result in iteration history, convergence status, coordinate matrix of each author (first column), and configuration of each author on multidimensional spaces (dimension 1-column 2, dimension 2-column 3, dimension 3-column 4, etc).

4.1 The Factor Procedure

The transposed, diagonal-cell-adjusted cocitation matrix of authors (Table 4) can be analyzed by the factor analysis program in most popular statistical software such as SAS (statistical analysis systems) or SPSS to ascertain the underlying intellectual structure of an academic discipline. Table 5 is converted from table 4 as an input file to the SAS system. The objective of this analysis is to group (condense) a large number of selected variables (authors) into a smaller set of composite dimensions (factors) representing research subspecialties/subdisciplines and contributing disciplines of a discipline. In doing so, each variable (author) is viewed as a dependent variable that is a function of a set of latent (underlying) factors.

```

DATA table5;
INPUT @3 Ackoff Applegate Bonczek Blanning Whinston;
CARDS;
1      2.5    1      1      2      2
2      1      1      0      0      1
3      1      0      2      2      1
4      2      0      2      3      2
5      2      1      1      2      2.5
;
PROC FACTOR METHOD=PRINCIPAL MINEIGEN=1 ROTATE=PROMAX;
RUN;

```

Table 5 PROC FACTOR SAS Program with Author's Names

4.2 The Cluster Procedure

Cluster analysis is a data reduction technique for grouping various entities (individuals, variables, objects) into clusters so that the entities in the same cluster are similar with respect to some predetermined selection criteria (Everitt 1980; Hair et al. 1992). Therefore, it is necessary to convert the raw cocitation frequency matrix into a measure of similarity or distance. To do so, SAS institute has developed the DISTANCE macro for computing various measures of distance, dissimilarity, or similarity between the observations of a SAS data set. The first section of this chapter explains the creation of a distance matrix, which is the input to the cluster procedure.

4.3 Mutidimensional Scaling

Multidimensional scaling (MDS) is a class of multivariate statistical techniques/procedures to produce two or three dimensional pictures of data (geometric configuration of points) using proximities among any kind of objects as input. Proximity data consist of one or more square symmetric or asymmetric matrices of similarities or dissimilarities between objects or stimuli (Kruskal et al. 1978, pp. 7-11). The MDS outputs consist of a spatial representation of data which shows underlying relationships on a two or three dimensional map. The MDS map helps clarify relationships visually using the ratio of distances on a map to corresponding data values such as a map of a country showing cities. The magnitude of number indicates how similar/dissimilar two objects are.

The purposes of MDS are to help researchers identify the "hidden structures" in the data and visualize relationships among/within the hidden structures to give clearer explanations of these relationships to others (Hair et al. 1987; Kruskal et al. 1990). Three SAS procedures (MDS, PLOT, and G3D) are necessary to convert the author cocitation frequency matrix to two or three dimensional pictures of data.

The distance matrix produced earlier by using xmacro.sas and distnew.sas programs should be converted to a coordinate matrix. The coordinate matrix is used to produce two-dimensional plots and annotated three-dimensional scatter diagrams. A distance matrix is the input to the multidimensional scaling procedure, PROC MDS, of the SAS system (release 8.0). The PLOT and G3D procedures process the coordinate matrix to visualize the similarity and dissimilarity within each group of an academic discipline as well as the similarity and dissimilarity among the various subspecialties within an academic discipline. The annotate facility in the SAS system produces figures with the name of the author on each data point.

5 CONCLUSION

This paper introduced an alternative approach to conduct author cocitation analysis (ACA) without relying on commercial citation databases, based on custom bibliographic database and cocitation matrix generation systems specifically developed to use the custom database. The alternative approach overcomes several weaknesses of commercial online data-based ACA research.

First, the alternative approach introduced here has the capability to access the non-primary authors of cited references. The non-primary authors refer to all authors other than the first author. The inability to access non-primary authors is a critical shortcoming of ACA research utilizing the commercial databases. Theoretically, the contributions made by non-primary authors must be counted when examining the intellectual structure of an academic discipline.

Second, strict criteria can be applied to the selection of citing articles. A researcher does not always write articles in a specialized field throughout his/her lifetime. Research interests can shift from one subspecialty area to other areas within an academic discipline. Custom bibliographic databases can be built to only include writings in a specific field.

Third, the alternative approach introduced here can be a more effective tool for identifying the intellectual structure of an academic field more accurately as well as its reference disciplines. All previous ACA studies except the ones conducted by Eom and his colleagues (Eom 2003) failed to identify the reference disciplines of an academic field. The method used for the selection of authors for ACA was the reason for failure. The method starts a predetermined list of authors selected by subjective judgments of researchers. It is impractical for ACA researchers to include all authors in the reference disciplines of an academic field prior to conducting ACA analysis. With the approach introduced in this paper, ACA becomes an exploratory tool. It can dig up the roots (reference disciplines), locate the trunk (foundations of an academic discipline), and sift through branches (subspecialties) of a tree (an academic discipline).

Fourth, many problems can also arise in relation to the sources of citation data and the mechanics of deriving citations from existing citation indexes. The problems may stem from multiple authorship, self-citations, homographs, synonyms, the unification problems, etc. (Lindsey 1980; Long 1980; Smith 1981). Use of SSI and SCI can raise a potential problem since these sources can exhibit English language bias (Baker 1990). Use of custom databases and the cocitation matrix generation system we developed can eliminate many of the problems discussed above such as multiple authorship, homographs, synonyms, etc.

References

- Baker, D.R. "Citation Analysis: A Methodological Review," *Social Work Research & Abstracts* (26:3), September 1990, pp 3-10.
- Bayer, A.E., Smart, J.C., and McLaughlin, G.W. "Mapping Intellectual Structure of a Scientific Subfield through Author Cocitations," *Journal of the American Society for Information Science* (41:6), September 1990, pp 444-452.
- Culnan, M.J. "The Intellectual Development of Management Information Systems, 1972-1982: A Co-Citation Analysis," *Management Science* (32:2), February 1986, pp 156-172.
- Eom, S.B. *Author Cocitation Analysis Using Custom Bibliographic Databases--an Introduction to the SAS Approach*, The Edwin Mellen Press, Lewiston, New York, 2003, p. 216.
- Everitt, B.S. *Cluster Analysis*, Heinemann Educational Books Ltd., London, 1980.
- Garfield, E. "Citation Indexes for Science," *Science* (122) 1955, pp 108-111.
- Garfield, E. *Citation Indexing: Its Theory and Application in Science, Technology, and Humanities*, Wiley, New York, 1979.
- Hair, J.F., Jr., Anderson, R.E., Tatham, R., and Black, W.C. *Multivariate Data Analysis with Readings*, (3rd ed. ed.), Macmillan Publishing, New York, 1992.

- Hair, J.F., Jr., Anderson, R.E., and Tatham, R.L. *Multivariate Data Analysis with Readings*, (2nd ed.), Macmillan Publishing Company, New York, 1987.
- Hoffman, D.L., and Holbrook, M.B. "The Intellectual Structure of Consumer Research: A Bibliometric Study of Author Cocitations in the First 15 Years of the Journal of Consumer Research," *Journal of Consumer Research* (19), March 1993, pp 505-517.
- Hulme, E.W. *Statistical Bibliography in Relation to the Growth of Modern Civilization*, Grafton, London, 1923.
- Kruskal, J.B., and Wish, M. *Multidimensional Scaling*, Sage University Paper Series on Quantitative Applications in the Social Sciences, Sage Publications, Beverly Hills, CA, 1978.
- Kruskal, J.B., and Wish, M. *Multidimensional Scaling*, Sage University Paper Series on Quantitative Applications in the Social Sciences, 07-011. Sage Publications, Beverly Hills and London, 1990.
- Lindsey, D. "Production and Citation Measures in the Sociology of Science: The Problem of Multiple Authorship," *Social Studies of Science* (10) 1980, pp 145-162.
- Long, J.S., et al. "The Problem of Junior-Authored Papers in Constructing Citation Counts," *Social Studies of Science* (10), May 1980, pp 127-143.
- McCain, K.W. "Cocited Author Mapping as a Valid Representation of Intellectual Structure," *Journal of the American Society for Information Science* (37:3) 1986, pp 111-122.
- McCain, K.W. "Mapping Authors in Intellectual Space: A Technical Overview," *Journal of the American Society for Information Science* (41:6), September 1990, pp 351-359.
- Prichard, A. "Statistical Bibliography or Bibliometrics?," *Journal of Documentation* (25:4), December 1969, pp 348-349.
- Raisig, L.M. "Statistical Bibliography in the Health Science," *Bulletin of Medical Library Association* (50:3), July 1962, pp 45-461.
- Smith, L.C. "Citation Analysis," *Library Trends* (30:1), Summer 1981, pp 83-106.