

# "REFERATION" DATABASE FOR DOCUMENT INFORMATION ANALYSIS

*Isao Asai*

*University of Osaka Prefecture, Osaka, Japan*

**Abstract.** Referations in a given document is defined as the join of citations (citing documents), references (cited documents) and document itself. It represents the relationships between a given document and the others, and produces the more valuable measures between documents than references or citations. Referation file of document collection in a specific topic can be easily obtained in microcomputers using reference data. Referation database, that is, document file with referation file, is a kind of knowledge one structured by relations. This paper gives a brief outline of some functions based on referations for the content and structure of bibliographic information.

## INTRODUCTION

It is well-known that document database with references, introduced by Garfield[1], has been contributed to the progress of information science and technology. The characteristics of the database is that the attribute of document has references(set of cited documents) and citations(set of citing documents, inverted file of references) instead of descriptors and abstract. At present, we are able to access SCISEARCH and SOCIAL SCISEARCH in online.

To find out the other usage of SCI database, a decade had elapsed since the first issue of the Science Citation Index. Co-citation proposed by Small[2] is a measure for automatically finding out the subject relatedness between documents that share one or more citations. Many of citation analysis using SCI database has been done for studying the science of science[3].

For the influence of valuable citation index, references had little part in SCI database. However, an analysis was recently done using references in SCI database[4]. It is based on similarity measure of documents, that is, bibliographic coupling by Kessler[5]. The revival of bibliographic coupling is quite proper under the nature of references.

On the usage of SCI database, the

following four points are summarized:

- 1) citation search  
cumulative citations of first author,
- 2) reference search
- 3) citation analysis  
co-citation between documents passed over five years,
- 4) reference analysis  
bibliographic coupling between documents less than five years.

The data structure of references in SCI database is made of first author, abbreviated journal name, and volume, page and year. The data structure of citations inverted is the same as that of references. This must be considered as an unavoidable consequences of the expanding amount of documents.

However, the usage of SCI database has some limitations for the characteristics of the data structure. There are, for example,

- 1) difficulty of identification of document,
- 2) insufficiency of document information,
- 3) separation of references and citations.

In particular, the realization of identifying document will enhance the value of SCI database.

This paper deals with a small scale of database with references. It mainly contains,

- 1) representation of reference data,
- 2) concept of referations,
- 3) new information retrieval system based on referations,
- 4) some functions of referation database.

Using microcomputers with higher performance, lower cost, and larger memory, I have been developed an user-friendly system designed for scientists and librarians[6-8]. This gives an outline of this system and some functions of referation database.

## REPRESENTATION OF REFERENCE DATA

It is very difficult to convert the reference data at the end of document into computer readable format. Although we are able to entry all bibliographic items of references, the

redundancy of data increases, and the problem of identification becomes impossible to solve. For extracting the professional information from document database, it is very important to identify document.

SCI database think out of decreasing the entry of reference data. The length of data is 30-40 bytes. It is considered to be a special coding system. Each code is marked when new document is collected. It is noticed that the code is not known in advance. Because of long bytes of code, it is difficult to identify and operate document.

If each document has a fixed code, many problems described previously will be solved. But it is impossible to mark an code corresponding to each document in all of the world. A feasible approach of coding is to restrict within documents of narrow limits in a specific field.

The procedure for coding document in a specific topic is as follows:

- 1) collection of documents in a specific topic,
- 2) coding of document,
- 3) selection of documents related to the topic from references,
- 4) coding of undefined document,
- 5) representation of reference data by code.

Therefore, a reference file is made under the following conditions:

- 1) reference file is composed of document code,
- 2) reference file is only the document cataloged to document file.

As the bibliographic items corresponding to document code are stored to document file, the entry of document code equals to input all bibliographic items.

A very easy method for the entry of reference data is developed. It takes two steps of procedure. First step, we search and display the list of document from document file by assisting of author's index. And second step, we catalog reference data and display the bibliographic items by inputting the corresponding code. These are executed on multi screen.

Under computer algorithm, citation file is obtained by inverting reference file. As reference data is made of document code, this inverted procedure is very quickly.

#### DEFINITION OF REFERATIONS

In a given document, it is possible to integrate citations, document itself,

and references. This new concept is named "referations"[6,8]. Under computer algorithm, referation file is obtained by merging three files, that is, citation file, document file, and reference file.

Figure 1 shows an example of referations of document number 857 published in 1981 having four documents of citations and five documents of references. In this case, the number of referations is ten. It is considered that referations is an important attribute of document.

```

+-----+-----+-----+-----+-----+
I DocN I Yr I           I           I
+-----+-----+-----+-----+-----+
I 1199 I 85 I           I           I
I 1193 I 85 I Citations I Referations I
I 1110 I 83 I           I           I
I 1097 I 83 I           I = Citations I
+-----+-----+-----+-----+
I 857 I 81 I Document I + Document I
+-----+-----+-----+-----+
I 863 I 78 I           I + References I
I 547 I 77 I           I           I
I 104 I 74 I References I           I
I 121 I 65 I           I           I
I 120 I 63 I           I           I
+-----+-----+-----+-----+

```

Figure 1 Example of Referations of Document Number 857

Now, this attempts to describe the concept of referations mathematically. Let  $n$  denote the total number of documents in a specific topic, and let  $R$  denote a matrix of reference file. Reference matrix  $R=(r_{ij})$  is a square matrix of  $n$  rows and  $n$  columns. The elements  $r_{ij}$  is defined by the relation

$$r_{ij} = 1 \text{ if document } i \text{ cites document } j$$

$$r_{ij} = 0 \text{ otherwise.}$$

Then citation matrix  $C$  is defined as follows.

$$C = (c_{ij}) = R^t$$

in where  $R^t$  denotes the transpose of the matrix  $R$ . Let  $E$  denote an unit matrix which each diagonal element is set equal to 1. Referation matrix  $F$  is therefore

$$F = (f_{ij}) = C + E + R.$$

This referation matrix is a square symmetric matrix of  $n$  rows and  $n$  columns.

#### REFERATION DATABASE

A document file has the table type of data structure in which a record

describes a single document and the fields describes the attributes of document. As the referentions of document is one of the attributes, it is natural that referation file is taken in document file. But in this paper, referentions does not take in document file. It is assumed here that referation file is an accompanying file to document file. Document file with referation file is named for referation database.

Document File(Table Data)				Referation File			
I	No	I	Author Title	I	No	I	Referat.
I	1	I	Luhn,H Keyword-in-con	I	1	I	1,3,7
I	2	I	Knox,D Effective Sear	I	2	I	2,5,9
I	3	I	Cuadra Annual Review	I	3	I	1,3,6
I	4	I	Lancas Information Re	I	4	I	4,8,11
I	5	I	Garfie The Permuterm	I	5	I	2,5,10,12
I	6	I	Vicker Classification	I	6	I	3,6,7
I	7	I	Borko, Indexing Conce	I	7	I	1,6,7,12
I	8	I	Meadow The Scientific	I	8	I	4,8,9,10
I	9	I	Bradfo Documentation	I	9	I	2,8,9
I	10	I	Garvey Communication	I	10	I	5,8,10
I	11	I	Cawkel The Paperless	I	11	I	4,11
I	12	I	Dewey Dewey Decimal	I	12	I	5,7,12

Figure 2 is an example of referation database. Each record in document file contains the document number, author name, and title to a single document. And each record in referation file is made of record numbers related to a single document. For example, record number three in referation file is related to record one, three, and six.

Figure 2 Example of Referation Database

The data model in database has been developed from complicated hierarchical model and network model to simple relational model. Knowledge base in expert system is now rule model and frame model. However, a simple approach in expert system could be made efficient use of databases and softwares accumulated until the present.

	DocN	Yr		DocN	Yr
	1122	83		1050	83
	959	81	--Co-Citation--	959	81
				930	80 Cit.
	781	78	--Co-Citation--	781	78
	674	77	--Co-Citation--	674	77
	104	74		400	75
	434	73	==Direct Link==	434	73 Doc.
				505	73
	32	72	--Cit-Ref.Link-	32	72
Cit.	405	72			
	31	71	--Cit-Ref.Link-	31	71 Ref.
	88	71		764	68
	38	70		428	68
Doc.	227	68	==Direct Link==	227	68
	406	68			
	37	67	--Bib.Coupling-	37	67
				96	67
Ref.	228	67	--Bib.Coupling-	228	67
	440	65		122	66
	454	65		52	64

Referation database described here has very simple data structure that is merely table type of file accompanied with referation file. What differs from Boolean search is that referation search makes out the weighted relationships between documents. Next, I describe about the measures.

Figure 3 Association Count Between Document Number 227 and 434

### ASSOCIATION MEASURES BETWEEN DOCUMENTS

It is possible to define an indicator of relationships between two documents using referentions. As a basic measure, the matching count between items is always used. For example, reference analysis measures bibliographic coupling between two documents based on references, and citation analysis counts co-citation based on citations. In the case of referentions, the same matching method is used.

and 227 do not belong to any categories. Direct link is found to document 227 and 434. This is indeed the reason that document itself was contained within referentions.

Bibliographic coupling is useful to latest documents, and co-citation count is adaptive to older documents. However, association count based on referentions is free to the publication year of document.

Figure 3 shows the association count between document number 227 and 434. The left side is the referentions of document 227 in 1968. And the right side is the referentions of document 434 in 1973. In these linkages, document 959, 781, and 674 is co-citation. And document 37 and 228 is bibliographic coupling. But document 434, 32, 31,

Three types of association measure between document A and B is defined as follows:

- 1) association count  
 $x = c$
- 2) association coefficient  
 $y = c / (a + b - c)$
- 3) association ratio  
 $z = 50 * (\text{SQR}(x/x_{\text{max}}) + \text{SQR}(y/y_{\text{max}}))$

in which a and b is each number of referentions in document A and B, and c

is the association count between A and B. The first type is association count itself. The second type is association coefficient normarized by a total of two referations. And the third type is association ratio adjusted by the association count and coefficient.

The merits of referation analysis is two points:

- 1) independence of publication year,
- 2) increase of the oppotunity of linkage.

It is desirable that computer readable referation database in various fields of science and technology realize.

### NEW INFORMATION RETRIEVAL SYSTEM

It became clearly that the representation of reference data by code was possible to list the referations of document and to measure the relationships between documents. Using referation database, a new information retrieval system will be developed.

In traditional information retrieval system, set of documents is serched by some query key such as descriptors, author names, journal name, and publication year. And an alterative set of documents through Boolean operation and/or sorting is displayed or printed out according to user's request. It is noticed that each document searched by Boolean techniques is no weight.

Referation search is the following three approachs:

- 1) search of referations in a given document,
- 2) search of documents associated with a given document,
- 3) search of documents, keywords, authors, and sources associated with set of documents obtained.

Table 1 shows an example of the third type of referation search. The set of 47 documents with keyword

"BIBLIOMETRIC" in title is retrieved from document file. Then the set of keywords associated with the 47 documents is searched using association measure based on referations. Table 1 is a part of screen in which three types of association measure are displayed in relevant order. Furthermore, the sets of documents, authors, and sources related to "BIBLIOMETRIC" can be searched by simple key operation.

### SOME FUNCTIONS OF REFERATION DATABASE

This gives a brief outline of some functions for referation database from the point of view of software design.

- (A) Input Procedure
  1. Definition of fields of document file  
The name, the number of bytes, and the attribute of each field are defined.
  2. Document file  
By assisting of indexes such as document, author, title, source, and publication year, document data is entried.
  3. Reference file  
By assisting of author index, reference data is entried using code.
- (B) Preparation of Auxilliary Files
  1. Field file
    - a. Merge of fields : author
    - b. one data a field : source, publication year
    - c. some data a field : title keyword.
  2. Index and name file  
Author, source, form, publication year, title, and keyword is prepared.
  3. Citation and referation file  
Citation : inverted of reference  
Referation : merge of citation, document, and reference.
  4. Rank file  
Referation, citation, reference, author, source, and keyword are

Table 1 Three Types of Keyword List Searched by 47 Documents with "BIBLIOMETRIC" in Title

I	No	I	Keyword	Count	I	Keyword	Coefft	I	Keyword	Ratio	I
I	1	I	CITATION	4866	I	BIBLIOME	0.0977	I	BIBLIOME	90.52	I
I	2	I	SCIENCE	4393	I	EMPIRICA	0.0879	I	BRADFORD	82.32	I
I	3	I	SCIENTIF	4161	I	BRADFORD	0.0864	I	LAW	79.18	I
I	4	I	BIBLIOME	3196	I	GENERAL	0.0852	I	DISTRIBU	75.07	I
I	5	I	BRADFORD	2826	I	LAW	0.0846	I	SCIENCE	66.48	I
I	6	I	LAW	2510	I	ZIPF	0.0841	I	SCIENTIF	66.17	I
I	7	I	LITERATU	2397	I	LOTKA	0.0833	I	CITATION	65.58	I
I	8	I	ANALYSIS	2354	I	DISTRIBU	0.0829	I	EMPIRICA	63.76	I
I	9	I	DISTRIBU	2073	I	APPLICAT	0.0808	I	ZIPF	63.63	I
I	10	I	INFORMAT	1947	I	STATIONA	0.0787	I	GENERAL	59.51	I

- ranked.
5. Bibliometric distribution file  
Bradford, Zipf, Lotka, growth, and  
life distributions are prepared.
- (C) Output Procedure
1. Referation list of a given  
document
  2. Index list  
KWIC, KWOC, author, source, title  
index
  3. Rank, name, and frequency list,  
and graph
  4. List of document searched.
- (D) Search and Association
1. Traditional search  
Keyword, author, source,  
publication year
  2. Referation search on a given  
document  
Referations and associated  
documents.
  3. Referation search on set of  
documents  
Associated documents, keywords,  
authors, and sources.
- (E) Auxillary Procedure
1. Sorting routine  
The usage to index and rank file.
  2. Boolean Logic  
AND, OR, NOT operation in search.
  3. Association measure  
Association count, coefficient,  
and ratio are obtained.

#### CONCLUSION

When the document related to a specific topic in references is represented by code, a new attribute of document, called for referations, can be easily made of citations, references, and document itself. And a new information retrieval system using referations can be realized on microcomputers.

Boolean search retrieves no weight and only documents. On the other hand, referation search produces the set of weighted documents, keywords, authors, and sources. These are some kinds of professional information in a specific topic.

This micro-based pilot system is very simple and feasible. The concept of referations and referation database could be applied to not only document information but also the patent and precedent information. It is desirable to develop a general software for referation database.

#### References

- (1) Eugene Garfield, "Citation Indexes for Science," Science, 122 (1955) 108-111.
- (2) Henry Small, "Co-citation in the Scientific Literature; A New Measure of the Relationship Between Two Documents," Journal of the American Society for Information Science, 27 (1973) 265-269.
- (3) Eugene Garfield, Citation Indexing: Its Theory and Application in Science, Technology and Humanities. (John Wiley & Sons, 1979).
- (4) G.Vladutz and J.Cook, "Bibliographic Coupling and Subject Relatedness," Proceedings of the 47th ASIS Annual Meeting, 21 (1984) 204-207.
- (5) M.M.Kessler, "Bibliographic Coupling Between Scientific Papers," American Documentation, 14-1 (1963) 10-25.
- (6) Isao Asai, "Development of a "Referation" Database Using Micro-computers," Proceedings of the 21th Annual Meeting on Information Science and Technology, (1984) 21-31 (in Japanese).
- (7) Isao Asai, "Analysis of Bibliographic Information on a Specific Topic Based on "Referation" Relationship," Proceedings of the 22th Annual Meeting on Information Science and Technology, (1985) 135-142 (in Japanese).
- (8) Isao Asai, "The Structure and Usage of A "Referation" Database on Micro-computers For Personal Document Management," (Lehmann, K.-D. ed. The Application of Micro-Computers in Information, Documentation, and Libraries, North Holland, Elsevier, 1986) 510-517.

# INFORMATION:

*The Transformation of Society*

ASIS '87

Proceedings of the  
50th Annual Meeting  
of the American Society  
for Information Science

Boston, Massachusetts  
October 4-8, 1987

Volume 24

*asis*