

### III.5. CITATION ANALYSIS OF SCIENTIFIC JOURNALS

The study of the use and relative impact of scientific journals is one of the important applications of citation analysis. Investigations on citations received a considerable impetus by the annual publication (since 1976) of the Journal Citation Reports (JCR) by the Institute of Scientific Information (ISI) (Garfield (1976a)) as part of the SCI and the SSCI. The first section of this chapter is devoted to the description of the JCR. We further study the statistical reliability of citation measures and some alternatives for the citation measures published in the JCR.

#### III.5.1. The Journal Citation Reports (JCR)

Generally speaking, the JCR is an annually published, statistical data set providing information on how often journals are cited, how many items were published, and how often, on the average, each item is cited. It also reports those source journals responsible for the references of each journal, the number of references each journal has published and the distribution of those references in time. We will now describe the different sections of the JCR (SCI edition) in more detail.

In addition to an introductory and explanatory part, some reprints on journal evaluation by citation analysis, a bibliography and a journal list (explaining the abbreviations used in the JCR), this volume consists of five parts entitled : Journal Rankings, Source Data Listing, Journal Half-Life Listing, Citing Journal Listing and Cited Journal Listing.

##### III.5.1.1. Journal Rankings

Section 1 of the *Journal Rankings* is an alphabetical listing of source journals according to abbreviated title. This listing consists of 13 columns. The first column contains sequential numbers and the second the journal-title abbreviation. The next four columns give the total number of times the journal was cited by SCI, SSCI and A&HCI source items during the year Y (year covered by this issue of the JCR), the portion of those total citations accounted for by articles published in the year Y-1, the portion of the total citations accounted for by articles the journal published in the year Y-2, and the sum of the two previous columns. For an illustration we refer the reader to Table III.5.1, taken from the 1987 edition of the JCR (Garfield (1988)).

Table III.5.1. First 6 columns of the Journal Rankings Section 1 of the 1987 JCR (Garfield (1988))

RANK	JOURNAL TITLE	CITATIONS IN 1987 TO			
		ALL YEARS	1986	1985	86+85
1	A VAN LEEUW J MICROB	809	47	34	81
2	AAPG BULL	3233	128	235	363

The next three columns give the number of source items published by the journal in the year Y-1, in the year Y-2 and in both years. The 10<sup>th</sup> column, headed 'impact factor', gives a figure for the relative frequency with which the journal's 'average paper' has been cited. This figure has been obtained as the ratio of all citations in the year Y to papers published in the years Y-1 and Y-2 (column 6) to the total number of source items published in the years Y-1 and Y-2 (column 9), see Table III.5.2.

Table III.5.2. Columns 7 to 10 from the Journal Rankings Section 1 of the 1987 JCR (Garfield (1988)), same lines as Table III.5.1

SOURCE ITEMS IN			IMPACT FACTOR	
1986	1985	86+85		
75	40	115	0.704	(= $\frac{81}{115}$ )
100	131	231	1.571	(= $\frac{363}{231}$ )

The impact factor (IPF) is a better measure for the scientific importance of a journal than the total number of citations because it takes the total number of publications into account.

The next two columns (columns 11 and 12) show, respectively, the number of times articles in the journal's year Y issues were cited in the references of SCI, SSCI and A&HCI source items in the same year Y and the number of source items the journal itself published during that year. The last column, headed 'immediacy index', is the result of dividing column 11 by column 12, see Table III.5.3. The immediacy index (IMI) is supposed to show how fast new ideas published in the journal are taken up by the scientific community. In practice, however, it is highly dependent on backlogs in the journals of the field. It is often simply impossible for papers to be cited in the same year

as they were published, except in the case of preprint distributions, informal contacts and other activities of invisible colleges.

Table III.5.3. Columns 11 to 13 from the Journal Rankings Section 1 of the 1987 JCR (Garfield (1988)), same lines as Tables III.5.1 and III.5.2

CITATIONS IN 1987 TO 1987 ITEMS	SOURCE ITEMS IN 1987	IMMEDIACY INDEX
2	40	0.050
19	112	0.170

Denoting the number of citations in the year  $Y$  to papers published in the year  $Z$  by  $CIT_Z(Y)$ , the number of publications in the year  $Y$  by  $PUB(Y)$ , the impact factor in the year  $Y$  by  $IPF(Y)$  and the immediacy index in the year  $Y$  by  $IMI(Y)$ , we have the following equations :

$$IPF(Y) = \frac{CIT_{Y-1}(Y) + CIT_{Y-2}(Y)}{PUB(Y-1) + PUB(Y-2)} \quad [III.5.1]$$

$$IMI(Y) = \frac{CIT_Y(Y)}{PUB(Y)} \quad [III.5.2]$$

Sections 2 to 6 contain the same information as Section 1 but listed differently. Section 2 gives a ranking according to the number of citations over all years, Section 3 according to the impact factor and Section 4 according to the immediacy index. Section 5 ranks journals by source items published in the year  $Y$  and Section 6 by the number of citations in the year  $Y$  to articles published in the years  $Y-1$  and  $Y-2$ . Section 7 is a listing of social sciences journals, arranged alphabetically according to abbreviated title and contains the same information as Section 1 does for science journals. Section 8 is a breakdown of SCI source journals by subject category, ranked by impact factor per category; cited half-life is also shown. The 'cited half-life' of a journal refers to the number of journal publication years, going back from the current year, which account for half of the total citations received by the cited journal during the current year. Finally, Section 9 is an alphabetical listing of all source journals and their Section 8 category listing. This section allows the user to quickly determine which category (or categories) a journal is listed under.

III.5.1.2. Source Data Listing

The second part of the JCR is the *Source Data Listing* which lists alphabetically the journals covered by the SCI with corresponding details on the number of articles published, the total number of references contained in these articles and the average number of references per article. Articles are subdivided into review and non-review articles and data are given for each category separately, see Table III.5.4.

Table III.5.4. Source Data Listing : an example taken from the 1987 JCR (Garfield (1988))

JOURNAL NAME			NON-REVIEW ARTICLES		
			SOURCE ITEMS (S)	REFERENCE ITEMS (R)	RATIO (R/S)
J ANTIMICROB CHEMOTH			245	4192	17.1
REVIEW ARTICLES			COMBINED TOTAL NON-REVIEW AND REVIEW		
SOURCE ITEMS (S)	REFERENCE ITEMS (R)	RATIO (R/S)	SOURCE ITEMS (S)	REFERENCE ITEMS (R)	RATIO (R/S)
5	259	51.8	250	4451	17.8

III.5.1.3. Journal Half-Life Listing

The third part is the *Journal Half-Life Listing*, containing three sections. Section 1 lists source journals alphabetically according to abbreviated title and shows the cumulative percentage of citations given by the citing journal in the year Y. The column on the upper left shows the citing half-life, which is defined as (see also Section III.6.2 and III.6.3) the number of journal publication years from the current year going back, which accounts for 50 % of the total number of references given by the citing journal, see Table III.5.5.

Table III.5.5. Journal Half-Life Listing Section 1 : an example taken from the 1987 JCR (Garfield (1988))

CITING HALF-LIFE	CITING JOURNAL					1987	1986	1985	1984	1983
4.9	LASER SURG MED					0.84	8.09	24.13	38.16	50.11
1982	1981	1980	1979	1978						
57.82	64.22	68.46	72.08	75.94						

Section 2 of the Journal Half-Life Listing shows the cumulative chronological distribution of cited references given by journals in the year Y to articles published in cited journals during the last 10 years. The column on the left shows the cited half-line. A small cited half-line may indicate that the journal is a rather recent one, or that it mainly publishes papers of immediate interest. In the latter case it might be a good idea for a librarian to relegate older issues to subordinate shelf-space. See Table III.5.6 for an example.

Table III.5.6. Journal Half-Life Listing Section 2 : an example taken from the 1987 JCR (Garfield (1988))

CITED HALF-LIFE 2.9	CITED JOURNAL LASER SURG MED	1987 4.72	1986 24.99	1985 50.44	1984 73.64	1983 89.18
1982 94.81	1981 96.61	1980 100.0	1979 100.0	1978 100.0		

Finally, Section 3 lists the journals in descending order of cited half-life.

#### III.5.1.4. Citing Journal Listing

The *Citing Journal Listing* lists all citing journals alphabetically according to their abbreviated titles. The first line of each entry shows the journal's impact factor, abbreviated title and total number of references. Succeeding columns of this row distribute the total number of references by year in which the articles cited in the references were published. Under this first line are listed the journals cited in the references of the citing journal named in the main entry line. These cited journals are listed in descending numerical order according to the frequency of their citation in references of the citing journal, see Table III.5.7.

Table III.5.7. Citing Journal Listing : an example taken from the 1987 JCR (Garfield (1988))

	TOTAL	1987	1986	1985	1984	1983	1982
1.15 J AM STAT ASSOC	2416	49	145	203	200	161	165
1.15 J AM STAT ASSOC	391	12	29	42	37	26	22
1.19 ANN STAT	170	5	9	18	22	12	17
1.00 BIOMETRIKA	153	1	9	8	8	13	7
...							
1981	1980	1979	1978	REST			
157	118	96	97	1025			
34	17	17	15	140			
18	8	10	8	43			
4	6	5	8	84			
...							

## III.5.1.5. Cited Journal Listing

The *Cited Journal Listing* lists cited journals in alphabetical order according to their abbreviated titles. Here the first line gives the journal's impact factor, abbreviated title and total number of citations received in the year Y. Succeeding columns distribute the citation total according to the year in which the cited articles were published. Under this main entry line are listed the journals in whose references citations to the cited journal appeared. These citing journals are listed in descending order according to the number of citations each contributed to the citation total, see Table III.5.8.

Table III.5.8. Cited Journal Listing : an example taken from the 1987 JCR (Garfield (1988))

	TOTAL	1987	1986	1985	1984	1983
0.53 J CHEM ENG DATA	1691	27	76	95	103	76
0.53 J CHEM ENG DATA	247	12	26	28	23	15
0.86 FLUID PHASE EQUILIBR	172	0	5	16	7	11
0.82 J CHEM THERMODYN	95	0	1	13	4	4
...						
1982	1981	1980	1979	1978	REST	
97	84	65	60	77	931	
18	16	5	6	5	93	
16	12	7	7	11	80	
9	4	2	1	3	54	
...						

From the Citing Journal Listing and the Cited Journal Listing we can derive the self-citing and the self-cited rates. The self-citing rate relates a journal's self-citation to the total number of references it gives. From Table III.5.7 we see that the 1987 self-citing rate of the J AM STAT ASSOC is  $391/2416 = 0.162$  or 16.2 %. The self-cited rate relates a journal's self-citations to the number of times it is cited by all journals, including itself. From Table III.5.8 we see that the 1987 self-cited rate of the J CHEM ENG DATA is  $247/1691 = 0.146$  or 14.6 %. A high self-cited rate is an indication of a journal's low visibility. A high self-citing rate is rather an indicator of the isolation of the field covered by the journal.

### III.5.2. Reliability of comparisons based on citation measures

Rankings of journals according to the number of citations received, the impact factor or the immediacy index are only meaningful as long as fluctuations reflect a real rise or drop in the importance or influence of the journal, and is not only the result of a purely random process. To account for the random effect on citation measures, Schubert and Glänzel (1983) devised a method for estimating the standard error of mean citation rates per publication and applied this method to find confidence intervals for the impact factor.

Schubert and Glänzel (1983) regard the publication of papers within a time period from  $s_1$  to  $s_2$ ,  $s_1 \leq s_2$ , as an action and a citation to these papers in a year  $s_2 + T$ ,  $T \geq 0$ , as a reaction. This action-reaction process is then modelled as a stochastic process. The function  $X_T$ ,  $T \geq 0$ , then denotes the stochastic variable which maps a paper published in a specified journal during the period from  $s_1$  to  $s_2$  to the number of citations it receives in the year  $s_2 + T$ .  $P(X_T = k)$  denotes the probability that a paper published in the period from  $s_1$  to  $s_2$  will receive exactly  $k$  citations in the year  $s_2 + T$ . For any  $T$ ,  $X_T$  is assumed to have a negative binomial distribution. So we have :

$$P(X_T = k) = \binom{n+k-1}{k} p_T^n q_T^k, \quad k = 0, 1, 2, \dots$$

(cf. Subsection I.2.4.4), where  $n > 0$  is a fixed parameter and  $p_T \in [0, 1]$  is a parameter depending on  $T$ ,  $q_T = 1 - p_T$ . From Subsection I.2.4.4 we know that  $E(X_T) = nq_T/p_T$  and  $\text{Var}(X_T) = nq_T/p_T^2$ .

To study the impact factor, we choose  $s_2 = s_1 + 1$  (so that we cover a two-year period) and  $T = 1$ ; the associated stochastic variable is then denoted as  $X_1$ . Let  $J$  be the stochastic variable that a paper will be published in journal  $j$ . Then the impact factor (IPF) can be defined as the conditional

expectation

$$Z = E(X_1 | J) .$$

Using the assumption that  $P(X_1 = k | J)$  follows a negative binomial distribution, we find :

$$Z = E(X_1 | J) = \sum_k k P(X_1 = k | J) = n_j q_{1,j} / p_{1,j} , \quad [III.5.3]$$

where  $n_j$  and  $p_{1,j}$  are parameters characteristic of journal  $j$ . The impact factor as published in the JCR (cf. Section III.5.1), being the empirical mean of a finite sample (the number of papers published during two years in journal  $j$ ), is an estimate of this conditional expected value. This estimate is affected by a certain error, which is characterised by the variance, or its square root, the standard error. From Section I.3.1 we know that  $\text{Var}(Z) = \text{Var}(X_1)/N = E(X_1)/(p_{1,j}N)$  (where  $N$  denotes the number of publications in journal  $j$  during the years  $s_1$  and  $s_1 + 1$ ).

In order to estimate the standard error of the impact factor, we need estimations for  $E(X_1)$  and  $p_{1,j}$ . The expectation  $E(X_1)$  is estimated by the sample mean, i.e. the impact factor (IPF) as found in the JCR. The parameter  $p_{1,j}$  is best found as the solution of the equation

$$\frac{\log(f_0)}{\text{IPF}} = \frac{p_{1,j}}{1 - p_{1,j}} \log(p_{1,j}) , \quad [III.5.4]$$

where  $f_0$  is the fraction of uncited papers. (This method is explained by Johnson and Kotz (1969)). If  $p_0$  is the solution of [III.5.4], we obtain as an estimation for the standard error of the impact factor the value

$$\left(\frac{\text{IPF}}{N p_0}\right)^{1/2} . \quad [III.5.5]$$

Schubert and Glänzel use this estimate to compare two impact factors, i.e. to see whether their difference is statistically significant. They also verify, successfully, the hypothesis that citations to journals are negative binomially distributed.

Their work, however, has one serious practical drawback : they need to know  $f_0$ , the fraction of uncited papers, but this fraction is not given in the JCR. (Schubert and Glänzel had access to the original tapes of the ISI database and could calculate this fraction.) Therefore, Nieuwenhuysen and Rousseau (1988) devised a quick and easy way, based on the Poisson distribution, to find a



lower bound on the size of fluctuations of the impact factor and the immediacy index. Under this assumption (a lower bound for) the length of a 95 % confidence interval for the impact factor of a journal is found, leading to the following minimal confidence interval :

$$\left[ \text{IPF} - 1.96 \frac{(\text{CIT}_1)^{1/2}}{N}, \text{IPF} + 1.96 \frac{(\text{CIT}_1)^{1/2}}{N} \right], \quad [\text{III.5.6}]$$

where  $\text{CIT}_1$  denotes the observed number of citations to journal  $j$  in the year  $s_2 + 1$ . Similarly, a minimal 95 % confidence interval for the immediacy index (IMI) is given by :

$$\left[ \text{IMI} - 1.96 \frac{(\text{CIT}_0)^{1/2}}{N_0}, \text{IMI} + 1.96 \frac{(\text{CIT}_0)^{1/2}}{N_0} \right]. \quad [\text{III.5.7}]$$

Here  $s_1 = s_2$  and  $T = 0$ ;  $N_0$  denotes the number of publications in the year  $s_2 + 0$  and  $\text{CIT}_0$  the number of citations to these publications in the same year.

### III.5.3. Proposals for citation measures other than those published in the JCR

#### III.5.3.1. The Pinski-Narin influence measure

The impact factor used in the JCR was defined by Garfield (1972). We note, however, that as early as 1960 this quantity was suggested (Raisig (1960)) as a measure for the impact of serials. Raisig called it the 'index of realised research potential'. Although the impact factor is a size-independent measure, since it is defined as a ratio, it suffers from other limitations (Pinski and Narin (1976)). According to these authors, the definition and calculation of the IPF does not contain any correction for the average length of individual papers. As a result of Garfield's method, journals which publish longer papers, especially review journals, tend to have higher impact factors.

A second limitation is that citations are not weighted. All citations are counted as equally important, regardless of the citing journal. A third limitation is that there is no normalisation for the different referencing characteristics of different scientific fields. To remedy these limitations, Pinski and Narin (1976) propose a new weighted measure for journals.

They start with an  $(n,n)$  citation matrix  $C = (C_{ij})$ , where  $C_{ij}$  indicates the number of references journal  $i$  gives to journal  $j$ . Then, they wish to extract from the citation matrix a measure of influence for each journal in the set.

$W_i$ , the influence weight of the  $i^{\text{th}}$  journal, is defined as :

$$W_i = \sum_{k=1}^n \frac{W_k C_{ki}}{S_i} \quad i = 1, \dots, n, \quad [\text{III.5.8}]$$

where  $S_i = \sum_{j=1}^n C_{ij}$  = total number of references contained in the  $i^{\text{th}}$  journal.

This yields a system of  $n$  equations, one for each  $i$ , that is solved iteratively by

$$W_i^{(m)} = \sum_{k=1}^n \frac{W_k^{(m-1)} C_{ki}}{S_i}, \quad [\text{III.5.9}]$$

where, as a first approximation to the weight of journal  $i$ , they use

$$W_i^{(1)} = \frac{\text{total number of citations to journal } i}{\text{total number of references to journal } i} \quad [\text{III.5.10}]$$

(during a fixed period).

The influence weights obtained in this way are a measure of influence per reference. The influence per publication is then defined as the weighted number of citations (each citation weighted by the weight of the journal it appears in) a publication receives. The total number of weighted citations for the  $i^{\text{th}}$  journal is :

$$\sum_{k=1}^n W_k C_{ki} = W_i S_i$$

(by means of [III.5.8]).

To get the influence per publication, one divides by the annual number of publications,  $\text{PUB}(i)$ . Multiplying  $W_i$  by  $(S_i/\text{PUB}(i))$  therefore yields the desired measure.

Thus, the Pinski-Narin influence measure for journal  $i$  in a network is defined as :

$$W_i \frac{S_i}{\text{PUB}(i)}, \quad [\text{III.5.11}]$$

where  $W_i$  and  $S_i$  only make sense with regard to the network under consideration. This measure has been refined by Geller (1978) (who used a Markov chain approach), Todorov (1984) and Noma (1988).

### III.5.3.2. Impact factors calculated over different periods

The Garfield impact factor in the year  $Y$  ( $\text{IPF}(Y)$ ) was defined (Garfield (1972)) as :

$$IPF(Y) = \frac{\sum_{i=1}^2 CIT_{Y-i}(Y)}{\sum_{i=1}^2 PUB(Y-i)} \quad [III.5.12]$$

It is now natural to define impact factors over different periods (Rousseau (1988a)) :

$$IPF_n(Y) = \frac{\sum_{i=1}^n CIT_{Y-i}(Y)}{\sum_{i=1}^n PUB(Y-i)} \quad [III.5.13]$$

This generalised impact factor satisfies the following difference equation :

$$IPF_n(Y) - \frac{SPUB(Y-n+1)}{SPUB(Y-n)} IPF_{n-1}(Y) = \frac{CIT_{Y-n}(Y)}{SPUB(Y-n)} \quad [III.5.14]$$

where  $SPUB(Y-k) = \sum_{i=1}^k PUB(Y-i)$ ,  $IPF_0(Y) = 0$ . Note that  $IPF_2(Y)$  is the Garfield impact factor (= IPF). We emphasise the fact that in our opinion

$$MAX(Y) = \max_{n=1,2,\dots} (IPF_n(Y)) \quad [III.5.15]$$

is a better measure of impact than the Garfield impact factor, for it is less field-dependent. Moreover, it has been shown (Dierick and Rousseau (1988)), using a random sample of 107 science journals, that the  $IPF_4$  and the  $IPF_3$  are generally larger than the Garfield impact factor ( $IPF_2$ ), indicating that the Garfield impact factor is also influenced by immediacy effects.

The relation between  $IPF_n$  and  $A_n(Y) = \frac{CIT_{Y-n}(Y)}{PUB(Y-n)}$  was studied in Rousseau (1988a) and further clarified, using a continuous approach, by Egghe (1988b).

### III.5.3.3. Other proposals

Schubert and Glänzel (1983) restrict the calculation of impact factors to papers classified as 'articles', 'reviews', 'notes' and 'letters to the editor'. Impact factors restricted to these types of publication are called 'corrected impact factors'.

Price (1981b) and Noma (1982) proposed a method to reduce the quantitative excess of journal self-citations in studying networks of journals.

Schubert and Glänzel (1986) proposed a measure of the citation speed of journals : the mean response rate (MRT) defined as :

$$\text{MRT} = -\log(f_0 + f_1 e^{-1} + f_2 e^{-2} + f_3 e^{-3} + f_4 e^{-4}) , \quad [\text{III.5.16}]$$

where  $f_i$  is the fraction of papers receiving their first citation in the  $i^{\text{th}}$  year after publication.

#### III.5.4. Notes and comments

Numerous papers have been published concerning citation measures and journal evaluation. We restrict ourselves to some examples. A review on journal citation measures has been published by Todorov and Glänzel (1988), to which the reader is referred for further information.

According to Line (1977) and Scarlan (1988), impact factors are of little value to special librarians, because, while users of journals read, many actually publish little or not at all. On the other hand, Anderson and Goldstein (1981) include the impact factor among their criteria of journal quality.

New measures of the relative standing of journals with respect to their subfields have been studied by Doreian (1987, 1988). The statistical validity of citation measures was studied by Schubert and Glänzel (1983), but also by Tomer (1986). The relationship between local use at the Antwerp State University Centre and citation measures was investigated as early as 1974 by Van Styvendaele (1974). Vervliet (1987) reports on the use of JCR data for a defensive collection policy for Belgian university libraries. In this context a 'defensive collection policy' means making a coordinated effort to reduce journal cancellations. From JCR data a new ranking was drawn up, which included the price per citation.

The term 'half-life' as a journal citation measure has been borrowed from the terminology of nuclear physics by Burton and Kebler (1960). The study of half-lives is an aspect of the concept of obsolescence, to be studied in the next chapter. We also mention an important earlier network study of psychological journals by Xhignesse and Osgood (1967).

There are significant differences in the citation potentials of different scientific fields, i.e. in the maximum number of times any given article - and hence also any journal - will be cited in its lifetime. The most widespread assumption is that the citation potential is a function of the number of research workers active in the discipline. Garfield (1979c) rejects this and claims that different research potentials are connected with different R/A (references per

article) numbers of the journal publications in the fields. From an investigation of biochemistry (with higher R/A's) and plant physiology journals (with lower R/A's), Marton (1983) found that the reason for the lower citation potentials in plant physiology are :

- 1 the readership other than specialists in the subject is narrower for plant physiology than for biochemistry journals;
- 2 plant physiologists have fewer thematically relevant new articles to cite than biochemists;
- 3 plant physiology research fields are relatively isolated, whereas biochemistry research fields are relatively integrated.

These factors explain the higher R/A numbers of biochemistry.

Studies of journals abound in the literature : among other studies we already mentioned Earle and Vickery (1969) for the social sciences, Pinski and Narin (1976) for physics, Keteleer (1986) for botany, Peritz (1986) for demography, Rousseau (1988a) for mathematics and Rousseau (1989b) for pharmacology. In his book on citation analysis Garfield (1979a) gives examples of a study on the relationship between two journals (Phytopathology and Virology), a study on one journal in relation to others (Journal of Clinical Investigation), studies of scientific fields (pediatrics, physics, botany-agriculture and engineering), and a study on the relation between fields (geology and geophysics) and the literature of one country (Soviet Union, France, Japan and Germany). Livestock periodicals with respect to a library selection policy were investigated by Adewole (1987); the citation behaviour of Indian phytopathologists was studied by Nagappa and Maheswarappa (1981). Jan Vlachý has published numerous papers on physics journals, books and individual papers (see, for example, Vlachý (1981,1983,1984)). An early study of the impact factor and the immediacy index applied to physics journals can be found in Inhaber (1974).

### III.6. OBSOLESCENCE

#### III.6.1. Generalities

A distinction has to be made between the 'general' or worldwide *obsolescence* of the literature on a given subject and the decline to the local use of documents in a particular library. The latter aspects of obsolescence was already treated in Part II. Here we will study the general obsolescence as measured by citation rates.

We can also take either a diachronous or a synchronous view of the obsolescence of scientific literature. In the diachronous view one considers a fixed group of documents, e.g. one year's issue of periodicals and studies the evolution in citations as time,  $t$ , increases. In the synchronous view one is concerned with the distribution of citations to documents of different ages during a given span of time. In this sense the JCR offers every year a synchronous view of the scientific literature. Nakamoto (1988) has observed a remarkable symmetry in diachronous and synchronous citation rates.

A review on obsolescence, although mainly focusing on local obsolescence, can be found in Gapen and Milner (1981).

#### III.6.2. The half-life analogy as applied to scientific literature

In nuclear physics the concept of half-life is used to describe the decay of radioactive substances. For physicists it means the time required for 50 % of the atoms in a sample of a radioactive source to disintegrate. These physical half-lives are of equal duration, that is, at any given time, the half-life of the remaining material is the same as the half-life of the original source.

Analogous to this physical concept, Burton and Kebler (1960) define the *half-life* of scientific literature as the time during which one-half of all the currently active (= cited) literature is published. When applied to journals as a source, this definition coincides with that of the citing half-life as used in the JCR (cf. Subsection III.5.1.3). It can be considered as a crude one-dimensional measure of obsolescence.

In a test Burton and Kebler (1960) found the following half-lives (Table III.6.1) :

Table III.6.1. Literature half-lives as calculated by Burton and Kebler (1960)

Chemical Engineering	4.8 years
Mechanical Engineering	5.2
Metallurgical Engineering	3.9
Mathematics	10.5
Physics	4.6
Chemistry	8.1
Geology	11.8
Physiology	7.2
Botany	10.0

This table shows that, generally speaking, the literature of more theoretical sciences show longer half-lives (e.g. mathematics) than the literature of those fields which are dependent on fresh data or new technological innovations (e.g. metallurgical engineering).

In recent studies the term 'half-life' has been used exclusively for diachronous studies (contrary to the original Burton-Kebler paper) and the term 'median citation age' has been used for synchronous studies (Wallace (1986), Stinson and Lancaster (1987)). In the next section, a relation between the half-life and the ageing rate will be derived.

### III.6.3. Determination of the ageing rate and the half-life

Avramescu (1979) proposes the following equation [III.6.1] for the diachronous citation distribution of individual papers or journals :

$$y(t) = C_0(e^{-\alpha t} - e^{-mt}) , \quad m > \alpha \quad . \quad [III.6.1]$$

Fig.III.6.1a illustrates this function for the case in which  $C_0 = 100$ ,  $\alpha = 0.4$  and  $m = 1$ . The citation distribution of brilliant papers suffering from delayed recognition (see also Garfield (1989b)) can be obtained from [III.6.1] by taking  $m$  almost zero and  $\alpha < 0$ . This is illustrated in Fig.III.6.1b for  $C_0 = 20$ ,  $\alpha = -0.1$  and  $m = 0.01$ .

When we use [III.6.1] for the synchronous citation distribution and assume  $m \gg \alpha$ , then  $y(t)$  is approximated by  $C_0(e^{-\alpha t}) = C_0 a^t$ , where we write  $e^{-\alpha} = a$  ( $a$  is then called the 'ageing factor'). Brookes (1970,1971) proposed the following elegant method to estimate the ageing factor from practical data. Determine the number of years  $i$  (usually  $i = 6, 7$  or  $8$ ) and set  $k$  equal to the number of citations to papers published at least  $i$  years ago. Let  $l$  be the

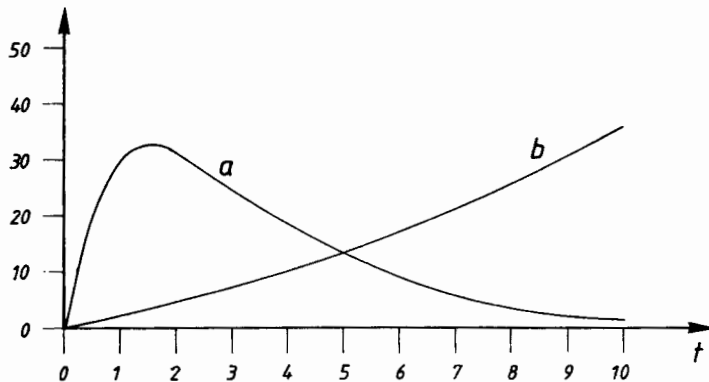


Fig.III.6.1 Citation distributions

number of citations to papers published less than  $i$  years ago. Then

$$\begin{aligned} k &= C_0(a^i + a^{i+1} + \dots) \\ &= a^i C_0(1 + a + a^2 + \dots) \\ &= a^i(k + \ell). \end{aligned}$$

Hence :

$$a^i = \frac{k}{k + \ell}$$

or

$$a = \sqrt[i]{\frac{k}{k + \ell}} \quad \text{[III.6.2]}$$

The same equation is found when time is considered as a continuous variable :

$$\begin{aligned} k &= \int_i^{\infty} C_0 a^{t'} dt' \\ &= a^i \int_0^{\infty} C_0 a^{t'} dt' \\ &= a^i(k + \ell). \end{aligned}$$

Hence, again :



$$a = \sqrt{\frac{k}{k+l}} \quad [III.6.3]$$

The continuous model is needed to determine from this the half-life  $t_0$  (since  $t_0$  is not necessarily an integer), defined mathematically as follows (cf. Section III.6.2) :  $t_0 \in \mathbb{R}^+$  is the unique number such that :

$$\int_{t_0}^{\infty} C_0 a^{t'} dt' = \int_0^{t_0} C_0 a^{t'} dt' = \frac{1}{2} \int_0^{\infty} C_0 a^{t'} dt' \quad .$$

This gives (since  $a < 1$ ) :

$$\frac{-a^{t_0}}{\log a} = -\frac{1}{2} \frac{1}{\log a} \quad .$$

Hence

$$a^{t_0} = \frac{1}{2}$$

Consequently :

$$t_0 = \frac{\log 0.5}{\log a} = -\frac{\log 2}{\log a} \quad [III.6.4]$$

Equation [III.6.4] can be used for any half-life calculation; hence we can use [III.6.4] for both cited and citing half-lives.

As the cited and citing half-lives for journals are given in the JCR (at least if  $T < 10$ ), equation [III.6.4] can be used as a test for the accuracy of Brookes' method (leading to [III.6.2]). Indeed [III.6.4] yields :

$$a = 2^{-1/t_0} \quad [III.6.5]$$

We have checked this for the botany journals studied by Keteleer (1986). This results in Table III.6.2, showing excellent agreement between both methods.

Table III.6.2. The relation between the cited half-life and the ageing factor

A : journal

B : cited half-life according to the 1983 JCR

C : a : calculated according to equation [III.6.2] (Brookes' method)

D : a : calculated according to equation [III.6.4]

A	B	C	D
PLANT PHYSIOL	6.1	0.89	0.89
PHYTOCHEMISTRY	6.6	0.89	0.90
PHYTOPATHOLOGY	9.6	0.93	0.93
PLANTA	5.8	0.88	0.89
CAN J BOT	7.7	0.91	0.91
PHYSIOL PLANT	7.3	0.91	0.91
AM J BOT	(12.7)	0.95	- *
NEW PHYTOL	6.9	0.90	0.90
ANNU REV PLANT PHYS	6.3	0.89	0.90
J EXP BOT	6.7	0.90	0.90
ANN BOT-LONDON	8.6	0.92	0.92
PLANT CELL PHYSIOL	5.3	0.87	0.88
Z PFLANZENPHYSIOL	4.8	0.84	0.87
PLANT SOIL	8.7	0.92	0.92
PLANT SCI LETT	4.2	0.78	0.85
J PHYCOL	7.1	0.90	0.91
WEED SCI	6.9	0.90	0.90
BOT GAZ	(11.9)	0.94	- *
CAN J PLANT SCI	8.1	0.92	0.92
AUST J PLANT PHYSIOL	5.1	0.78	0.87
PHYSIOL PLANT PATHOL	5.4	0.85	0.88

\* The JCR does not give the exact cited half-life when it is greater than 10 years. In these cases we could only apply Brookes' method and calculate the cited half-life from the ageing factor.

III.6.4. 'Real' versus 'apparent' - 'synchronous' versus 'diachronous'

Line and Sandison (1974) have argued that there is no reason to suppose that obsolescence measured synchronously will be the same as obsolescence measured diachronously. When citations are used to measure obsolescence, the growth in the number of scientists and/or publications is a complicating factor. Suppose that the 1990 literature on some subject cites the 1985 literature twice as much as the 1975 literature (in absolute numbers). However, if between 1975 and 1985 the number of papers on this subject has doubled, then the 1990 citations merely reflect probability and show no evidence of decline in use with age. So it is not surprising that people such as Line (1970) and Sandison (1971) have stated that ageing curves based on actual data only show an apparent decline and that obsolescence should be corrected for the growth

of the literature.

Brookes (1970,1980b) finds the whole idea of a 'real' obsolescence in contrast to an 'apparent' obsolescence a rather 'mystical' notion. For him the uncorrected measure reflects the reality. Moreover, Brookes (1970) showed that the growth in the number of contributors to the literature on some subject will cancel the effect of the growth of the literature if both rates are equal. This can be understood by the following reasoning. If the literature increases, the chance that a particular paper will be cited decreases (because papers are competing with each other to be incorporated in the reference lists of new papers). On the other hand, the more authors that contribute to the subject, the greater the chance that a particular paper will be discovered and chosen for citation.

To test Brookes' hypothesis that both growth rates will tend to balance each other in obsolescence studies and to compare synchronous and diachronous measures of obsolescence, Stinson (1981) and Stinson and Lancaster (1987) conducted an extensive study in the area of human and medical genetics.

They compared 13734 citations in the diachronous study over a time span of 19 years, with 3669 references in the synchronous study. When the first two years of the synchronous study were excluded, it was found that the rate of obsolescence measured diachronously and corrected for the growth of the SCI (this is not the same as correcting for the growth of the literature!) was statistically equivalent to that found in the uncorrected synchronous study. Hence, this investigation suggests that Brookes is right and that the easier uncorrected synchronous study provides an accurate measure of the decline in use with age.

#### III.6.5. Notes and comments

Coughlin and Baran (1988) studied several stochastic models of information obsolescence. Burrell's model with ageing (1985b), discussed in Subsection II.6.2.2, although constructed for local use, receives strong empirical support from citation histories. Other obsolescence models were studied by MacRae (1969) and Krauze and Hillinger (1971). Griffith et al. (1979) have found that journal self-citations show a faster obsolescence than references given to other sources. Further, they distinguished archival journals (cf. Price (1970)), which age slowly - theoretically they should show no ageing - and research front journals which exhaust their utility within a few years. The idea of over-citation of recent material (and no real decline in the utility of older literature) was put forward by Price (1965) and corroborated by Marton (1985).

Motylev (1986,1989) disclaims the idea that scientific literature ages rapidly, that publications in rapidly developing fields age faster and that the maximum book use occurs only a few years after its publication. He claims that these ideas have cropped up because a wrong methodology was used in the treatment of empirical data and because the meaning of the term 'ageing of scientific and technical literature' has been incompletely (and wrongly) understood. According to Motylev (1989), the use of scientific and technical literature depends heavily on the development of modern science and technology (some factors are the acceleration of knowledge increase, interdisciplinarity, intensification of the science-technology interaction). Moreover, the practical organisation of documentation and information services also exerts a great influence on actual usage data.

The relationship between journal productivity and obsolescence was studied by Wallace (1986). It was found that highly productive journals tend to have low journal median citation ages and that high journal median citation ages were always associated with journals that were unproductive in terms of the number of references to those journals in the data set studied by this author. The remaining journals, which were not highly productive and did not have high journal median citation ages appeared to be distributed randomly.

Case studies on obsolescence have been made, among other researchers, by Ravichandra Rao (1973) for the Proceedings of the IRE and the Journal of the American Chemical Society; by Kohut (1974) for geoscience literature; by Clark (1976) for US patent literature; and by Queiroz and Lancaster (1981) for the field of thermoluminescent dosimetry.

### III.7. SCIENCE POLICY APPLICATIONS

#### III.7.1. Generalities

Citation data play an important role in quantitative studies of science and technology. However, extracting citation data from the ISI's database and combining them with data taken from other sources requires great skill and a thorough knowledge of the internal structure of computer files. In this context we refer the reader to Moed (1988), Moed and Vriens (1989), Anderson et al. (1988) and de Bruin and Moed (1989).

We will restrict our presentation of *science policy* applications to a few cases in which we were personally involved, referring readers interested in the science policy aspect of informetrics to start with to Van Raan's Handbook (1988). This book contains 22 contributions written by leading figures in the field, giving a clear overview of all important aspects. The handbook contains not only contributions on methods and techniques to develop indicators for the measurement of research and technological performance, but also papers on techniques to study cognitive processes in the development of scientific fields and in the interaction between science and technology. It includes examples and applications of data-analysis methods going further than the simple methods presented in Part I of this book.

Another important source for science policy applications is the topical issue of *Scientometrics* (May 1989) on the relations between qualitative theory and scientometric methods in science and technology studies (Leydesdorff et al. (1988)). Further, we would like to mention the journals *Scientometrics*, *Science Policy* and *Science and Public Policy*, where the interested reader will find many contributions on science policy issues. Moreover, several talks during the International Conferences on Informetrics (Diepenbeek, London (Ontario), Bangalore (to be held in 1991),...) deal with mappings of science and science policy aspects (Egghe and Rousseau (1988a)).

#### III.7.2. Comparing three weighting methods

In Subsection III.2.4.2 we already mentioned the problem involved in counting citations to multi-authored papers. That different counting procedures can have a dramatic effect on the assessment of the relative scientific strength of different nations will be shown in this section. We will follow the reasoning of an unpublished paper by Eda and Evangelos Kranakis (1988).

III.7.2.1. Three weighting methods

Let  $C$  be a set of countries;  $n$  will denote a number of articles,  $a_i$  the number of co-authors with contributions to the  $i^{\text{th}}$  article and  $a_i(c)$  the number of co-authors from country  $c$  with contributions to the  $i^{\text{th}}$  article.

Clearly, for all  $i = 1, 2, \dots, n$  :

$$a_i = \sum_{c \in C} a_i(c) . \quad [\text{III.7.1}]$$

We are interested in comparing the following three ways of counting the total contribution of a country  $c$  in the given set of  $n$  papers.

1. Straight counting. Only the first author's contribution is acknowledged, with weight 1. In this case the total weight assigned is  $W_1 = n$ . The total weight assigned to country  $c$  depends on the number of times a scientist in country  $c$  is a first author. The probability that country  $c$  will be assigned a weight 1 to paper  $i$  is denoted as  $p_i(c)$ , and hence the total weight assigned to country  $c$  is :

$$W_1(c) = \sum_{i=1}^n p_i(c) . \quad [\text{III.7.2}]$$

Then, the contribution of country  $c$  as a fraction of the overall assignment of weights is :

$$Q_1(c) = \frac{W_1(c)}{W_1} = \frac{1}{n} \sum_{i=1}^n p_i . \quad [\text{III.7.3}]$$

Note that in reality  $Q_1(c)$  may have any value between 0 and the number of papers where at least one of the authors belongs to country  $c$ . Moreover, we have already pointed out that taking the first author is not a random sample.

2. Normal (or unit) counting. Every co-author's contribution in each article is weighted exactly 1 unit. In this case, the total weight assigned is :

$$W_2 = \sum_{i=1}^n a_i . \quad [\text{III.7.4}]$$

The total weight assigned to country  $c$  is

$$W_2(c) = \sum_{i=1}^n a_i(c) , \quad [\text{III.7.5}]$$

and the contribution of country  $c$  as a fraction of the overall assignment of weights is :

$$Q_2(c) = \frac{W_2(c)}{W_2} = \frac{\sum_{i=1}^n a_i(c)}{\sum_{i=1}^n a_i} . \quad [\text{III.7.6}]$$

3. Adjusted (or fractional) counts. Every co-author's contribution in the  $i^{\text{th}}$  article is weighted exactly  $1/a_i$  units. In this situation, the weight of country  $c$  in the  $i^{\text{th}}$  article must be  $a_i(c)/a_i$ . In this case the total weight assigned is :

$$W_3 = \sum_{i=1}^n \sum_{c \in C} \frac{a_i(c)}{a_i} = n .$$

The total weight assigned to country  $c$  is then :

$$W_3(c) = \sum_{i=1}^n \frac{a_i(c)}{a_i} , \quad [\text{III.7.7}]$$

and the contribution of country  $c$  in the overall assignment of weights is :

$$Q_3(c) = \frac{W_3(c)}{W_3} = \frac{\sum_{i=1}^n (a_i(c)/a_i)}{n} . \quad [\text{III.7.8}]$$

#### III.7.2.2. An example

We consider the situation of three cited papers : the first and the second are co-authored by two scientists of countries  $a$  and  $b$ , where the author of country  $a$  is first author twice. The third paper is written by a scientist of country  $c$ .

This yields the following results :

##### 1. First author counting

	a	b	c	Relative contribution
paper 1	1	0	0	Country a : 2/3
2	1	0	0	Country b : 0/3
3	0	0	1	Country c : 1/3

##### 2. Normal counting

	a	b	c	Relative contribution
paper 1	1	1	0	Country a : 2/5
2	1	1	0	Country b : 2/5
3	0	0	1	Country c : 1/5

3. Fractional counting

	a	b	c	Relative contribution
paper 1	1/2	1/2	0	Country a : 1/3
2	1/2	1/2	0	Country b : 1/3
3	0	0	1	Country c : 1/3

Even this elementary example shows the unequal treatment countries get, depending on the weighting method used.

III.7.2.3. Comparisons of weighting methods

As straight counts depend too much on chance, we will concentrate here on the two other methods. Suppose  $\Delta(c) = |Q_2(c) - Q_3(c)|$ ,  $m = \min(a_1, a_2, \dots, a_n)$  and  $M = \max(a_1, a_2, \dots, a_n)$ . Then the following theorem holds.

III.7.2.3.1. Theorem (Kranakis and Kranakis (1988))

For any country c,

$$\Delta(c) = |Q_2(c) - Q_3(c)| \leq \left(\frac{1}{m} - \frac{1}{M}\right) \frac{\sum_{i=1}^n a_i(c)}{n} . \quad \text{[III.7.9]}$$

Proof. Let

$$A_i(c) = \frac{a_i(c)}{a_i} , \quad i = 1, \dots, n.$$

Then :

$$Q_3(c) = \frac{1}{n} (A_1(c) + \dots + A_n(c)) . \quad \text{[III.7.10]}$$

However, for all  $i = 1, \dots, n$  :

$$\frac{a_i(c)}{M} \leq A_i(c) \leq \frac{a_i(c)}{m} . \quad \text{[III.7.11]}$$

Hence,

$$\frac{1}{nM} \sum_{i=1}^n a_i(c) \leq \frac{1}{n} \sum_{i=1}^n A_i(c) \leq \frac{1}{nm} \sum_{i=1}^n a_i(c) . \quad \text{[III.7.12]}$$

But, we also have, by the very definition of m and M :



$$\frac{1}{nM} \sum_{i=1}^n a_i(c) \leq \frac{\sum_{i=1}^n a_i(c)}{n} \leq \frac{1}{nm} \sum_{i=1}^n a_i(c) . \quad [\text{III.7.13}]$$

Hence, by the definition of  $Q_2(c)$  [III.7.6] and  $Q_3(c)$  [III.7.10] we find :

$$\Delta(c) \leq \left( \frac{1}{m} - \frac{1}{M} \right) \frac{\sum_{i=1}^n a_i(c)}{n} .$$

This proves the theorem.  $\square$

#### III.7.2.3.2. Examples

1. If all the  $a_i$ 's are equal, then  $m = M$  and  $Q_2(c) = Q_3(c)$ .

2. If for every  $c$  there exists a constant  $\phi(c)$  such that, for every  $i = 1, \dots, n$   $a_i(c)/a_i = \phi(c)$ , then  $Q_2(c) = Q_3(c)$ . This shows that the upper bound obtained in the theorem is not optimal.

3. Determine  $c$  and assume for every  $i = 1, \dots, n$   $a_i(c) = 1$  and  $a_i = i$ . Then  $W_2 = n(n+1)/2$  and  $W_3 = n$ ;  $W_2(c) = n$  and  $W_3(c) = 1 + 1/2 + 1/3 + \dots + 1/n \approx \log(n) + \gamma$  (where  $\gamma$  is Euler's constant). Hence  $Q_2(c) = 2/(n+1)$  and  $Q_3(c) \approx (\log(n) + \gamma)/n$ . Consequently :

$$\Delta(c) \approx \left| \frac{2}{n+1} - \frac{\log(n) + \gamma}{n} \right| \quad [\text{III.7.14}]$$

and if  $n$  is high :

$$\frac{Q_2(c)}{Q_3(c)} \approx \frac{2}{\log(n)} . \quad [\text{III.7.15}]$$

#### III.7.2.4. A two-country example

In this example we make the following assumptions :

A1  $C = \{c, c'\}$ , i.e. there are two countries.

A2  $1 \leq a_i(c) \leq a_i \leq 2$ ; i.e. country  $c$  contributes to every article and every article has exactly one or two authors.

A3 There exist integers  $k$  and  $\ell$  such that  $a_1 = \dots = a_k = 1$ , i.e. the first  $k$  articles have exactly one author, necessarily belonging to country  $c$ ;  $a_{k+1} = \dots = a_n = 2$ , i.e. the next  $n-k$  papers have exactly two authors. Among these  $n-k$  articles the first  $\ell$  have exactly one author of country  $c$ , and the remaining  $n-k-\ell$  have two authors from country  $c$ .

Now, for the situation described above we get the following results :

$$Q_2(c) = \frac{a_1(c) + \dots + a_n(c)}{a_1 + \dots + a_n} = \frac{k + \ell + 2(n-k-\ell)}{k + 2(n-k)} = 1 - \frac{\ell}{2n-k} ;$$

$$Q_2(c') = \frac{\ell}{n} ;$$

$$Q_3(c) = \frac{1}{n} \left( \frac{a_1(c)}{a_1} + \dots + \frac{a_n(c)}{a_n} \right) = \frac{1}{n} \left( k + \frac{\ell}{2} + \frac{2(n-k-\ell)}{2} \right) = 1 - \frac{\ell}{2n} ;$$

$$Q_3(c') = \frac{\ell}{2n} .$$

If we assume that in those cases where there are two co-authors from different countries, the first belongs to country  $c$  once in three times, then, as follows from the definition of straight counting :

$$Q_1(c) = \frac{k + \frac{\ell}{3} + (n-k-\ell)}{n} = 1 - \frac{2\ell}{3n} ;$$

and

$$Q_1(c') = \frac{2\ell}{3n} .$$

Then  $Q_3(c') \leq Q_1(c') \leq Q_2(c')$ ;  $Q_1(c) \leq Q_3(c)$  and  $Q_2(c) \leq Q_3(c)$ .  
Finally :

$$\Delta(c) = |Q_3(c) - Q_2(c)| = \frac{k\ell}{2n(2n-k)} .$$

Note that for fixed  $n$ , the higher  $\ell$  is (the rate of international co-authorship is increasing), the lower  $Q_1(c)$  will be, for every  $i = 1, 2, 3$ .

Now we come to the main question : are there any conditions under which the percentage of a country's contribution will appear to decline because its rate of international co-authorship is going up, even though its rate of publications also increases?

Yes, this is possible. Suppose that we have two situations  $(\ell^{(1)}, n^{(1)}, k)$  and  $(\ell^{(2)}, n^{(2)}, k)$  (where  $k$  is fixed and smaller than  $n^{(1)}$ ) such that  $\ell^{(1)}/n^{(1)} < \ell^{(2)}/n^{(2)}$ . Then always  $Q_1^{(1)}(c) > Q_1^{(2)}(c)$  and  $Q_3^{(1)}(c) > Q_3^{(2)}(c)$ . If moreover  $\ell^{(1)}/\ell^{(2)} < (2n^{(1)} - k)/(2n^{(2)} - k)$ , also  $Q_2^{(1)}(c) > Q_2^{(2)}(c)$ . For this case, we take, say,  $\ell^{(1)} = n^{(1)}/4$ ,  $\ell^{(2)} = n^{(2)}(1/4 + 1/10)$ . Then the

above condition is reduced to  $k < \frac{4n^{(1)}n^{(2)}}{7n^{(2)} - 5n^{(1)}}$ . Now  $k < n^{(1)} < \frac{4n^{(1)}n^{(2)}}{7n^{(2)} - 5n^{(1)}}$  as soon as  $5n^{(1)} - 3n^{(2)} > 0$  (e.g.  $n^{(1)} = 1000$ ,  $n^{(2)} = 1500$ ).

So, in general (in the above simple example), it is true (in most cases), that the relative contribution of a country declines when its rate of international co-authorship increases. This example can easily be extended in different situations and more countries (Kranakis and Kranakis (1988)).

### III.7.3. Kinematic statistics of scientific output (Rousseau (1989c,d))

#### III.7.3.1. Introduction

Scientific output in all its forms (such as scientific papers, reference lists or patent applications) belong to Popper's World 3 (Popper (1974)), the world of human knowledge as expressed not only in documents, but also in all other human artefacts. According to Brookes (1981), researchers in the field of information science should explore this man-made cognitive world. As a result of such an exploration we will present a new visualisation of the output and the change in output of a system consisting of  $n$  scientific production units or similar entities. (For other representations we refer the reader to Rousseau (1989c).)

The relative share in scientific output of the  $n$  components ( $n \geq 3$ ) will be represented in a regular  $n$ -angle. The  $n$ -angle itself is thought to be inscribed in a circle with a radius of one. So every vertex, each representing one of the  $n$  components, lies at a distance of one from the centre of the circumscribed circle. The orientation of the  $n$ -angle is of no importance as we are only studying relative shares. The order in which the vertices are assigned to the different components, however, does matter.

#### III.7.3.2. The centre of publication/citation

The *centre of publication* (citation), denoted  $\vec{c} = (c_x, c_y)$ , is defined by the following equations :

$$c_x = \frac{\sum_{i=1}^n m_i L_{i,x}}{M}, \quad c_y = \frac{\sum_{i=1}^n m_i L_{i,y}}{M}, \quad \text{[III.7.16]}$$

where  $n$  is the number of elements in the system,  $\vec{L}_i = (L_{i,x}, L_{i,y})$  is the location vector of the  $i^{\text{th}}$  element (i.e. the coordinates of the  $i^{\text{th}}$  vertex of the polygon) and  $m_i$  is the number of publications (or received citations) of the  $i^{\text{th}}$  element (during a given period).

Finally,  $M = \sum_{i=1}^n m_i$  is the total number of publications (or citations)

of the system. As  $\vec{c}$  belongs to the convex hull of the  $\vec{L}_i$ 's,  $\vec{c}$  always belongs to the closed  $n$ -angle.

Since components are represented by the vertices of a regular polygon, the centre of the circumscribed circle is the equilibrium point, i.e. the place where the centre of publication (citation) is situated when all components have an equal share.

The year-to-year displacement,  $\vec{\Delta}$ , of the centre of publication (citation) is an important dynamic parameter of the system under study. Further, we define the standardised momentum vector of the system, denoted as  $\vec{Q}$ , by  $\vec{Q} = \vec{M}(Y) \cdot \vec{\Delta}$  (where  $\vec{M}(Y)$  in the year  $Y$  is the average between the total number of publications (or received citations) in the year  $Y-1$  and the year  $Y$ ). The norm of  $\vec{Q}$  (denoted as  $||\vec{Q}||$ ) is then a measure of the total change in relative importance within the system. The vector  $\vec{Q}$ , or its norm, is, in our opinion, a very important parameter in the system. It combines two key variables: one that characterises the main quantity under study, namely the number of publications (or citations), and one that characterises the behaviour of the system, namely the displacement of the centre of publication (citation).

When there are more than three elements in the system, we chose that representation (i.e. order of vertices) which shows the largest total displacement of the centre. This leads to an interesting combinatorial problem.

With  $n$  components there are  $n!$  possible different assignments of vertices to production units (say, countries). But as the representation is circular, the starting point as well as the orientation does not matter (1-2-3-4-5-6 is equivalent for our purposes to 3-4-5-6-1-2; also 1-2-3-4-5-6 and 1-6-5-4-3-2 are equivalent). This finally leaves  $(n-1)!/2$  cases (compare with Section II.2.5).

### III.7.3.3. An example

As an example we have studied the six element publication system consisting of Norway, Denmark, Finland, Sweden, the Netherlands and Belgium. To find the optimal arrangement we had to consider  $(6-1)!/2 = 60$  cases. The optimal arrangement (i.e. where the total distance travelled by the centre of publication is the largest) was found to be No-De-Ne-Sw-Fi-Be. In this case the total displacement of the centre of publication is almost 60 % more than in the other extreme case in which the total distance is the smallest. This is the case for the arrangement No-Fi-De-Sw-Be-Ne. All data concerning this

system can be found in Tables III.7.1, III.7.2 and III.7.3. Figs.III.7.1 and III.7.2 illustrate the results. We note a general trend in the direction of the two stronger countries : Sweden and the Netherlands.

Table III.7.1

A. Year B. Total number of publications for the six element system (according to the ISI-database) C. Relative share of Norway D. Relative share of Denmark E. Relative share of Finland F. Relative share of Sweden G. Relative share of the Netherlands H. Relative share of Belgium							
A	B	C	D	E	F	G	H
1977	26995	0.0930	0.1464	0.0892	0.2537	0.2495	0.1681
1978	27187	0.0894	0.1418	0.0963	0.2589	0.2493	0.1643
1979	28678	0.0905	0.1394	0.0959	0.2520	0.2599	0.1622
1980	30300	0.0906	0.1374	0.0973	0.2506	0.2559	0.1682
1981	31969	0.0838	0.1411	0.0961	0.2549	0.2621	0.1620
1982	34986	0.0809	0.1339	0.0996	0.2583	0.2625	0.1648
1983	37531	0.0847	0.1322	0.0998	0.2571	0.2742	0.1520
1984	39936	0.0807	0.1284	0.1001	0.2554	0.2813	0.1541
1985	41782	0.0818	0.1277	0.0976	0.2603	0.2852	0.1474
1986	43292	0.0758	0.1342	0.0977	0.2536	0.2913	0.1475
1987	43396	0.0745	0.1271	0.0973	0.2505	0.2944	0.1561

Table III.7.2

A. Year B. x-coordinate of the centre of publication (in the case of the optimal arrangement) C. y-coordinate of the centre of publication (in the case of the optimal arrangement)					
A	B	C	A	B	C
1977	0.1200	-0.1728	1983	0.1339	-0.2173
1978	0.1130	-0.1893	1984	0.1347	-0.2242
1979	0.1223	-0.1886	1985	0.1454	-0.2324
1980	0.1107	-0.1838	1986	0.1561	-0.2315
1981	0.1257	-0.1987	1987	0.1456	-0.2303
1982	0.1143	-0.2091			

Table III.7.3

A	B	C	D	E	F	G
1978	-0.0070	-0.0165	27091	-189.6	-447.0	486
1979	+0.0093	+0.0007	27933	+259.8	+ 19.6	261
1980	-0.0116	+0.0048	29489	-342.1	+141.5	371
1981	+0.0150	-0.0149	31135	+467.0	-463.9	658
1982	-0.0114	-0.0104	33478	-381.6	-348.2	517
1983	+0.0196	-0.0082	36259	+710.7	-297.3	770
1984	+0.0008	-0.0069	38734	+ 31.0	-267.3	269
1985	+0.0107	-0.0082	40859	+437.2	-335.0	551
1986	+0.0107	+0.0009	42537	+455.2	+ 38.3	457
1987	-0.0105	+0.0012	43344	-455.1	+ 52.0	458

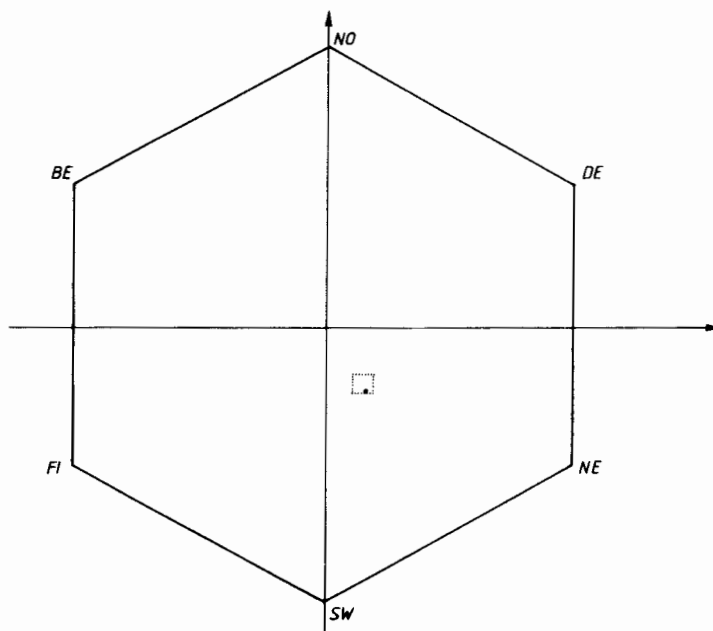


Fig.III.7.1 Optimal graphical representation of the publication system consisting of Norway, Denmark, the Netherlands, Sweden, Finland and Belgium. The point in the small square indicates the publication centre of 1987.

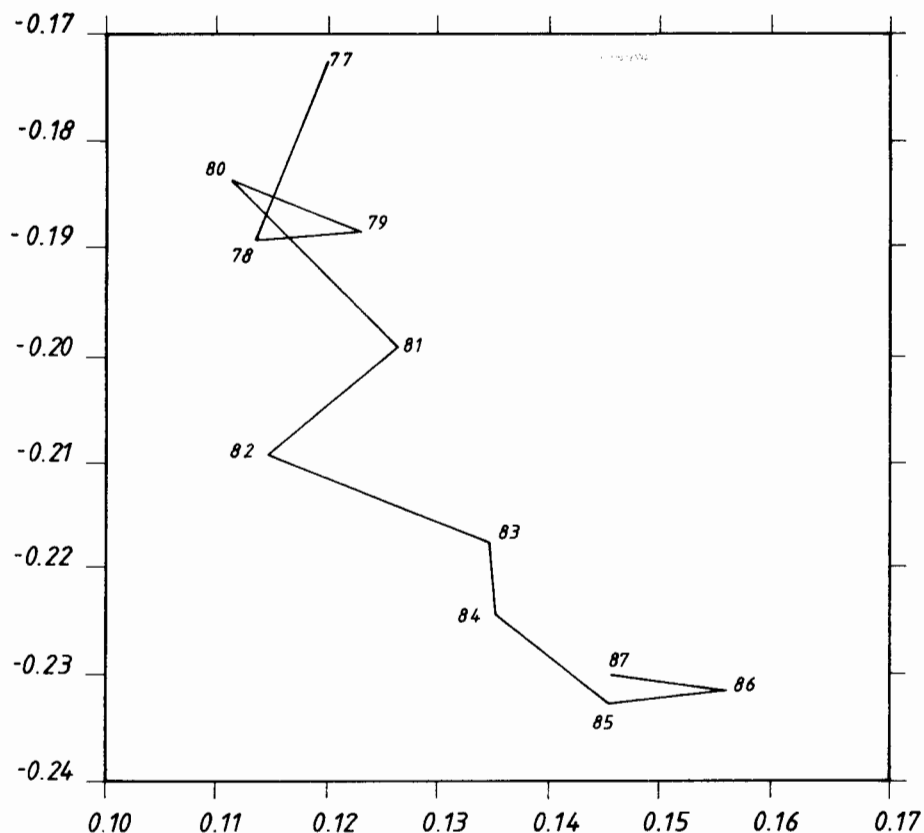


Fig.III.7.2 Year-to-year displacement of the publication centre of the No-De-Ne-Sw-Fi-Be system; close-up of the small square of Fig.III.7.1

#### III.7.4. Notes and comments

The ISI database plays an important role in the development of indicators of international scientific activity. In this context we refer the reader to a paper by Carpenter and Narin (1981), who investigated the adequacy of the SCI as a data source of scientific activity. They found that in general the SCI was a balanced source for the core of physical and biological sciences, but that in particular the coverage of the Soviet literature seemed incomplete.

In all fairness we should also mention that citations are just one of several indicators to evaluate science.

Martin and Irvine wrote a series of papers on the evaluation of Big Science and on recent developments in British science (e.g. Irvine and Martin (1985a,b), Irvine et al. (1985), Irvine et al. (1986)). Their work has put them among the leading European scientists in their field (together with the Leiden group under the direction of Van Raan). They also published a review on the experience of a number of developed countries in identifying and exploiting research that is likely to yield economic and social benefits (Martin and Irvine (1989)).

In the United States science research is reviewed from a quantitative point of view in the biennially published 'Science Indicators'. The bulk of the data presented in Science Indicators is listed in tables, figures and charts. The citation data for this review are furnished by Francis Narin and the staff of Computer Horizons, Inc. Of course, these data are ultimately derived from the ISI database.

Publishing patterns and the citation impact of 32 countries were studied by Braun et al. (1985). This study is regularly updated and extended through a number of 'World Flashes on Basic Research', appearing in *Scientometrics*. Relative indicators such as the activity index and the attractivity index for the comparative assessment of publication output and citation impact were defined and used by Braun, Schubert and Glänzel (Schubert et al. (1986), Braun et al. (1985)).

The activity index (AI) was first proposed by Frame (1977) and is defined as

$$AI = \frac{\text{the country's share in the world's publication output in the given field}}{\text{the country's share in the world's publication output in all science fields}}$$

When  $AI = 1$ , the country's research effort in the given field corresponds exactly to the country's general research effort in the world;  $AI > 1$  reflects a higher than average effort (measured in publications) and  $AI < 1$  a lower than average effort.

The attractivity index (AAI) characterises the relative impact of a country's publications in a given subject field as reflected in the citations they attract :

$$AAI = \frac{\text{the country's share in citations attracted by publications in the field}}{\text{the country's share in citations attracted by publications in all science fields}}$$



AAI = 1 indicates that the country's citation impact in the given field corresponds precisely to the world's average.

A detailed follow-up on scientometric indicators was published by Schubert et al. (1989). Indicators of 96 countries in 114 major fields of science are reported in a special issue of *Scientometrics*.

Citation-based measures of research interactivity were also studied by Pinski (1980). Problems with the evaluation of Big Science are discussed by Moravcsik (1988b).

Citation-based techniques have been shown to be effective in the analysis of research funding (Pao and Goffman (1986)). Grantees of the Edna McConnell Clark Foundation in the field of schistosomiasis research received an exceptionally high average of 4.88 citations per paper per year, showing that Clark grantees have dominated the field.

Moed et al. (1985b) found a serious lack of agreement between past performance analysis by peer judgment and by bibliometric indicators in Holland. It seemed that Dutch National Survey Committees compared Dutch groups only, regardless of their international level.

### III.8. OTHER USES OF CITATION MEASURES

1. Citation analysis is also a tool to be applied in the history of science (Garfield (1979a)). Although major achievements in science can easily be recognised, minor, though sometimes crucial contributions, are often overlooked. Moreover, it is not always easy to retrace the chronology of scientific events. For both problems citation analysis can be helpful.

To illustrate the feasibility of a citation-based approach, Garfield studied the history of the breaking of the genetic code and compared his findings with an account of the facts as written by Asimov (1963). His study showed that citation analysis indeed provides a way of identifying key events and their chronology, of finding relationships between scientific ideas and discoveries, and of judging their relative importance. Moreover, Garfield's study uncovered one event of importance that had been overlooked by Asimov.

Hurt (1985) compared an historical approach (i.e. lists of references used by historians of science) with a bibliometric study to find and rank important authors in the field of quantum mechanics. However, he found a significant difference between the two rankings and concluded that present citation analysis techniques only yield an approximate measure of the importance of literature. More refined models should be investigated.

2. The study of recent citation and referencing data convinced Narin and Noma (1985) that science and technology are far more closely linked today than is generally perceived.

3. To compare the work of scientists who have been publishing in similar fields for different lengths of time, Geller et al. (1978, 1981) devised a model to predict the life-time citation rate per author.

4. Rousseau (1987b) proposed a mathematical method - stemming from the field of operational research - to determine influences on a scientific publication. References in papers cited by the article under study are said to have an indirect influence on this paper.

5. Hayes (1983) analysed citation statistics for more than 400 tenure-level faculty in American schools of library and information sciences. The top 40 faculty were identified and examined in more detail. The most clear-cut result was the evident importance of a Ph.D. program in creating an environment that encourages publication. Of the 23 schools with a Ph.D. program, 16 are

among the group of the top 20 schools.

#### 6. The Ortega hypothesis.

It is often thought that the growth of science owes a lot to the work of average scientists. This is formulated by Jose Ortega y Gasset (1932) as follows :

For it is necessary to insist upon this extraordinary but undeniable fact : experimental science has progressed thanks in great part to the work of men astoundingly mediocre, and even less than mediocre. That is to say, modern science, the root and symbol of our actual civilization, finds a place for the intellectually commonplace man and allows him to work therein with success. In this way the majority of scientists help the general advance of science while shut up in the narrow cell of their laboratory, like the bee in the cell of its hive, or the turnspit of its wheel.

However, based on citation data, Cole and Cole (1972) found that it is primarily elite scientists who contribute to scientific progress. Their findings thus contradicted the Ortega hypothesis. This result of the Coles was corroborated by Oromaner (1985). On the other hand, Turner and Chubin (1976) argued that the Coles' results could also be explained by the fact that science badly misuses the talent at its disposal. The whole problem, set off by a paper of M. and B. MacRoberts (1987), gave rise to a heated debate, resulting in a series of comments published in *Scientometrics* (volume 12 (5-6), 1987).

7. Braam et al. (1988) investigated the influence of citation and co-citation thresholds and the extent to which cluster structures depend on the equation used to compute the co-citation strength. Therefore, they compared the use of Salton's equation (see equation [III.8.2]) with that of the Jaccard index (equation [III.8.1]).

In the context of citation analysis the Jaccard index (Jaccard (1901)) is defined as :

$$S_j(i,j) = \frac{\text{coc}(i,j)}{\text{cit}(i) + \text{cit}(j) - \text{coc}(i,j)} \quad , \quad \text{[III.8.1]}$$

where  $S_j(i,j)$  denotes the co-citation strength between documents  $i$  and  $j$ , as calculated according to Jaccard's equation,  $\text{cit}(k)$  denotes the number of citations received by document  $k$  ( $k = i$  or  $j$ ) and  $\text{coc}(i,j)$  is the number of co-citations received by  $i$  and  $j$ . This number can also be described as the number of items in the intersection of the set of all citations to  $i$  and the set of all citation to  $j$ , divided by the number of items in the union of these

two sets. Jaccard's index is also called the 'relative co-citation frequency'.

Salton's cosine equation (Salton and McGill (1984)) can be defined as :

$$S_s(i,j) = \frac{\text{coc}(i,j)}{(\text{cit}(i).\text{cit}(j))^{1/2}} \cdot \quad \text{[III.8.2]}$$

Both Jaccard's index and Salton's cosine equation are examples of similarity measures. They both take values in the interval [0,1] and, if documents  $i$  and  $j$  are cited at least once, both measures are zero if the documents are not co-cited. Finally, if  $\text{cit}(i) = \text{cit}(j) = \text{coc}(i,j)$ , both measures attain their maximum value of 1.

Although results obtained by using the Jaccard index and Salton's equation were different (Braam et al. (1988)), especially concerning the interrelationships between clusters, there was also a great amount of similarity when the Salton strength value was twice the strength value as calculated by the Jaccard index. As a result, Hamers et al. (1989) conducted an investigation into the relation between  $S_j$  and  $S_s$ . They found that, although  $S_s/S_j$  can take any value between 1 and infinity, in most practical cases it has a value close to 2.

8. The language barrier in the humanities has been studied by Yitzhaki (1988), who introduced the language self-citation index as the proportion of references made by authors in a specific field which are written in the same language as the citing source. The language self-citation index is similar to the notion of journal self-citation and is an indicator of the degree to which researchers on a specific subject field draw upon the literature published in their own language.

9. Some investigations study the use of local journal collections as measured by citations in theses and local study projects. Thus, citation analysis becomes a direct tool in local library collection development (Walther (1972), Chambers and Healey (1973), Laborie and Halperin (1976), Mancall and Drott (1979), Crissinger (1981), McCain and Bobick (1981), Vandegehuchte (1988)).

10. Finally, we mention that citation data can often be described by the informetric laws to be studied in the next part.