

## QUASI-CORRESPONDENCE ANALYSIS ON SCIENTOMETRIC TRANSACTION MATRICES

R. J. W. TIJSEEN,\* J. DE LEEUW,\*\* A. F. J. VAN RAAN\*

*\*Science Studies Unit, LISBON-Institute*

*\*\*Department of Data Theory, University of Leiden, Leiden  
(The Netherlands)*

(Received December 2, 1986)

In principle, a scientometric transaction matrix can be modelled by assuming that the number of transactions is the result of independent row and column contributions. More often one is primarily interested in the cross-structural relations between the participating entities, whereas the row and column marginals are of lesser or no importance. The values of the residuals after fitting an independence model to a complete transaction matrix can be analyzed by correspondence analysis to investigate the structure of the transactions between the rows and columns, after correcting for their marginal frequencies. Recently a modification of correspondence analysis has been developed, quasi-correspondence analysis, which seems quite suitable for the analysis of citation-based transaction matrices which are incomplete or in which the incorporation of certain transactions may seem inappropriate. An illustration of both data analysis-techniques will be given using a journal-to-journal citation matrix.

### Introduction

The rapidly increasing accumulation of scientific knowledge has been an important incentive in the development and application of so-called "science indicators", i.e. research methods to assess, categorize and measure specific characteristics of science, such as the effectiveness of scientific work and research performance. The so-called bibliometric indicators of science use quantified characteristics of scientific literature as, for instance, the number of references within scientific publications. These citations to related publications can be seen as an recognition of such work.

The validity of a citation as such an elementary unit of communication between science entities (e.g., publishing entities like authors, journals or research groups/institutes) remains a matter of controversy: one must not only consider the existence of negative citations, irrelevant citations and self-citations, but also the possible citation-database limitations as well as existing (field- and time-dependent) citation practices (cf. *Moed et al.*<sup>1</sup>). *MacRoberts* and *MacRoberts*<sup>2</sup> have also given evidence that relevant publications are often missing in the reference list and references, in

general, do not accurately cover the topic(s) discussed. Consequently, extreme opinions are to be found on the use of citations in constructing quantitative measures of science, extending from those who would employ citation-based measures as a reliable indication of the peer recognition of scientific work to those who only advocate the use of citations in literature search. Citation-based measures have nevertheless become important science indicators because citation data often yield relatively unobtrusive information on the relations between scientific entities, especially when one is dealing with citations of higher levels of aggregation. In practice such citation-based quantitative measures evaluating the "impact" of scientific activities in a given subfield of science, or between various subfields, have not only proved to be an valuable tool in assessing characteristics or even effectiveness of scientific activities (cf. *Garfield*<sup>3</sup>), but also for providing data in science policy studies (cf. *Moed et al.*<sup>4-5</sup>). In view of these pros and cons mentioned it will assumed in this paper that citations yield an indication of the scientific merit and of the utility of the cited object, and citation-transactions thus provide a useful measure to evaluate interrelationships between scientometric entities.

Aggregated transaction data are often displayed in a matrix, in which sets of entities are assigned to the rows and columns. Each cell of such a transaction-matrix contains citation-values indicating the level of transaction between a row and a column entity, in general the observed number of citations. When one is interested in the interrelationships between the same set of entities, the rows classify the cited-mode and the columns the citing-mode of an entity, or vice versa. The elements in the main diagonal of such a matrix represent the self-transactions. In general, the number of self-transactions is (much) higher than the off-diagonal transaction values. They are often the result of specific features in the transaction process. For example, in the case of journal-to-journal citation-transactions the self-citations within scientific journals are partly due to the fact that authors within a specific field tend to utilize a selection of the scientific journals covering that field. This mechanism can result in a high concentration of citations to publications in the journal in which both the cited author as well as the citing author publish.

There are a number of data analysis-techniques, in which the structure of the transaction matrix can be investigated by modelling the citation-data. An adequate model of the proportionally high number of self-citations will generally lead to a non-adequate modelling of the other elements in the matrix and thus obscure the off-diagonal relationships. Since our interest is mainly on the (off-diagonal) citation relations between entities these diagonal elements are of less or no importance in the modelling-process. In order to discard these self-transactions one can proceed by adjusting or eliminating these values to minimize their effect on the analysis results.

In this paper such a data analysis-technique is discussed, yielding information on the relations between entities of both the row and column-mode, while ignoring (the values in) specified matrix elements. Such elements may be ignored because they are unobserved or missing, because they pertain events that cannot occur, or because they are unusual in other respects. We discuss the case in which the self-citations of a journal-to-journal citation matrix are adjusted, but the formalism to be developed is completely general and can be applied to any element or a group of elements in any rectangular (scientometric) matrix.

### Modelling the transaction data

A large influence of one entity on another entity or a strong interrelationship between scientometric entities can be expressed by a relatively high number of citation-transactions. In general, one can consider the strength of a relationship as a function of the number of these transactions. Scientometric transaction matrices with entities from different domains will generally contain a transaction structure quite different from the expected structure if one only considers the marginals as in the case of the so-called *independence model*. In this model one assumes that the existing transaction levels between the entities are only determined by the "size" of the entities: the rows  $i$  ( $i = 1, \dots, i', \dots, I$ ) and columns  $j$  ( $j = 1, \dots, j', \dots, J$ ) have an independent contribution to the cell-frequencies, i.e. the probability of a transaction from object  $i$  to object  $j$  is equal to the probability of a transaction from column  $i$  multiplied by the probability to receive a transaction in row  $j$ , thus in terms of the population parameters  $\pi$

$$\pi_{ij} = \pi_{i+} \pi_{+j} \quad (1)$$

with the '+' denoting the summation over the omitted index.

It is well known that the maximum likelihood estimator (MLE) of the expected number of transactions, obtained from the datamatrix with the observed number of transactions  $X = \{x_{ij}\}$  is equal to

$$e_{ij} = x_{i+} x_{+j} / n \quad (2)$$

with  $x_{++} = n$

The significance of the difference between the observed value and expected values derived from this model can be evaluated with the Pearson chi-square statistic:

$$\chi^2 = \sum_i \sum_j (x_{ij} - e_{ij})^2 / e_{ij} \quad (3)$$

This statistic has an asymptotic chi-square distribution with  $(I-1)(J-1)$  degrees of freedom (df). If the value of  $\chi^2$  has a probability ( $P$ ) near zero according to the chi-square distribution, the expected and observed transaction values are significantly different. In such cases one might consider an alternative model with additional parameters to account for the variation in the citation frequencies due to first- and/or higher-order interactions between rows and columns of the matrix.

Of course the independence model is often only a baseline-model and it is obvious that the expected values based on this model will generally not fit the citation data. The differences between the original citation data and the expected values will yield useful information on the citation-relations, because the "size-effects" will then have been ruled out. Since we are primarily interested in an analysis of residuals after fitting a suitable restrictive model we focus on the independence model. More sophisticated models (i.e. with more model-parameters) will generally yield a better fit of the data, but often at the cost of problems when interpreting the multitude of parameters and often leaving non-informative residuals. A serious drawback prevents fitting the independence model to a transaction matrix without, for example, involving self-transactions: the ML-estimates of the values in the off-diagonal elements of the matrix cannot be computed directly without the main diagonal. A solution to this problem can be found by introducing the *quasi-independence model*; a generalisation of the independence model to incomplete matrices.

### Quasi-independence

The quasi-independence model enables one to ignore elements in the matrix and still fit an independence model on the remaining observation values  $\{(i, j) \text{ in a given subset of the index pairs } L\}$  based on a MLE-procedure, while estimating the expected values for observations which are not modelled  $\{(i, j) \text{ not in } L\}$ . A brief outline of the quasi-independence model will be given for a two-way matrix. For a more detailed discussion of the concepts the reader is referred to *Goodman*<sup>6</sup>.

Contrary to the independence model, a direct estimation of the expected frequencies is impossible when one fits a model based on quasi-independence. In this case an iterative maximum likelihood algorithm of the following 'iterative proportional fitting' type (cf. *Deming, Stephan*<sup>7</sup>) can be applied: if the matrix with the observed frequencies  $x_{ij}$  is complete, one can fit a quasi-independence model  $\pi$  which assumes that  $\pi_{ij} = \alpha_i \beta_j$  for all  $(i, j)$  in  $L$ . The  $\pi_{ij}$  with  $(i, j)$  not in  $L$  are unrestricted and not estimated; these values are found by substituting the observed values, thus  $\pi_{ij} = x_{ij}$ . The multinomial likelihood equations are

$$\sum_j \{x_{ij} \mid j \in J(i)\} = \sum_j \{\pi_{ij} \mid j \in J(i)\} \quad (4a)$$

$$\sum_i \{x_{ij} \mid i \in I(j)\} = \sum_i \{\pi_{ij} \mid i \in I(j)\} \quad (4b)$$

with  $I(j)$  and  $J(i)$  respectively denoting the index  $i$  and the index  $j$  for which the cells  $(i, j)$  are in  $L$ . It follows from Eqs (4a) and (4b) that the expected marginals of the restricted cells have to be equal to the observed marginals.

The quasi-independence algorithm starts by setting the elements  $(i, j)$  in  $L$  equal to the parameter product  $\alpha_i \beta_j$ , for an arbitrary choice of  $\alpha$  and  $\beta$ . A convenient choice is  $\alpha_i \beta_j = 1$ . For each  $i$  all elements  $\pi_{ij}$  with  $j \in J(i)$  are multiplied with a constant, satisfying Eq. (4a). Only the row-sums add up to the correct marginal numbers. Subsequently the same procedure is applied to the columns  $j$  with  $i \in I(j)$ , resulting in correct column-sums, but undoing the correct row-sums. This process is repeated until all values  $\pi_{ij}$  converge to stable values with an acceptable level of accuracy.

*Price*<sup>8</sup> made a first attempt to fit a quasi-independence model on a square transaction matrix. In *Price*'s procedure the diagonal elements are considered missing and initial values are assigned to the diagonal elements via a multiplicative model based only on the off-diagonal elements. The transaction matrix with the estimated self-transactions is subsequently used to compute the final estimates of both the diagonal and off-diagonal elements based on the independence model. A major shortcoming of this two-step method is the fact that it lacks a proper conceptual basis and each step leads to different expected values of the diagonal elements. *Noma*<sup>9</sup> elaborated on this procedure by fitting a quasi-independence model and subsequently a quasi-symmetric model. The latter model is an extension of the quasi-independence model in which an additional interaction-parameter is incorporated for each row and column-combination. The values of the interaction-parameters were used as input in a multi-dimensional scaling technique to compute a spatial indicator of similarity between the objects. *Noma*'s use of the quasi-independence model and a quasi-symmetric model proved to be a conceptual improvement when modelling transaction data with dominating diagonal values. However, two critical aspects of this procedure to represent the relations between objects based on quasi-independence are to be considered. First, the spatial results clearly depend on the extent to which the data are fitted by the model. Using the parameter values of model with a lack of fit can only result in a rough approximation of the similarities between the entities involved. A more promising approach would be to fit a (highly) restrictive model, such as the quasi-independence model, and analyze the residuals to investigate the remaining structure between the entities. The residuals contain the information on the relations between the entities after correcting for the "size" of the entities and the size of the unproportionally high numbers of self-citations. Secondly, and more important, the rows and columns are treated symmetrically, whereas a transaction matrix is often highly asymmetric,

reflecting large differences between row and column-mode of an entity. For example, the cited and citing characteristics between two journals can be of an entirely different nature, because the journals emphasize on different features of scientific research within a (sub)field, e.g. applied research versus basic research. This asymmetry is not accounted for in the symmetric model and remains hidden in the pattern of the residuals. These points can however be adequately dealt with by the data analysis-techniques: correspondence analysis and, more in particular, quasi-correspondence analysis.

### Correspondence analysis

Quasi-correspondence analysis (abbreviated to QCA) can be used to fit a quasi-independence model to a square transaction matrix and subsequently investigate the relations between the residuals in terms of the relations between scores assigned to the rows and columns. This technique is a generalisation of correspondence analysis (CA in the following), which is basically a standard eigenvector-eigenvalue decomposition of the matrix of residuals after fitting the independence model. In short, CA can thus be seen as a technique which analyses a structure of values after correcting for the marginal frequencies. It can therefore be used complementary to loglinear modelling (cf. *Van der Heijden, De Leeuw*<sup>10</sup>), but it is also possible to interpret CA as a technique able to find a multidimensional representation of the dependence between rows and columns (cf. *Benzécri*<sup>11</sup>). The CA-results can be displayed in a simultaneous spatial representation of scores assigned to the rows and columns of a matrix.

Before describing QCA, a brief discussion of CA must be given. CA can be defined in terms of deviations from the independence model. Let  $X$  be the matrix with the observed number of standardized citation transactions, with entries  $x_{ij}$  adding up to  $n$ . The row marginals  $x_{i+}$  are contained in the diagonal matrix  $D_r$ , and  $D_c$  contains the column marginals  $x_{+j}$ . The vector  $t$  contains elements equal to one. The matrix  $E$  is equal to  $D_r t t' D_c/n$  with elements  $e_{ij}$ . The  $'$ -sing denotes the transpose of a vector or matrix. The matrix containing the standardized residuals after accounting for the row and column effects is decomposed by computing the singular value decomposition

$$D_r^{-1/2}(X - E)D_c^{-1/2} = U\Omega V' \quad (5)$$

where  $U'U = I$  and  $V'V = I$ .  $\Omega$  is a diagonal matrix containing the descending singular values  $\omega_s$ , where  $s$  ( $s = 1, \dots, s$ ) is the index for the orthogonal solutions or, in

geometrical terms, the dimensions. If the independence model fits well, the residuals  $(X - E)$  are small resulting in small singular values. In this case, the CA-results are obviously of not much value; the row and column parameters of the independence model are sufficient to approximate the number of citation transactions.

If one uses a limited number of CA-dimensions to describe the residual-structure each row and column is quantified by  $s$  quantifications using the corresponding elements of the eigenvectors. The resulting row and column scores are normalised by

$$R = D_r^{-1/2} U n^{1/2} \quad (6a)$$

$$C = D_c^{-1/2} V n^{1/2} \quad (6b)$$

with a weighted average equal to 0 and a weighted variance equal to 1. Furthermore,  $R'D_rR = nI$ ,  $C'D_cC = nI$ ,  $t'D_rR = 0$  and  $t'D_cC = 0$ .

The row and column scores are normalised in such a way that the Euclidean distance between a row  $i$  and a row  $i'$  of  $R^* = R\Omega$  is equal to the chi-squared distance  $\delta^2$ , which is defined as the distance between the respective row/column-profiles, were e.g. the profile of a column  $j$  is the column of the values  $x_{ij}/x_{+j}$ :

$$\begin{aligned} \delta^2(i, i') &= (I_i - I_{i'})' D_r^{-1} X' D_c^{-1} X' D_r^{-1} (I_i - I_{i'})' n = (r_i - r_{i'})' \Omega^2 (r_i - r_{i'}) = \\ &= (r_i^* - r_{i'}^*)' (r_i^* - r_{i'}^*) \end{aligned} \quad (7)$$

$I_i$  and  $I_{i'}$  are unit vectors from the identity matrix  $I$ . Approximations of the chi-square distances are found by only considering the columns of  $R^*$  corresponding to the largest singular values of  $\Omega$ . A similar approximation of the columns can be given by the Euclidean distances between the columns of  $C^* = C\Omega$ .

To facilitate the interpretation one can integrate the separate plots of the row and column scores into a single plot. This is done with the aid of a centroid principle expressed in a so-called transition formula, either

$$R\Omega = D_r^{-1} XC \quad (8a)$$

or

$$C\Omega = D_c^{-1} X'R \quad (8b)$$

Depending on the choice of the centroid principle the distances between a row point  $i$  and a column point  $j$  in such a joint plot can be interpreted by regarding the row points as the weighted average—or centroids—of the column points or vice versa. The row and column profiles  $D_r$  and  $D_c$  are used as weights.

If one interprets the relations between row scores (or column scores) using the chi-square distances, one must bear in mind the fact that rows or columns with similar profiles will have small distances between them, whereas large distances indicate considerably different profiles. The profiles of the marginal frequencies of  $X$  are always located in the origin of the plot; points near the origin correspond to profiles resembling the mean profiles. Row and column scores with profiles very deviant from the mean-profile, adding significantly to the chi-square total [cf. Eq. (3)], are found in the periphery of the plot. Using the transition formulas one can roughly interpret the distance between a row  $i$  and a column  $j$ ; the distance between the points is small if  $x_{ij} \gg e_{ij}$ , points are far apart if  $x_{ij} \ll e_{ij}$ .

The relation between the row and column scores and the original data is found through substituting Eq. 6(a, b) in Eq. (5), obtaining the so-called reconstitution formula:

$$D_r^{-1} (X - E) D_c^{-1} n = R\Omega C'$$

leading to

$$X = E + D_r R \Omega C' D_c n^{-1} \quad (9)$$

which clearly shows that correspondence analysis decomposes the departure from independence. The elements of  $R\Omega C'$  are equal to  $(x_{ij} - e_{ij})/e_{ij}$ .

The Pearson-statistic  $\chi^2$  can be defined in terms of the so-called total 'inertia' of a CA-solution:

$$\text{tr}\Omega^2 = \sum_s \omega_s^2 = \chi^2/n \quad (\text{tr} = \text{trace; the sum of the diagonal elements}) \quad (10)$$

The importance of a dimension can now be interpreted as the ratio of the inertia in a dimension and the total inertia  $\omega_s^2/\text{tr}\Omega^2$ , or more simply, as the proportion of  $\chi^2$ , which is decomposed in a dimension  $s$ .

It has already been pointed out that the row and column scores can be represented in different ways as coordinates in a joint plot, each method with its own specific advantages and disadvantages (cf. *Van der Heijden*<sup>1,2</sup>). Using the centroid principle one can choose between a plot of  $(R, C\Omega)$  or  $(R\Omega, C)$ . In both cases distances



between row and column scores approximate chi-squared distances [cf. Eq. (7)]. A symmetric interpretation can be applied by constructing a joint plot of  $(R\Omega^{1/2}, C\Omega^{1/2})$ , which is an approximation of the two centroid-representations. The approximation will be better as the values of the elements of  $\Omega$  become less different. This option will thus spread the distortion of the approximated chi-squared distances equally over the rows and columns, but a clear interpretation of the results in terms of chi-square distances or the centroid principle is lost. The residuals however, can be represented in terms of scalar products between the row and column vectors; the nature and strength of the spatial relationship between coordinates is determined by the length of the vectors from the origin to the points and the angle between them: e.g. a small angle between relatively long vectors indicates a strong relationship between the corresponding objects, whereas orthogonal vectors indicate unrelated rows or columns.

### Quasi-correspondence analysis

Generalizing CA to a technique capable of decomposing residuals after fitting a quasi-independence model is now quite straightforward. It will only be discussed briefly. Detailed information on the maximum likelihood algorithm and other specifics of the technique can be found in *De Leeuw & Van der Heijden*.<sup>13</sup> Starting from the datamatrix  $X$ , a matrix  $Y$  is computed containing the maximum likelihood estimates under the quasi-independence model. The marginals  $D_r$  and  $D_c$  of  $X$  and  $Y$  are identical. The values of the non-modelled elements after fitting the quasi-independence model are similar to those in  $X$ , hence these elements have residuals equal to zero. Analogous to Eq. (5) the singular value decomposition on the matrix of the residuals  $(X - Y)$  is computed:

$$D_r^{-1/2} (X - Y) D_c^{-1/2} = U\Omega V' \tag{11}$$

The centroid principles are now written as

$$R\Omega = D_r^{-1} XC - D_r^{-1} YC \tag{12a}$$

$$C\Omega = D_c^{-1} X'R - D_c^{-1} YR \tag{12b}$$

The relation between the chi-square distances and the singular values [cf. Eqs (3) and (10)] is now lost because the trace of the singular values is equal to

$$\sum_i \sum_j (x_{ij} - y_{ij})^2 / d_i^r d_j^c \tag{13}$$

with  $d_i^r$  and  $d_j^c$  denoting the corresponding elements  $i$  and  $j$  in the diagonal matrices  $D_r$  and  $D_c$ .

A remedy for this unpleasant feature can be found by replacing the elements in the diagonal weighting matrices  $D_r$  and  $D_c$  [Eq. (11)] with maximum likelihood estimators ( $\alpha_i, \beta_j$ ) of the corresponding rows  $i$  and columns  $j$ . In this case the trace of the singular values results in the familiar chi-square statistic for testing quasi-independence:

$$\sum_i \sum_j \{ (x_{ij} - y_{ij})^2 / y_{ij} \mid (i, j) \text{ in } L \} \quad (14)$$

The centroid principles of this normalisation can be found as follows: suppose  $P$  is a matrix with elements  $p_{ij}$  equal to  $x_{ij}$  for all  $(i, j)$  in  $L$ , while  $p_{ij}$  is equal to  $\alpha_i \beta_j$  for all  $(i, j)$  not in  $L$ . Note that this treatment of the elements is the reverse of the previous one: the quasi-independence algorithm now iterates on the diagonal elements. It converges to the same point as the iterative proportional fitting-algorithm in which both algorithms yield the same value for the chi-square test of quasi-independence. Let  $Q = Pt t' P / t' P t$  and  $t$  is a vector with unit-elements. If  $P$  is the matrix of observed values,  $Q$  is the matrix of expected values based on the independence model. In this case  $P - Q = X - E$  and  $D_r$  and  $D_c$  are the marginal frequencies of  $P$  and  $Q$ . The centroid principles are now be defined as

$$R\Omega = D_r^{-1} P C \quad (15a)$$

$$C\Omega = D_c^{-1} P' R \quad (15b)$$

Comparing these centroid principles with Eq. (8) shows that quasi-correspondence analysis of the matrices  $X$  and  $Y$  is identical to CA of the matrix  $P$ .

The characteristics of the different options to plot the row and column scores mentioned in the foregoing, also apply to the QCA-results. For interpretative reasons the symmetric joint plot ( $R\Omega^{1/2}, C\Omega^{1/2}$ ) of the CA and QCA-results is used in the following application to a journal-to-journal transaction matrix.

### Application to astrophysical and astronomical journals

The data consist of citation counts between seven highly cited scientific journals from the United States and Europe, with publications on topics in the fields of astronomy and astrophysics. These two specific scientific subfields were chosen

because journal articles are a predominant form of scientific communication within these subfields. All journals used mainly consist of 'normal' article; letters or other types of short publications often appear in supplements or separate journals. The citation counts were collected by manual search from the 1983 *Journal Citation Reports*.<sup>15</sup> The result is given in Table 1.

An a priori differentiation between the journals is already apparent from the journal titles: two US journals, *Astronomical Journal (AN)* and *Proceedings of the Astronomical Society of the Pacific (PASP)*, and the British journal *Monthly Notices of the Royal Astronomical Society (MN)* emphasize astronomical subjects, while the US journal *Astrophysical Journal (AP)* concentrates on astrophysical topics. The European journal *Astronomy and Astrophysics (AA)* covers both fields. In addition to the publications on astrophysics, the US journal *Astrophysics and Space Science (APSS)* also contains publications on space physics and related topics on the solar system. The *Annual Review of Astronomy and Astrophysics (ARAA)* contains papers in which an overview is given of the past and current developments in various subfields of astronomy and astrophysics; these papers generally contain of large amount of references. Due to summarizing characteristics, the review papers are often highly cited by papers which deal with topics within the reviewed subfield(s).

The elements in the main diagonal of the non-review journals are dominating; 44.9% of the total amount of citations are self-citations. In terms of the citing (column) totals the *AP* has the largest proportion of self-citations, nameiy 68.3%. The *PASP* has the lowest self-citing percentage with 11.0%. Obviously, the review characteristics of *ARAA* prevent a high self-citation rate; the self-citations only amount to 3.1%. The values of the self-citations of the non-review journals all

Table 1  
Journal-to-journal data from the 1983 *Journal Citation Reports*

| Cited  | Citing journal |        |      |      |      |     |     |
|--------|----------------|--------|------|------|------|-----|-----|
|        | 1              | 2      | 3    | 4    | 5    | 6   | 7   |
| 1 AA   | 2714           | 2 009  | 867  | 296  | 297  | 129 | 159 |
| 2 AP   | 4506           | 16 079 | 3383 | 1358 | 1186 | 895 | 975 |
| 3 MN   | 1163           | 2 327  | 1959 | 268  | 315  | 175 | 141 |
| 4 AN   | 454            | 965    | 433  | 651  | 100  | 126 | 58  |
| 5 APSS | 424            | 978    | 261  | 191  | 464  | 94  | 33  |
| 6 PASP | 282            | 576    | 208  | 170  | 69   | 183 | 63  |
| 7 ARAA | 237            | 601    | 140  | 64   | 74   | 32  | 45  |

exceed the expected values on the basis of independent row and column contributions (cf. Table 1 and Table 2).

Before we analyze the matrix with quasi-correspondence analysis it is illustrative to show the effect of those high diagonal values by computing a CA-solution on the complete matrix. A 4-dimensional solution was computed with the CA-program

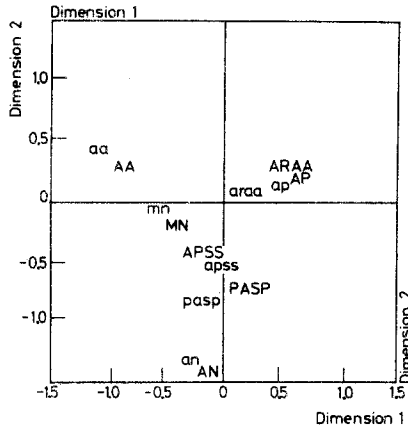


Fig. 1. Correspondence analysis of the journal-to-journal citations from the 1983 *Journal Citation Reports*. Large label - citing journal; small label - cited journal. Labels are centred at the location of the journal-mode

Table 2  
Expected number of self-citations based on the independence model (rounded numbers)

| Journal        | AA   | AP     | MN  | AN  | APSS | PASP | ARAA |
|----------------|------|--------|-----|-----|------|------|------|
| Expected value | 1287 | 13 583 | 936 | 170 | 125  | 52   | 36   |

*ANACOR* (*Gifi*<sup>15</sup>). The chi-square was equal to 7393.6 (df = 36;  $p < 0.001$ ), with singular values  $\omega_1 - \omega_4$  equal to 0.056, 0.039, 0.028 and 0.021, respectively. The magnitude of the first and second singular value compared to the third and fourth singular value suggest a display of the row and column scores of first two dimensions as a parsimonious representation of the results with a relatively small loss of information. The two-dimensional solution decomposes 63% of the inertia; 37% in the first dimension and 26% in the second dimension. The plot of the row and column scores of the CA-solution is given as Fig. 1. As a result of the dominance of the diagonal elements over the off-diagonal elements the row and column profiles

tend to become more similar and, consequently, relatively little differentiation is found between the row and column scores. In the case of the *AP*, *AN* and *MN* the row and column points are located very close to each other as a result of the proportionally high number of self-citations within these journal.

The first dimension of the *CA*-results—the horizontal axis—reveals a relationship which accounts for the largest amount of variance found between the journals: the duality between *AP* and *ARAA* (which have a strong citation-relationship—cf. Table 2) and *AA*. The journals *AA* and *AP* thus have a weaker citation-interrelationship than would be expected, after correcting for the large row and column sums for both journals. *ARAA* cites *AP* much more than expected, whereas the cited *ARAA* is located near the centre indicating a mean cited pattern i.e. the other journals cite *ARAA* in a column-proportional manner. Considering the contents of a review-journal such a result is likely to occur.

If one projects the scores of the other journals on the axes of the first dimension these journals have an intermediate position. The second dimension—the vertical axis—is used mainly to separate *AN* from the other journals.

Considering the contents of the journals, an overall interpretation seems to lead to the conclusion that the position of the three journals in the centre of the triangle (*MN*, *APSS* and *PASP*) is not so much the result of similarities between these journals, but a result of the differentiation between the *US* and European journals (in particular *AA* versus *AP*) in the first dimension, whereas the second dimension tends to differentiate between the (astro)physics-oriented journals (in particular *AN*) and the other journals.

The deviant values of the self-citations of the non-review journals clearly suggest a quasi-correspondence analysis on the data-matrix. The following *QCA*-analysis was computed with the use of a program written in *APL*. The display of the resulting row and column scores is given in Fig. 2.

Of course the remaining chi-square after fitting the quasi-independence model to the data is lower than in the case of the independence model ( $\chi^2 = 354.7$ ;  $df = 36$ ;  $p < 0.001$ ). However, this model still doesn't fit the data adequately. The remaining inertia is decomposed with a two-dimensional *QCA*-solution with singular values equal to 0.062 and 0.049, accounting for respectively 42% and 26% of the inertia. A two-dimensional solution was mainly chosen to simplify the interpretation of the results. The three additional singular values of a 5-dimensional solution (0.041, 0.033 and 0.012, respectively) would have justified a 3-dimensional or even a 4-dimensional solution.

Eliminating the effect of the large diagonal values has resulted in a drastic change in the display of the structure of transaction between the journals: the positions of the citing and cited modes of the journals are now no longer located near each other.

These positions of the rows and columns yield a more accurate reflection of the structure of the journal-interrelations. Although the clear-cut differentiation between the separate journals from the *CA*-solution is lost, the row and column scores of the *QCA*-solution still display an overall structure which is still comparable with the *CA*-results, but after a 90-degree rotation of the axes. The distinction between the

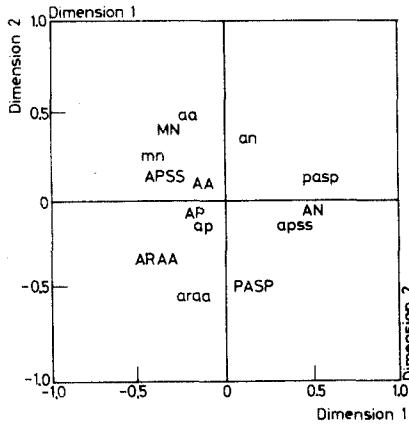


Fig. 2. Quasi-correspondence analysis of the journal-to-journal citations from the 1983 *Journal Citation Reports*. Large label - citing journal; small label - cited journal. Labels are centered at the location of the journal-mode.

Table 3  
The standardized residuals after fitting the quasi-independence model

| Cited journal | Citing journal |        |        |        |        |        |        |
|---------------|----------------|--------|--------|--------|--------|--------|--------|
|               | 1              | 2      | 3      | 4      | 5      | 6      | 7      |
| 1 AA          | 0              | -0.008 | +0.023 | -0.002 | +0.011 | -0.006 | -0.022 |
| 2 AP          | -0.003         | 0      | -0.006 | -0.007 | -0.005 | +0.028 | +0.005 |
| 3 MN          | +0.013         | +0.004 | 0      | -0.018 | +0.008 | -0.019 | -0.011 |
| 4 AN          | -0.004         | -0.006 | +0.020 | 0      | -0.012 | -0.015 | +0.019 |
| 5 APSS        | -0.002         | +0.009 | -0.019 | +0.025 | 0      | -0.026 | +0.006 |
| 6 PASP        | -0.002         | -0.008 | -0.004 | +0.043 | -0.005 | 0      | -0.014 |
| 7 ARAA        | -0.005         | +0.014 | -0.019 | -0.008 | +0.004 | +0.011 | 0      |

astronomical-oriented journals (*AN*, *PASP* and *APSS*) and the other more physics-oriented journals has now become the most important feature within the structure of residual (cf. Table 3), after fitting a quasi-independence model. The rotation is specifically caused by the elimination of the self-citation values of the two highest

(self-)cited journals: *AP* and *AA*. Without the self-citations the profiles of these journals are less dominating in the transaction structure. The more prominent citing relationships of *AN* now become the most important feature in the first dimension.

This leads to a number of differences between the solutions: For example, the first dimension now focusses on the relationship between *AN* and *PASP*, *APSS*. Notice that *PASP* is cited more often by *AN* than expected, whereas in the reverse citation-process this is not the case. In fact, the orthogonality of the cited and citing modes of *AN* and *PASP* in the structure of Fig. 2 indicate that these citation-processes are relatively unrelated to each other. The high citing journals *AP*, *AA* and *MN* have lost their peripheral position—their adjusted profiles now have a larger resemblance with the marginal profiles. *AP* takes a more central position in the plot, indicating that the other journals refer to *AP* in more or less proportional way and in their turn are also proportionally cited by *AP*. Both *AA* and *APSS* also have a central position as citers, spreading their references more or less proportionally over all journals, with *AA* still having a slight emphasis on *MN*.

On the lower side of the figure, the second dimension displays the relationship between the citing *PASP* and the cited *ARAA*; inspection of the matrix of standardized residuals reveals a relatively high positive citation-excess between these modes of *PASP* and *ARAA*. The strong *CA*-relationship between *AP* and *ARAA* is still visible in the *QCA*-solution, but is now considerably weaker. In the upper part of the figure the relatively strong citing relationship between the European journals *MN* and *AA* still exists, especially for the citations from *MN* to *AA*. Clearly both modes of *ARAA* have also lost their core-position, which is surprising considering the characteristics of such a journal. It turns out that *ARAA* takes a position of its own, largely determining the third *QCA*-dimension.

\*

This research is supported by the Netherlands organization for the advancement of pure research (z.w.o.).

### References

1. H. F. MOED, W. J. M. BURGER, J. G. FRANKFORT, A. F. J. VAN RAAN, The application of bibliometric indicators: Important field-and time-dependent factors to be considered, *Scientometrics*, 8 (1985) 177.
2. M. H. MacROBERTS, B. R. MacROBERTS, Quantitative measures of communication in science: A study of the formal level, *Social Studies of Science*, 16 (1986) 151.
3. E. GARFIELD, *Citation Indexing*, New York, Wiley-Interscience, 1979.
4. H. F. MOED, W. J. M. BURGER, J. G. FRANKFORT, A. F. J. VAN RAAN, On the measurement of research performance: the use of bibliometric indicators, Research Policy and Science Studies Unit, University of Leiden, 1983.

5. H. F. MOED, W. J. M. BURGER, J. G. FRANKFORT, A. F. J. VAN RAAN, The use of bibliometric data for the measurement of university research performance, *Research Policy*, 14 (1985) 131.
6. L. A. GOODMAN, The analysis of cross-classified data: independence, quasi-independence, and interactions in contingency tables with or without missing entries, *Journal of American Statistical Association*, 63 (1968) 1091.
7. W. E. DEMING, F. F. STEPHAN, On the least squares adjustment of a sampled frequency table when the expected marginal totals are known, *Annals of Mathematical Statistics*, 11 (1940) 427.
8. D. J. PRICE, The analysis of square matrices of scientometric transactions, *Scientometrics*, 3 (1981) 55.
9. E. NOMA, An improved method for analyzing square scientometric transaction matrices, *Scientometrics*, 4 (1982) 297.
10. P. G. M. VAN DER HEIJDEN, J. DE LEEUW, Correspondence analysis used complementary to loglinear analysis, *Psychometrika*, 50 (1985) 429.
11. J. P. BENZÉCRI et al., *Analyse des donnees* (2 vols), Paris, Dunod, 1973.
12. P. G. M. VAN DER HEIJDEN, Correspondence analysis of transition matrices, *Kwantitatieve Methoden*, 19 (1985) 19.
13. J. DE LEEUW, P. G. M. VAN DER HEIJDEN, *Quasi-Correspondence analysis*, Research Report RR-85-19, Department of Data Theory, University of Leiden, 1985.
14. E. GARFIELD (ed.), *SCI Journal Citation Reports*, 18, Philadelphia, PA; Institute for Scientific Information, 1983.
15. A. GIFI, *ANACOR*, Department of Data Theory, University of Leiden, 1985.