

TRANSITION MATRICES, MODEL FITTING AND CORRESPONDENCE ANALYSIS

Peter G.M. van der Heijden  
Dept. of Research Methods and Psychometrics  
Hooigracht 15  
2312 KH Leiden  
The Netherlands

~~Transparencie~~ Did appear in:

Diday (Ed.), Data Analysis & Informatics IV.  
Amsterdam: North Holland. 1986

Transition matrices are frequently analysed using the independence and quasi-independence model. In the context of transition matrices this last model is used to eliminate the influence of diagonal elements. Usually the independence models do not fit. In this paper we propose to analyze the residuals from these models with (generalized) correspondence analysis. It is also discussed how to deal with structural zeros in the context of correspondence analysis.

1.0 Introduction. Longitudinal categorical data are frequently summarized in two-way transition matrices with elements  $f(i,j)$ , where  $f(i,j)$  is the number of individuals in state  $i$  at time  $t$  ( $i, \dots, t, \dots, T$ ) and in state  $j$  at time  $t+1$ . For each individual the number of transitions equals  $T-1$ . When  $k$  is the number of individuals, the total frequency on the transition matrix equals  $f(+)+k(T-1)$ , where '+' denotes summing over the corresponding index. In this paper we propose to analyse two-way transition matrices with correspondence analysis, and a generalization of correspondence analysis. We will show that this has several advantages for applications in ethology and social mobility respectively. We distinguish these two different fields of application because these fields have different traditions in the analysis of transition matrices, with which correspondence analysis should be connected.

2. Usual ways to analyze transition matrices. A starting point in the analysis of sequences of observations (also called 'sequential analysis') is often to compare the observed transition frequencies with the transition frequencies expected under the assumption that the independence model would hold, i.e. when the state on time  $t+1$  does not depend on the state on time  $t$ . Under the independence model expected frequencies  $m(i,j)$  have the form  $m(i,j) = m(i+)m(+j)/m(+)$ . The difference between the observed frequencies and expected frequencies can be tested with Pearson's chi-square statistic  $\chi^2 = \sum (f(i,j) - m(i,j))^2 / m(i,j)$ . When the independence model does not hold (which is usually the case) one either fits a less restrictive model to all cells, or one tests individual cells for independence. This last approach is often used in ethology. Although the testing of all cells individually is often done in applications, it should not be recommended because the tests are not independent, and one may lose sight of the structure which may exist between significant cells. Often the less restrictive quasi-independence model is fitted. This model states that for some cells in the matrix the expected frequencies are to be equal to the observed frequencies, and for other cells an independence model should hold. This model can be useful when the elements on the diagonal of the matrix are either extremely large, or extremely low. Extreme diagonal elements can be the result from the chosen sampling strategy: for example when time sampling is used, the observed behavior is recorded after each fixed time-interval; when only the transitions to other states are notified, the diagonal is zero by design (structural zeros). In this case the quasi-independence model can be written as

$$(1) \quad m(i,j) = f(i,j) \text{ when } i=j; \quad m(i,j) = a(i)b(j) \text{ when } i \neq j,$$

where  $a(i)$  and  $b(j)$  can be estimated iteratively (cf. [2]). The quasi-independence model can also be used when there are off-diagonal structural zeros. When the quasi-independence model does not fit, less restricted models can be fitted (see section 4.2). Of course, it is possible to test individual cells for departure from model (1), but generally this is not done.

In the sequel transition probabilities, which add row-wise up to one, will be called 'profiles'. Profiles are important since, when the state at time  $t$  is known, the state's profile specifies the probabilities that some other state will follow at time  $t+1$ . Often one is interested in differences between profiles. Correspondence analysis facilitates the study of these differences.

**3.0 Correspondence analysis.** We will treat correspondence analysis briefly here, emphasizing geometrical and quantification aspects. For details and proofs we refer to [1], [5] and [6]. First we discuss classical correspondence analysis, and secondly a generalization of correspondence analysis, proposed in [4].

**3.1 Classical correspondence analysis.** Correspondence analysis is a technique with which it is possible to construct a multi-dimensional representation of the dependence between the row and column variables of a two-way contingency table: scores found for the categories can be used as coordinates for the respective row and column points. These scores can be normalized in such a way that distances between row points or between column points in Euclidean space are equal to chi-square distances. The solution can be found as follows: let  $F$  be the matrix to be analyzed;  $D$  and  $D'$  diagonal matrices with marginal frequencies  $f(i+)$  and  $f(+j)$ ;  $E$  the matrix with expected frequencies computed under the independence model. First the singular value decomposition of the matrix  $D^{-1/2}(F-E)D^{-1/2}$  is computed. Elements of this matrix are equal to standardized residuals, scaled by  $(1/n)^{1/2}$ . These residuals are decomposed with (2):

$$(2) D^{-1/2}(F-E)D^{-1/2} = UV^T,$$

where  $UV^T = I$ ,  $V^T V = I$ , and  $\Lambda$  is a diagonal matrix with singular values  $\lambda(\alpha)$  in descending order;  $\alpha$  is the index for dimension. The dimensionality of the solution is  $\min(j-1, i-1)$ . The scores for rows and columns are normalized as in (3):

$$(3a) R = D_r^{-1/2} U n^{-1/2}; \quad (3b) C = D_c^{-1/2} V n^{-1/2}.$$

The so-called reconstitution formula is found by substituting (3) in (2):

$$(4) F = E + D R A C' D' n^{-1}.$$

Formula (4) shows that the departure from independence is decomposed. This decomposition has the following relation with the  $\chi^2$ -statistic:  $\text{trace } \Lambda^2 = \chi^2/n$ . The importance of dimension  $\alpha$  can be evaluated by the ratio of the inertia of dimension  $\alpha$  and the total inertia:  $\lambda(\alpha)^2 / \sum \lambda(\alpha)^2$ . This quantity can be interpreted as the proportion of  $\chi^2$  that is decomposed in dimension  $\alpha$ .

Clouds of points can be interpreted using the chi-square distances: when two rows (or columns) are near each other, their profiles are similar. Row  $i$  and column  $j$  will be near each other when  $f(ij) > e(ij)$ . An important aid in the interpretation can be found in the property that the sum of the weighted squared distances of the row points (or column points) to the origin, is equal to  $\lambda(\alpha)^2$  for dimension  $\alpha$ . Using this one can evaluate the relative contribution of row  $i$  to dimension  $\alpha$ . Correspondence analysis is formally identical to canonical correlation analysis of contingency tables.  $\lambda(1)$  and the product-moment correlation are related: the correlation between the row and the column variable is, under all possible rescalings of the categories, maximal and equal to  $\lambda(1)$ , when as quantification for the categories the scores for the first dimension are taken.  $\lambda(2)$  is the second maximal correlation of the (orthogonally quantified) variables, etc. Correspondence analysis thus finds the maximal canonical correlations between the quantified row and column variables.

**3.2 A generalization of correspondence analysis.** In this section we will describe a generalization of correspondence analysis proposed in [4] (see also [7]). In [4], Escofier generalizes correspondence analysis by computing the singular value decomposition of the matrix  $S^{-1/2}(G_1 - G_2)S^{-1/2}$  to find scores  $R$  and  $C$ , and singular values  $\Lambda$ . Here  $G_1$  and  $G_2$  matrices of the same size, and  $S_r$  and  $S_c$  are diagonal

matrices with weights for row and column categories.  $S_r$ ,  $S_c$ ,  $G_1$  and  $G_2$  and not necessarily related in the way that  $D_r$ ,  $D_c$  and  $E$  are to the matrix  $F$ . This generalization is difficult to interpret in its most general form. It is advised to use this generalization only in cases that  $G_1$  and  $G_2$  have identical marginal frequencies, and to take these as diagonal elements of  $S_r$  and  $S_c$ . Thus the generalization simplifies to the situation that for  $E$  a matrix different from the independence model is taken. When we denote this matrix as  $G$ , formulas (2), (3) and (4) remain unchanged (apart from replacing  $E$  by  $G$ ). Using the appropriate normalization, the Euclidean distances between the points are still equal to chi-square distances. A row point represents the difference between the profiles of the row in  $F$  and  $G$ . Interpretation of solutions remains basically the same.

**4.0 Correspondence analysis of transition matrices.** The rationale for applying correspondence analysis to transition matrices is the following. First, we saw in section 2 that the profile concept is an important one in the analysis of transition matrices. The difference between two rows of transition probabilities can be studied using correspondence analysis since, with the appropriate normalization, Euclidean distances between the rows are equal to chi-square distances. Secondly, one of the arguments in section 2 for not recommending the study of significant cells was that one might lose sight of the relationship that may exist between these cells. Correspondence analysis can be used to find this structure. A third reason to apply correspondence analysis is to find the maximal canonical correlations between the row and column variable. Finally, in ethology there is a tradition of factor analysis of transition matrices. However, this approach can be criticized for various reasons. Correspondence analysis circumvents many of these criticisms.

We need Escofier's generalization of correspondence analysis to compare observed transition frequencies with expected frequencies following the quasi-independence model (1). Thus for the diagonal cells there is no difference to be reconstituted. We think this is an elegant way to get rid of the usual dominating influence of these elements in classical solutions.

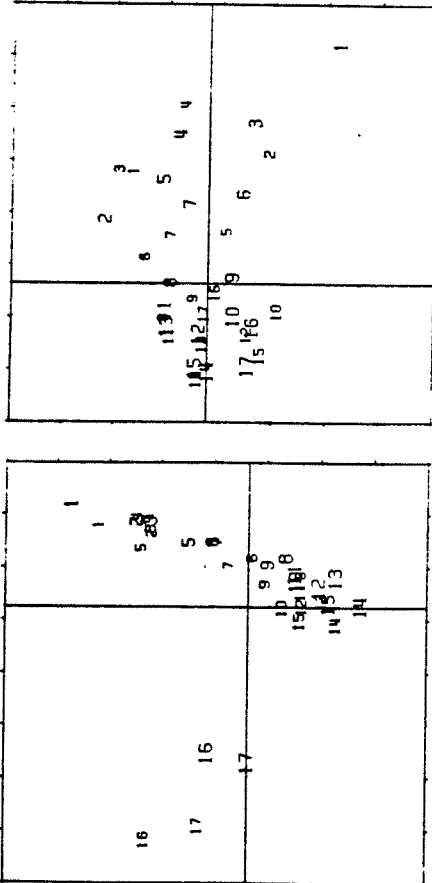
**4.1 An example from ethology.** A special procedure sometimes used in the analysis of transition matrices is factor analysis performed on correlation matrices constructed from a transition matrix. This procedure is thought to be useful when the transition matrix is large, and when it is assumed that the manifest behavior states are triggered by a smaller number of latent motives, drives, or a comparable concept. References can be found in [8]. In this procedure the transition matrix is transformed to a matrix with elements  $f(ij)/m(ij)$ , where  $m(ij) = m(i+)m(+j)/m(++)$ . Subsequently Spearman rank correlations are computed between rows, and columns, yielding two correlation matrices. On these matrices factor analysis is performed. A varimax rotation is used to make interpretation easier. In some later applications this approach is modified by computing standardized residuals  $(f(ij) - m(ij))/m(ij)$ . The following criticisms were raised against this procedure. Firstly, factor analysis of the correlation matrix for the rows and for the columns yields different results, due to the asymmetry of the matrix. A second criticism is that correlations do sometimes not reflect the observed associations in the data: it can occur that a correlation is around zero or negative while two behaviors trigger each other.

We propose to analyze transition matrices with correspondence analysis rather than with factor analysis. It seems that the two criticisms can be circumvented using correspondence analysis: the same number of factors for the rows and for the columns is produced, and the chi-square distance seems a more meaningful measure for the similarity between two rows or two columns. Furthermore, correspondence analysis is a method in which standardized residuals can be decomposed: it stays 'closer' to the original data than factor analysis does.

We will now discuss the analysis of a transition matrix, that we have taken from [9]. This matrix has the special property that there are off-diagonal structural zeros. We will show how to deal with such a situation. The behavior of isolated male zebra finches was recorded for several birds. The transition matrix for bird 27 is shown in table 1. Only transitions to other states were recorded: the diagonal cells are structural zeros. Furthermore, the matrix contains off-diagonal



Figure 3: Correspondence analysis of table 2  
 $\lambda_1 = .51 (.53)$ ,  $\lambda_2 = .31 (.20)$ ; Father is large  
 $\lambda_1 = .25 (.68)$ ,  $\lambda_2 = .08 (.08)$



5.0 Conclusions and discussion. It is shown that (generalized) correspondence analysis is a suitable method for the analysis of transition matrices, especially when the number of categories is large. The reasons for this are outlined in section 4.0, and illustrated in 4.1 and 4.2. Correspondence analysis is probably also a suitable method to study the departure from models different from the quasi-independence model. Furthermore, the study of the departure from the quasi-independence model will probably also be a good way to tackle similar problems with diagonal cells in for instance confusion matrices and import-export tables.

References

- [1] Benzécri, J.P. et al.. Analyse des données (2 vols.) (Dunod: Paris, 1973)
- [2] Bishop, Y.M.M., Fienberg, S.E., & Holland, P.W.. Discrete multivariate analysis: theory and practice (MIT-press: Cambridge, 1975)
- [3] Breiger, R.L.. The social class structure of occupational mobility, Amer. Journ. of Soc., 87, (1981) 578-611.
- [4] Escofier, B. Analyse de la différence entre deux mesures sur le produit de deux mêmes ensembles, Cahiers de l'analyse des données, 3, (1983) 325-329.
- [5] Gifi, A. Non-linear multivariate analysis (R.U.L./F.S.W., Department of Data Theory: Leiden, 1981).
- [6] Greenacre, M.J. Theory and applications of correspondence analysis. (Academic Press: London, 1984).
- [7] Heijden, P.G.M. van der, & Leeuw, J. de. Correspondence analysis used complementary to loglinear analysis. Psychometrika (in press).
- [8] Hooff, J.A.R.A.M. van, Categories and sequences of behavior: methods of description and analysis, in: Scherer, K.R. & Ekman, P. (eds.) Handbook of methods in non-verbal behavior research (Cambridge Univ. Press: Cambridge, 1982).
- [9] Slater, P.J.B., & Ollason, J.C. The temporal patterning of behavior in isolated male zebra finches: transition analysis, Behaviour, 42, (1972) 248-269.