



Diagnostics in a simple correspondence analysis model: An approach based on Cook's distance for log-linear models



Nirian Martín

Department of Statistics and Flores de Lemus Institute, Carlos III University of Madrid, 28903 Getafe (Madrid), Spain

ARTICLE INFO

Article history:

Received 17 December 2013

Available online 20 January 2015

AMS subject classifications:

62J20

62H17

Keywords:

Two-way contingency table

Multinomial sampling

Correspondence model

Log-linear model

Pearson's Chi-square residuals

Cook's distance

ABSTRACT

Diagnostics have not received much attention in the literature of simple correspondence analysis models. Since Cook's distance was defined to identify influential observations of the linear regression model, it has been extended to different models, in particular to log-linear models. In this paper we provide the asymptotic distribution of Cook's distance of any kind of log-linear models and also a method for diagnostics, based on it. By using Goodman's $RC(K)$ model as a log-linear model to approximate the ordinary simple correspondence analysis procedure, we follow a Cook's distance approach to identify influential cells and three examples illustrate the performance of this method.

© 2015 Elsevier Inc. All rights reserved.

1. From the ordinary simple correspondence analysis to Goodman's RC model: a log-linear approach

The original development of the simple correspondence analysis approach, popular from the publication of Benzécri [5], aims to describe association patterns between two categorical variables as well as using graphical procedures for that, once some departures from the model of independence are identified. A variety of extensions have been subsequently considered with respect to correspondence models, introduced by Goodman [11–13] for the first time, to make statistical inference when it is assumed that the data are generated under a specific sampling design. More thoroughly, when reducing the dimensions of the original correspondence analysis model, a new nested model is checked through a goodness of fit test in order to avoid losing representation of row and column profiles of the contingency table. In the following lines we shall describe the basic simple correspondence analysis model, but more details can be found in [3, Section 4]. It is worth mentioning that these models may also be useful for descriptive purposes and graphical displays.

Let X and Y be two categorical variables, with I and J categories respectively and n individuals cross classified in the two way contingency table with multinomial sampling. Without loss of generality, each possible outcome of (X, Y) may be identified by $(i, j) \in \{1, \dots, I\} \times \{1, \dots, J\}$. The probability of a specific outcome, denoted by $p_{ij} = \Pr(X = i, Y = j)$, is unknown and depends on θ through a model ($p_{ij} = p_{ij}(\theta)$) and it is assumed to be non null ($p_{ij} > 0$). The contingency table of unknown probabilities is an $I \times J$ matrix $\mathbf{P} = \mathbf{P}(\theta) = (p_{ij}(\theta))_{i=1, \dots, I; j=1, \dots, J}$ which can be expressed in terms of the transposed row vectors as $\mathbf{P}^T(\theta) = (\mathbf{p}_1(\theta), \dots, \mathbf{p}_I(\theta))$ and stacking these vectors we can express the whole contingency in a unique column vector through the vec operator (see [19, p. 343]) as

$$\mathbf{p} = \mathbf{p}(\theta) = \text{vec}(\mathbf{P}^T(\theta)) = (p_{11}(\theta), \dots, p_{1J}(\theta), \dots, p_{I1}(\theta), \dots, p_{IJ}(\theta))^T.$$

E-mail address: nirian.martin@uc3m.es.

By following the same scheme of the previous contingency table, the observed frequencies may be denoted by a unique column vector, $\mathbf{n} = (n_{11}, \dots, n_{1j}, \dots, n_{I1}, \dots, n_{Ij})^T$, and also its generator, $\mathbf{N} = (N_{11}, \dots, N_{1j}, \dots, N_{I1}, \dots, N_{Ij})^T$, which is a multinomial IJ -dimensional random vector with parameters n and $\mathbf{p}(\boldsymbol{\theta})$, $\mathcal{M}(n, \mathbf{p}(\boldsymbol{\theta}))$.

The “ordinary or basic” simple correspondence analysis (CA) model establishes that

$$p_{ij} = p_{i\bullet} p_{\bullet j} \left(1 + \sum_{k=1}^K \lambda_k u_{ki} v_{kj} \right) \tag{1}$$

$$= p_{i\bullet} p_{\bullet j} \left(1 + \mathbf{u}_i^T \mathbf{D}_\lambda \mathbf{v}_j \right),$$

where $p_{i\bullet} = \Pr(X = i) = \sum_{j=1}^J p_{ij}$, $p_{\bullet j} = \Pr(Y = j) = \sum_{i=1}^I p_{ij}$, $K = \min\{I - 1, J - 1\}$, u_{ik} is the k -th canonical score of category i of variable X , v_{jk} is the k -th canonical score of category j of variable Y , λ_k is the k -th canonical correlation such that $1 \geq \lambda_1 \geq \dots \geq \lambda_K \geq 0$,

$$\boldsymbol{\lambda} = \boldsymbol{\lambda}(K) = (\lambda_1, \dots, \lambda_K)^T, \quad \mathbf{u}_i = \mathbf{u}_i(K) = (u_{i1}, \dots, u_{iK})^T, \quad \mathbf{v}_j = \mathbf{v}_j(K) = (v_{j1}, \dots, v_{jK})^T,$$

and \mathbf{D}_* is the diagonal matrix of $*$. The value of K represents the dimension of the “ordinary or basic” simple correspondence analysis model, the values of the \mathbf{u}_i and \mathbf{v}_j vectors, the i th row and j th column canonical scores respectively, while the values of the matrices $\mathbf{U} = \mathbf{U}(K) = (\mathbf{u}_1, \dots, \mathbf{u}_I)_{K \times I}$ and $\mathbf{V} = \mathbf{V}(K) = (\mathbf{v}_1, \dots, \mathbf{v}_J)_{K \times J}$, all the row and column canonical scores respectively and the values of the $\boldsymbol{\lambda}$ vector, all the canonical correlations. If the summation term $\sum_{k=1}^K \lambda_k u_{ki} v_{kj}$ is null, (1) is equivalent to the independence model, otherwise the value of $\sum_{k=1}^K \lambda_k u_{ki} v_{kj}$ indicates the deviation from independence of cell (i, j) . For identification purposes, since (1) is overparameterized, several constraints have to be considered on parameters. One way to establish these constraints is based on maximizing the correlation between the k -th row of the matrices \mathbf{U} and \mathbf{V} , by following successively the order $k = 1, \dots, K$, that is

$$\sum_{i=1}^I \sum_{j=1}^J u_{ki} v_{kj} p_{ij} = \lambda_k, \quad k = 1, \dots, K \tag{2}$$

is maximized given that \mathbf{U} and \mathbf{V} are centered and standardized,

$$\sum_{i=1}^I u_{ki} p_{i\bullet} = \sum_{j=1}^J v_{kj} p_{\bullet j} = 0, \quad \sum_{i=1}^I u_{ki}^2 p_{i\bullet} = \sum_{j=1}^J v_{kj}^2 p_{\bullet j} = 1, \quad k = 1, \dots, K.$$

Maximizing (2) subject to the previous $4K$ restrictions, \mathbf{U} and \mathbf{V} are obtained, which are uncorrelated in different dimensions, i.e.

$$\sum_{i=1}^I u_{ki} u_{k'i} p_{i\bullet} = \sum_{j=1}^J v_{kj} v_{k'j} p_{\bullet j} = 0, \quad k, k' \in \{1, \dots, K\}, \quad k \neq k'.$$

The previous expressions in matrix notation are given by

$$\mathbf{U} \mathbf{p}_X = \mathbf{V} \mathbf{p}_Y = \mathbf{0}_K, \quad \mathbf{U} \mathbf{D}_{\mathbf{p}_X} \mathbf{U}^T = \mathbf{V} \mathbf{D}_{\mathbf{p}_Y} \mathbf{V}^T = \mathbf{I}_K, \tag{3}$$

where $\mathbf{p}_X = \mathbf{P} \mathbf{1}_J = (p_{1\bullet}, \dots, p_{I\bullet})^T$, $\mathbf{p}_Y = \mathbf{P}^T \mathbf{1}_I = (p_{\bullet 1}, \dots, p_{\bullet J})^T$ are marginal probability vectors ($\mathbf{1}_*$ is the $*$ -th dimensional vector of 1's). Note that $\mathbf{U} \mathbf{D}_{\mathbf{p}_X}^{\frac{1}{2}}$ and $\mathbf{V} \mathbf{D}_{\mathbf{p}_Y}^{\frac{1}{2}}$ are orthogonal matrices. Hence, performing the singular value decomposition of the Pearson chi-square residuals for the independence model

$$\frac{p_{ij} - p_{i\bullet} p_{\bullet j}}{\sqrt{p_{i\bullet} p_{\bullet j}}} = \sqrt{p_{i\bullet} p_{\bullet j}} \sum_{k=1}^K \lambda_k u_{ki} v_{kj}, \quad \mathbf{D}_{\mathbf{p}_X}^{-\frac{1}{2}} (\mathbf{P} - \mathbf{p}_X \mathbf{p}_Y^T) \mathbf{D}_{\mathbf{p}_Y}^{-\frac{1}{2}} = \mathbf{D}_{\mathbf{p}_X}^{\frac{1}{2}} \mathbf{U}^T \mathbf{D}_\lambda \mathbf{V} \mathbf{D}_{\mathbf{p}_Y}^{\frac{1}{2}},$$

we obtain the values of \mathbf{U} , \mathbf{V} and $\boldsymbol{\lambda}$. This idea serves as linking between the correspondence analysis and the correlation analysis. It is worth noting that (1) is a saturated model, that is the adjusted probabilities coincide with the relative frequencies and so in practice the singular value decomposition is

$$\mathbf{D}_{\mathbf{p}_X}^{-\frac{1}{2}} (\widehat{\mathbf{P}} - \widehat{\mathbf{p}}_X \widehat{\mathbf{p}}_Y^T) \mathbf{D}_{\mathbf{p}_Y}^{-\frac{1}{2}} = \mathbf{D}_{\mathbf{p}_X}^{\frac{1}{2}} \mathbf{U}^T \mathbf{D}_{\widehat{\boldsymbol{\lambda}}} \mathbf{V} \mathbf{D}_{\mathbf{p}_Y}^{\frac{1}{2}}. \tag{4}$$

Taking into account that the chi-square statistic for testing independence is defined as n times the sum of the square of the left hand side of (4), i.e., $X^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{1}{p_{i\bullet} p_{\bullet j}} (\widehat{p}_{ij} - \widehat{p}_{i\bullet} \widehat{p}_{\bullet j})^2$, the chi-square statistic for independence might be also expressed in terms of the square of the Euclidean norm of the vector of canonical correlations as

$$X^2 = n \|\widehat{\boldsymbol{\lambda}}\|^2 = n \sum_{k=1}^K \widehat{\lambda}_k^2. \tag{5}$$

Download English Version:

<https://daneshyari.com/en/article/1145614>

Download Persian Version:

<https://daneshyari.com/article/1145614>

[Daneshyari.com](https://daneshyari.com)